



Lecture 6 Network Layer – part 1: Data Plane Chapter 4 (7/e) (6/e Ch4-first part)

EDA344/DIT 423, CTH/GU

Based on the book Computer Networking: A Top Down Approach, Jim Kurose, Keith Ross, Addison-Wesley.

Marina Papatriantafilou – Network layer part 1 (Data Plane)

From last week's review questions

Q: Can an application have reliable data transfer over UDP? If yes how? If no, why?

Student correct A: Yes, if the application layer itself takes care of all the extra functions needed to add reliability on top of UDP's service

Q: How can an ACK have the same number in the same TCP session?

Student correct A: The ACK will have the same number if the package after [the first acknowledged one] it was lost [or arrives out of order]. The ACK is then retransmitted.

Q: Can a TCP's session sending rate increase indefinitely?

Student correct A: ...there are several limits, how fast the receiver can receive, if a sequence is lost then the sending rate will be halved, and also the limits of the hardware

Student Q: Some people argue that the network layer is more correctly called the internet layer. Why do we use the term network layer instead?



The Internet: bridging networks

Internetwork layer (IP main "actor"):

- internetwork appears as single, uniform entity, despite underlying local network heterogeneity
- network of networks

Gateway:

- "embed internetwork packets in local packet format"
- route (at inter-network level) to next gateway



Network layer

Consider transporting a segment from Sender to Receiver

- S: encapsulates segments into datagrams
- R: delivers segments to transport layer
- network layer protocols in *every* host, router
 - examine headers in all datagrams passing through



NW layer's actual job - routing <u>and</u> forwarding

Interplay between the two:



Marina Papatriantafilou – Network layer part 1 (Data Plane)

Roadmap Network Layer

- Forwarding versus routing
- Network layer service models
 - Network layer architecture (shift):
 Software-Defined Networks
 - Inside a router: switching fabrique
 - The Internet Network layer: IP, Addressing & related



- (Next) Control, routing
 - path selection
 - instantiation, implementation in the Internet

Network service model

Q: What service model for "channel" carrying packets from sender to receiver? (general networking scope, ie not only Internet-scope)

example services for **individual** packets:

- guaranteed delivery
- guaranteed delivery with less than 40 msec

example services for a **flow** of packets:

- in-order delivery
- guaranteed min/avg/max arrival rate to flow
- datagram network provides network-layer connectionless service
 - classic Internet model
- virtual-circuit network can provide network-layer connection-oriented service
 - not present in classic Internet protocols but efforts to simulate behaviour are being made

Virtual circuits:

"source-to-dest path behaves almost like telephone circuit"

- call setup, teardown for each call *before* data can flow
 - signaling protocols to setup, maintain, teardown VC (ATM, frame-relay, X.25; not in IP)
- each packet carries VC identifier (not destination host)
- *every* router maintains "state" for each passing connection
- resources (bandwidth, buffers) may be *allocated* to VC (dedicated resources = predictable service)



Virtual Circuits forwarding table



Marina Papatriantafilou – Network layer part 1 (Data Plane)

Datagram networks (the Internet model)

- no call setup at network layer
- routers: no state about end-to-end connections
 - no network-level concept of "connection"
- packets forwarded using only destination host address



Datagram (IP) forwarding table



Datagram or VC network: why?

"Classic" Internet (datagram)

- data exchange among computers
 - "elastic" service, no strict timing req.
- many link types
 - different characteristics
 - uniform service difficult
- "smart" end systems (computers)
 - can adapt, perform control, error recovery
 - simple network core, complexity at edge

VC (eg ATM: a past's vision of the future's ww-network)

- evolved from telephony
- human conversation:
 - strict timing, reliability requirements
 - need for guaranteed service
- "dumb" end systems
 - "Classic" telephones
 - complexity in the core of network

Re-shaping in progress Software-Defined Networks

Roadmap Network Layer

- Forwarding versus routing
- Network layer service models
 - Network layer architecture (shift): Software-Defined Networks
- How a router works: switching fabrique
- The Internet Network layer: IP, Addressing & related



- (Next) Control, routing
 - path selection
 - instantiation, implementation in the Internet

Per-router ("classic" Internet) control plane

Individual routing algorithm (control) components *in each and every router* interact (in the control plane)



Marina Papatriantafilou – Network layer part 1 (Data Plane)

