

# Network Layer - Data plane -

Kyunghan Lee Networked Computing Lab (NXC Lab) Department of Electrical and Computer Engineering Seoul National University https://nxc.snu.ac.kr kyunghanlee@snu.ac.kr







# Study Goals

- Understanding principles behind network layer services, focusing on data plane:
  - network layer service models
  - forwarding versus routing
  - how a router works
  - generalized forwarding

□ Understanding principles behind network control plane:

- traditional routing algorithms
- SDN controllers
- Internet Control Message Protocol
- network management





## Network Layer

- Transport segment from sending to receiving host
- on sending side encapsulates segments into datagrams
- on receiving side, delivers
  segments to transport layer
- network layer protocols in every host, router
- router examines header fields in all IP datagrams passing through it







## Two key network-layer functions

#### Network-layer functions:

- forwarding: move packets from router's input to appropriate router output
- routing: determine route taken by packets from source to destination
  - routing algorithms

#### Analogy: taking a trip

- forwarding: process of getting through single interchange
- routing: determine path from source to destination by a navigation algorithm





### Network layer: data plane, control plane

#### Data plane

- local, per-router function
- determines how datagram arriving on router input port is forwarded to router output port
- forwarding function



#### Control plane

- network-wide logic
- determines how datagram is routed among routers along end-to-end path from source host to destination host
- two control-plane approaches:
  - traditional routing algorithms: implemented in routers
  - software-defined networking (SDN): implemented in (remote) servers





### Per-router control plane

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







## Logically centralized control plane

A distinct (typically remote) controller interacts with local control agents (CAs)





Introduction to Data Communication Networks, M2608.001200, 2021 FALL SEOUL NATIONAL UNIVERSITY



### Network service model

**Q**: What *service model* for "channel" transporting datagrams from sender to receiver?

example services for individual datagrams:

- guaranteed delivery
- guaranteed delivery with less than 40 msec delay

example services for a flow of datagrams:

- □ in-order datagram delivery
- guaranteed minimum bandwidth to flow
- restrictions on changes in inter-packet spacing











### Router architecture overview

□ High-level view of generic router architecture:















11







## Destination-based forwarding

forwarding table					
Destination Address Range	Link Interface				
11001000 00010111 00010000 00000000 through 11001000 00010111 00010111 11111111	0				
11001000 00010111 00011000 00000000 through 11001000 00010111 00011000 11111111	1				
11001000 00010111 00011001 00000000 through 11001000 00010111 00011111 11111111	2				
otherwise	3				

What if destination addresses are not well split?





### Longest prefix matching

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 which interface? DA: 11001000 00010111 00011000 10101010 which interface?





### Longest prefix matching

- We will see why longest prefix matching is used shortly, when we study addressing
- Longest prefix matching: often performed using ternary (i.e., 0, 1, X) content addressable memories (TCAMs)
  - content addressable: present address to TCAM: retrieve address in one clock cycle, regardless of table size
  - Cisco Catalyst: can up ~1M routing table entries in TCAM



source: www.thenetworksherpa.com



Introduction to Data Communication Networks, M2608.001200, 2021 FALL SEOUL NATIONAL UNIVERSITY



## Switching fabrics

- Transfer packet from input buffer to appropriate output buffer
- Switching rate: rate at which packets can be transferred 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: high speed for access and enterprise routers



bus





## Switching via interconnection network

- overcome bus bandwidth limitations
- <u>banyan networks</u>, <u>crossbar</u>, other interconnection nets initially developed to connect processors in multiprocessor
- advanced design: fragmenting datagram into fixed length cells, switch cells through the fabric.
- Cisco 12000: switches 60 Gbps through the interconnection network (1999 model)
- Cisco NEXUS 9500 (datacenter switch): switches up to 115 Tbps (2020 model)
  - Supports 10G/40G/50G/100G/400G ports









### 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 ports



 buffering required when datagrams arrive from fabric faster than the transmission rate

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

 scheduling discipline chooses among queued datagrams for transmission

Priority scheduling – who gets best performance, network neutrality





### Output port queueing



 buffering when arrival rate via switch exceeds output line speed (or output service rate)

queueing (delay) and loss due to output port buffer overflow!