Logically separated control plane

A distinct (can be centralized/remote/distributed) controller interacts with local control agents (CAs)

- this architecture (Software Defined Networking) can enable new functionality
- (will be discussed more later in the course)



Marina Papatriantafilou – Network layer part 1 (Data Plane)

Roadmap Network Layer

- Forwarding versus routing
- Network layer service models
 - Network layer architecture (shift):
 Software-Defined Networks
- Inside a router
- The Internet Network layer: IP, Addressing & related



- (Next) Control, routing
 - path selection
 - instantiation, implementation in the Internet

Router architecture overview



Input port functions



Switching fabrics

- transfer packet from input buffer to appropriate output buffer
- switching rate: rate at which packets can be transfer from inputs to outputs
 - often measured as multiple of input/output line rate
 - N inputs: switching rate N times line rate desirable
- three types of switching fabrics:



Switching via memory

first generation routers:

- traditional computers with switching under direct control of CPU
- packet copied to system's memory
- speed limited by memory bandwidth (2 bus crossings per datagram)



Switching via a bus

datagram from input port memory

to output port memory via a shared bus

- bus contention: switching speed limited by bus bandwidth
- 32 Gbps bus, Cisco 5600: sufficient speed for access and enterprise routers



bus

Switching Via an Interconnection Network

- Overcome bus bandwidth limitations
- Banyan networks, other interconnection nets (also used in multi-processorsmemory interconnects)
 - Cisco 12000: switches at 60 Gbps
 - Example Banyan interconnect: using 3-bit link address













- *buffering* required when datagrams arrive from fabric faster than the transmission rate
- scheduling discipline chooses among queued datagrams for transmission

Datagram (packets) can be lost due to congestion, lack of buffers

Priority scheduling – who gets best performance (vs network neutrality)

Example contemporary routers

Cisco Catalyst 3750E

Stackable (can combine units)64 Gbps bandwidth13 Mpps (packets per second)12,000 address entries

Price: from 100 kSEK





HP ProCurve 6600-24G-4XG Switch 10 Gbps

Up to 75 Mpps (64-byte packets) Latency: < 2.4 μs (FIFO 64-byte packets) 10,000 address entries

Price approx. 50 kSEK

Roadmap Network Layer

- Forwarding versus routing
- Network layer service models
 - Network layer architecture (shift):
 Software-Defined Networks
- How a router works
- The Internet Network layer: IP, Addressing & related



- (Next) Control, routing
 - path selection
 - instantiation, implementation in the Internet

The Internet network layer



IPv4 datagram format



Roadmap Network Layer

- Forwarding versus routing
- Network layer service models
 - Network layer architecture (shift): Software-Defined Networks
- How a router works
- The Internet Network layer: IP, Addressing & related
 - Hierarchical addressing



- (Next) Control, routing
 - path selection
 - instantiation, implementation in the Internet

IP addressing: introduction

- *IP address:* 32-bit id for host/router *interface*
- interface: connection between host/router and physical link
 - routers have multiple interfaces
 - end-host typically has 1-2 interfaces (e.g., wired Ethernet and wireless 802.11)



223.1.1.1 =	11011111	00000001	00000001	00000001
	223	1	1	1

Subnets

• IP address:

- -subnet part: high order bits (variable number)
- -host part: low order bits
- what 's a subnet ?
 - -Devices that can physically reach each other without intervening router
 - device interfaces have same subnet-part (prefix) of IP address



network consisting of 3 subnets

Subnets

recipe

- to determine the subnets: detach each interface from its host or router, i.e. create islands of isolated networks
- each isolated network is a *subnet*



223.1.3.0/24

subnet mask: eg /24

defines how to find the subnet part of the address ...

Marina Papatriantafilou – Network layer part 1 (Data Plane)

CIDR: Classless InterDomain Routing

- subnet portion of address of arbitrary length
- address format: a.b.c.d/x, where x is # bits in subnet portion of address



Subnets, masks, calculations



	Binary form	Dot-decimal notation
IP address	11000000.10101000.00000101.10000010	192.168.5.130
Subnet mask	11111111.11111111.1111111100000000 24 first bits set to 1	255.255.255.0
Network prefix: <i>bitwise AND of</i> (address, mask)	11000000.10101000.00000101 00000000	192.168.5.0
Host part (obtained with similar calculation, with a mask having the 8 last bits = 1)	0000000.0000000.00000000.10000010	0.0.0.130

CIDR Address Masks

CIDR Notation	Dotted Decimal	CIDR Notation	Dotted Decimal
/1	128.0.0.0	/17	255.255.128.0
/2	192.0.0.0	/18	255.255.192.0
/3	224.0.0.0	/19	255.255.224.0
/4	240.0.0.0	/20	255.255.240.0
/5	248.0.0.0	/21	255.255.248.0
/6	252.0.0.0	/22	255.255.252.0
/7	254.0.0.0	/23	255.255.254.0
/8	255.0.0.0	/24	255.255.255.0
/9	255.128.0.0	/25	255.255.255.128
/10	255.192.0.0	/26	255.255.255.192
/11	255.224.0.0	/27	255.255.255.224
/12	255.240.0.0	/28	255.255.255.240
/13	255.248.0.0	/29	255.255.255.248
/14	255.252.0.0	/30	255.255.255.252
/15	255.254.0.0	/31	255.255.255.254
/16	255.255.0.0	/32	255.255.255.255

Marina Papatriantafilou Network layer part 1 (Data Plane)

Classless Address: example

An ISP has an address block 122.211.0.0/16

A customer needs max. 6 host addresses,

□ ISP can e.g. allocate: 122.211.**176.208**/29

□ 3 bits enough for host part

subnet mask 255.255.255.248

Reserved addresses		Dotted Decimal	Last 8 bits	RFC 3021 "The network address itself
	Network	122.211.176. 208	11010000	{ <network-number>, 0} is an obsolete form of directed broadcast,</network-number>
	1st address	122.211.176.209	11010001	but it may still be used by older hosts."
		•••••	•••••	
	6th address	122.211.176. 214	11010110	
	Broadcast	122.211.176. 215	11010111	

Marina Papatriantafilou – Network layer part 1 (Data Plane)

Roadmap Network Layer

- Forwarding versus routing
- Network layer service models
 - Network layer architecture (shift): Software-Defined Networks
- How a router works
- The Internet Network layer: IP, Addressing & related
 - Hierarchical addressing
 - How to get addresses
- (Next) Control, routing
 - path selection
 - instantiation, implementation in the Internet



IP addresses: how to get one (for an end-host)?

hard-coded by system admin in a file

- Windows: control-panel->network->configuration->tcp/ip->properties;
- UNIX: /etc/rc.config

DHCP: Dynamic Host Configuration Protocol: dynamically get address:

- 1. host broadcasts "DHCP discover" msg
- 2. DHCP server (in same subnet) responds with "DHCP offer" msg
- 3. host requests IP address: "DHCP request" msg
- 4. DHCP server sends address: "DHCP ack" msg

DHCP returns more than just allocated IP address:

- address of first-hop router for client
- name and IP address of DNS sever
- network mask (indicating network versus host portion of address)



IP addresses: how to get one (net-part)?

Q: how does *network* get subnet part of IP addr?

A: gets allocated portion of its provider ISP's address space; eg:



Hierarchical Addressing: Route Aggregation

- □ Hierarchical addressing allows efficient advertisement of routing information
- □ The "outside" does not need to know about subnets.



Forwarding: longest prefix matching

r longest prefix matching

when looking for forwarding table entry for given destination address, use *longest* address prefix that matches destination address

Destination Address Range	Link interface
11001000 00010111 00010*** ********	0
11001000 00010111 00011000 ********	1
11001000 00010111 00011*** ********	2
otherwise	3

examples:

DA: 11001000 00010111 00010110 10100001 DA: 11001000 00010111 00011000 10101010

which interface? which interface?

IP Addressing: the last word...

<u>Q</u>: How does an ISP get block of addresses?

<u>A:</u> ICANN: http://www.icann.org/

Internet Corporation for Assigned Names and Numbers

- allocates addresses
- manages DNS
- assigns domain names, resolves disputes

ISPs obtain IP addresses from a

- Local Internet Registry (LIR) or
- National Internet Registry (NIR),
- their appropriate Regional Internet Registry (RIR, 5 worldwide).