## How much buffering?

□ RFC 3439 rule of thumb

- average buffering equal to "typical" RTT (say 250 msec) times link capacity C
- e.g., C = 10 Gpbs link: 2.5 Gbit buffer

### Recent recommendation

- with N flows, buffering equal to  $\frac{RTT \cdot C}{\sqrt{N}}$ 
  - Think about deep buffer vs. shallow buffer





## Scheduling mechanisms

- scheduling: choose next packet to send on link
- FIFO (first in first out) scheduling: send in order of arrival to queue
  - real-world example?
  - discard policy: if packet arrives to full queue: who to discard?
    - *tail drop:* drop arriving packet
    - *priority:* drop/remove on priority basis
    - *random:* drop/remove randomly







# Scheduling policies

- Priority scheduling: send highest priority queued packet
- multiple *classes*, with different priorities
  - class may depend on marking or other header info, e.g. IP source/dest, port numbers, etc.
  - real world example?







# Scheduling policies

#### Round Robin (RR) scheduling:

- multiple classes
- cyclically scan class queues, sending one complete packet from each class (if available)
- real world example?







# Scheduling policies

#### Weighted Fair Queuing (WFQ):

- generalized Round Robin
- □ each class gets weighted amount of service in each cycle
- □ real-world example?







## Reading Assignment #3 – Chapter 4&5

Written Quiz #3: Nov. 25<sup>th</sup> (4~5 questions)



Chapter 4	The	Netwo	rk Layer: Data Plane	333
	4.1 Overview of Network Layer			334
		4.1.1	Forwarding and Routing: The Network Data and Control Planes	334
		4.1.2	Network Service Models	339
	4.2	What's	Inside a Router?	341
		4.2.1	Input Port Processing and Destination-Based Forwarding	344
		4.2.2	Switching	347
		4.2.3	Output Port Processing	349
		4.2.4	Where Does Queuing Occur?	349
		4.2.5	Packet Scheduling	353
	4.3 The Internet Protocol (IP): IPv4, Addressing, IPv6, and More			357
		4.3.1	IPv4 Datagram Format	358
		4.3.2	IPv4 Datagram Fragmentation	360
		4.3.3	IPv4 Addressing	362
		4.3.4	Network Address Translation (NAT)	373
		4.3.5	IPv6	376
	4.4	Genera	lized Forwarding and SDN	382
		4.4.1	Match	384
		4.4.2	Action	380
		4.4.3	OpenFlow Examples of Match-plus-action in Action	386
	4.5	Summa	ary	389
	Homework Problems and Questions			
	Wireshark Lab			
	Interview: Vinton G. Cerf			
Chapter 5	Th	e Netw	ork Layer: Control Plane	40
	5.1	Introd	luction	40
	5.2	Routi	ng Algorithms	40
		5.2.1	The Link-State (LS) Routing Algorithm	40
		5.2.2	The Distance-Vector (DV) Routing Algorithm	41
	5.3	Intra-	AS Routing in the Internet: OSPF	41
	5.4	Routi	ng Among the ISPs: BGP	42
		5.4.1	The Role of BGP	42
		5.4.2	Advertising BGP Route Information	42
		5.4.3	Determining the Best Routes	42
		5.4.4	IP-Anycast	43
		5.4.5	Routing Policy	43
		5.4.6	Putting the Pieces Together: Obtaining Internet Presence	43
	5.5	The S	DN Control Plane	43
		5.5.1	The SDN Control Plane: SDN Controller and SDN Control	
			Applications	43
		5.5.2	OpenFlow Protocol	44
		5.5.3	Data and Control Plane Interaction: An Example	44
		5.5.4	SDN: Past and Future	44

5.6 ICMP: The Internet Control Message Protocol
 5.7 Network Management and SNMP
 5.7.1 The Network Management Framework
 5.7.2 The Simple Network Management Protocol (SNMP)
 5.8 Summary



447

449

450

452

454