Roadmap Network Layer

- Forwarding versus routing
- Network layer service models
 - Network layer architecture (shift): Software-Defined Networks
- How a router works
- The Internet Network layer: IP, Addressing & related
 - Hierarchical addressing
 - How to get addresses
 - NAT
- (Next) Control, routing
 - path selection
 - instantiation, implementation in the Internet



(Well, it was not really the last word...) NAT: network address translation



NAT: network address translation



NAT: network address translation

- 16-bit port-number field:
 - 64k simultaneous connections with a single LAN-side address!
- NAT is controversial:
 - routers should in principle process up to layer 3
 - violates end-to-end argument
 - NAT possibility must be taken into account by app designers, e.g., P2P applications
 - address shortage should instead be solved by IPv6

Roadmap Network Layer

- Forwarding versus routing
- Network layer service models
 - Network layer architecture (shift): Software-Defi Networks
- How a router works
- The Internet Network layer: IP, Addressing & related
 - Hierarchical addressing
 - How to get addresses
 - NAT
 - IPv6
- (Next) Control, routing
 - path selection
 - instantiation, implementation in the Internet



IPv6: motivation

- *initial motivation:* 32-bit address space almost completely allocated.
- additional motivation: header format must help to:
 - speed processing/forwarding
 - distinguishing types of traffic (for potentially different services)

IPv6 datagram format:

- fixed-length 40 byte header
- no fragmentation allowed
- 128-bit addresses (2¹²⁸ = 10³⁸ hosts)
- Standard subnet size: 2⁶⁴ hosts

IPv6 datagram format

priority: identify priority among datagrams in flow flow Label: identify datagrams in same "flow." (concept of "flow" not well defined).

checksum: removed entirely to reduce processing time at each hop *options:* allowed, but outside of header, indicated by "Next Header" field

ver	pri		flow labe		
þ	bayloac	l len	next hdr	hop limit	
source address (128 bits)					
destination address (128 bits)					
data					
← 32 bits					

Marina Papatriantafilou – Network layer part 1 (Data Plane)

Network Layer

Transition from IPv4 to IPv6

- not all routers can be upgraded simultaneously
 - how can the network operate with mixed IPv4 and IPv6 routers?
- tunneling: IPv6 datagram carried as payload in IPv4 datagram among IPv4 routers



Tunneling (6in4 – static tunnel)



Marina Papatriantafilou – Network layer part 1 (Data Plane)

- Google: 8% of clients access services via IPv6
- NIST: 1/3 of all US government domains are IPv6 capable
- Long (long!) time for deployment, use
 - -20 years and counting!

-think of application-level changes in last 20 years: www, streaming media, social media, skype, ...

-Why?

Roadmap Network Layer

- Forwarding versus routing
- Network layer service models
 - Network layer architecture (shift): Software-Defined Networks
- How a router works
- The Internet Network layer: IP, Addressing & related
 - Hierarchical addressing
 - How to get addresses
 - NAT
 - IPv6
- (Next) Control, routing
 - path selection
 - instantiation, implementation in the Internet



Reading instructions Network Layer (incl. next lecture)

 KuroseRoss book 	
Careful	Quick
6/e: 4.1-4.6 7/e: 4.1-4.3, 5.2-5.4, 5.5, 5.6, [new- SDN, data and control plane 4.4, 5.5: in subsequent lectures, connecting to multimedia/streaming Study material available through the pingpong-system]	6/e: 4.7 7/e: 5.7

•	network	layer	servi	ce	mod	el	S
		•					

 Contrast virtual circuit and datagram routing (simplicity, cost, purposes, what service types they may enable)

forwarding versus routing

- Explain the interplay between routing and forwarding
- how a router works
 - What is inside a router? How/where do queueing delays happen inside a router? Where/why can packets be dropped at a router?
- What is subnet? What is subnet masking?
 - Practice masking calculations, dividing address spaces
- Explain how to get an IP packet from source to destination
- Explain how NAT works.

Some complementary material /video-links

- IP addresses and subnets <u>http://www.youtube.com/watch?v=ZTJIkjgyuZE&list=PLE9F3F05C381ED8E8&featu</u> <u>re=plcp</u>
- How does PGP choose its routes <u>http://www.youtube.com/watch?v=RGe0qt9Wz4U&feature=plcp</u>

Some taste of layer 2: no worries if not all details fall in place, need the lectures also to grasp them.

- Hubs, switches, routers <u>http://www.youtube.com/watch?v=reXS_e3fTAk&feature=related</u>
- •
- What is a broadcast + MAC address <u>http://www.youtube.com/watch?v=BmZNcjLtmwo&feature=plcp</u>
- Broadcast domains: <u>http://www.youtube.com/watch?v=EhJO1TCQX5I&feature=plcp</u>

Extra slides

Network	Network Service		Guarantees ?			
Architecture	Model	Bandwidth	Loss	Order	Timing	feedback
Internet	best effort	none	no	no	no	no (inferred via loss)
ATM	CBR	constant rate	yes	yes	yes	no congestion
ATM	VBR	guaranteed rate	yes	yes	yes	no congestion
ATM	ABR	guaranteed minimum	no	yes	no	yes
ATM	UBR	none	no	yes	no	no

- Internet models being extented
- (will study these later on)

Input port queuing

- fabric slower than input ports combined -> queueing may occur at input queues
 - queueing delay and loss due to input buffer overflow!
- Head-of-the-Line (HOL) blocking: queued datagram at front of queue prevents others in queue from moving forward



output port contention: only one red datagram can be transferred. lower red packet is blocked one packet time later: green packet experiences HOL blocking

switch

fabric/

Marina Papa

Network Layer

Output port queueing



- buffering when arrival rate via switch exceeds output line speed
- queueing (delay) and loss due to output port buffer overflow!

DHCP client-server scenario



Marina Papatriantafilou – Network layer part 1 (Data Plane)

DHCP client-server scenario



Marina Papatriantafilou – Network layer part 1 (Data Plane)

NAT traversal problem

- client wants to connect to server with address 10.0.0.1
 - server address 10.0.0.1 local to LAN (client client can't use it as destination addr)
 - only one externally visible address: 138.76.29.7
- solution1: statically configure NAT to forward incoming connection requests at given port to server
 - e.g., (123.76.29.7, port 2500) always
 forwarded to 10.0.0.1 port 25000
- *Solution 2:* automate the above through a protocol (universal plug-and-play)
- Solution 3: through a proxy/relay (will discuss in connection to p2p applications)

NAT traversal problem

- solution 2: Universal Plug and Play (UPnP) Internet Gateway Device (IGD) Protocol. Allows NATed host to:
 - learn public IP address (138.76.29.7)
 - add/remove port mappings (with lease times)
 - i.e., automate static NAT port map configuration

NAT traversal problem

- solution 3 (application): relaying (used in Skype)
 - NATed server establishes connection to relay
 - External client connects to relay
 - relay bridges packets between two connections

IP fragmentation, reassembly

- network links have MTU (max.transfer size) - largest possible link-level frame
 - different link types,
 different MTUs
- large IP datagram divided ("fragmented") within net
 - one datagram becomes several datagrams
 - "reassembled" only at final destination
 - IP header bits to identify + order related fragments

IP fragmentation, reassembly

	length	ו ID	fragf	lag offs	et	
example:	=4000) =x	=0) =()	
 4000 byte datagram MTU = 1500 bytes 	one larg several s	e data smalle	ngram er data	become agrams	S	
1480 bytes in		length	ID	fragflag	offset	
data field		=1500) =x	=1	=0	
	/				-	
offset =		length	ID	fragflag	offset	
1480/8		=1500) =X	=1	=185	
		length	ID	fragflag	offset	
	:	=1040	=x	=0	=370	

- as it travels source to destination
- addr fields of interest here

misc fields	223.1.1.1	223.1.1.3	data
----------------	-----------	-----------	------

Starting at A, given IP datagram addressed to B:

- Iook up net. address of B
- find B is on same net. as A (B and A are directly connected)
- link layer will send datagram directly to B (inside link-layer frame)

misc			data
fields	223.1.1.1	223.1.2.3	uala

Starting at A, dest. E:

- Iook up network address of E
- **E** on *different* network
- routing table: next hop router to E is 223.1.1.4
- link layer is asked to send datagram to router 223.1.1.4 (inside link-layer frame)
- datagram arrives at 223.1.1.4
- continued.....

-			
misc			
fields	223.1.1.1	223.1.2.3	data

- Arriving at 223.1.4, destined for 223.1.2.2
- Iook up network address of E
- E on *same* network as router's interface 223.1.2.9
 - o router, E directly attached
- link layer sends datagram to 223.1.2.2 (inside link-layer frame) via interface 223.1.2.9
- datagram arrives at 223.1.2.2!!! (hooray!)

