

L4: Internet Routing and Controls



Sebastian Magierowski
York University

Outline

- Routing
- Error Detection
- Retransmission
- Congestion Control
- Flow Control
- Medium Access Control

Forwarding Tables

- How do routers know what to do with their packets?
- Their forwarding tables tell them:
- Forwarding: The process of taking a packet from an input and sending it out the appropriate output
- Forwarding tables need to contain every detail of a link

Destination	Interface	MAC Address
128.208.128.0/17	if0	8a:0c:1f:e4:6b:1c
128.208.0.0/18	if0	8a:0c:bb:e4:3b:a1
128.208.96.0/19	if2	8a:0c:7b:a9:b2:fc

- They are often implemented in VLSI hardware
 - high-speed memories

- know the IP of where you want to send
- the router port from which you want to send
the physical address of the destination (if applicable, not all protocols require this)

CSE 3213, W14

L4: Routing & Control

3

Routing Tables

- Where do forwarding tables come from?
- From routing tables:
- Routing: The process of building the tables that determine the correct destinations for packets

Destination	Next Hop
128.208.128.0/17	171.69.245.10
128.208.0.0/18	171.69.245.10
128.208.96.0/19	178.45.23.124

- Simpler than forwarding tables
 - typically just a data structure in a computer

but logically contain the information we want to need to move packets through network to destinations

CSE 3213, W14

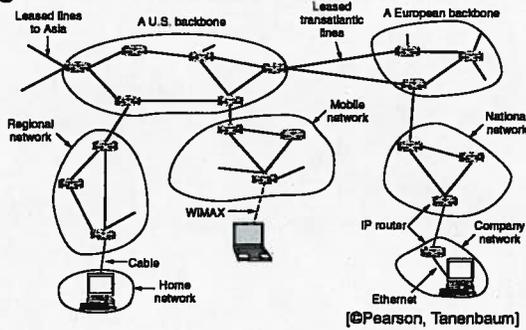
L4: Routing & Control

4

Routing

how do you organize routing

- How do you build routing tables? How do you route?
- Routers run algorithms that update their knowledge of the network every few seconds to hours
- They do this by sending out queries for information and by responding to queries
- Routing seeks to find the cheapest path from any source to any destination
- Minimize link costs



CSE 3213, W14

L4: Routing & Control

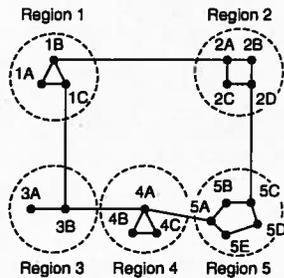
5

first a reminder
↳ Hierarchical Routing

- Routing over two-levels substantially cuts on complexity

- Within AS
- And between them
 - Region-region comms condensed to single router
 - Increased path length a common penalty (e.g. 1A to 5C)

sacrifice some efficiency



CSE 3213, W14

L4: Routing & Control

6

non-hierarchical

1A: Full Table

Dest.	Line	Hops
1A	-	-
1B	1B	1
1C	1C	1
2A	1B	2
2B	1B	3
2C	1B	3
2D	1B	4
3A	1C	3
3B	1C	2
4A	1C	3
4B	1C	4
4C	1C	4
5A	1C	4
5B	1C	5
5C	1B	5
5D	1C	6
5E	1C	5

hierarchical

1A: Hierarchical Table

Dest.	Line	Hops
1A	-	-
1B	1B	1
1C	1C	1
2	1B	2
3	1C	2
4	1C	3
5	1C	4

hops to individual

hops to group

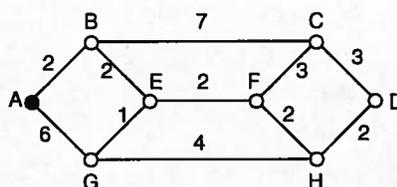
much smaller

©Pearson, Tanenbaum

Link Metrics (Costs)

- Routing often seeks to identify the shortest path between destinations in a graph, the smallest link metric
- Easiest is to just treat all links the same and just count hops (RIP: Routing Information Protocol, does this)
- But many options exist

- mean delay (latency)
- distance
- bandwidth
- average traffic
- communication cost
- political/economic policy



Bellman-Ford, like BGP - distance-vector approach

[©Pearson, Tanenbaum]

CSE 3213, W14

L4: Routing & Control

7

Routing Types

- Methods of building routing tables?
- They do this by sending out queries for information and by responding to queries

link state routing

distance vector routing

- Intradomain Routing
 - Routing within an AS
 - OSPF (Campus), IS-IS (ISPs)
- Interdomain Routing
 - Routing between AS
 - BGP

OPEN shortest path first Intermediate System

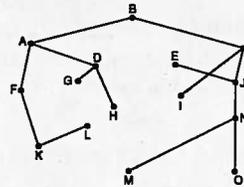
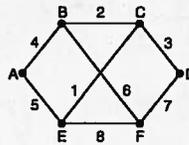
CSE 3213, W14

L4: Routing & Control

8

Intradomain Routing

- **OSPF: Open Shortest Path First**
 - Link-state routing
 - Every router builds up a knowledge of the whole network topology
 - 1) Find link quality to all next-hop neighbors (local view formed)
 - 2) Send this information throughout the whole network (global view formed)
 - 3) Compute shortest path to every router (Dijkstra's algorithm)
- Details...



[©Pearson, Tanenbaum]

CSE 3213, W14

L4: Routing & Control

9

Intradomain Routing

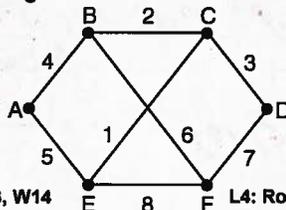
- **Learning about neighbours**
 - When router is introduced it sends out HELLO packet on each line
 - Router on other line sends back its name
- **Setting link costs**
 - Connecting routers can construct costs by sending their bandwidth limits
 - Delay can also be constructed by sending ECHO packets
- **Building link state packets**
 - Aggregate the info

getting to know who you are connected to

point-to-point

find numbers

build packet containing all data learned



A		B		C		D		E		F	
Seq.	Age										
B	4	A	4	B	2	C	3	A	5	B	6
C	2	C	2	D	3	F	7	C	1	D	7
F	6	F	6	E	1			F	8	E	8

CSE 3213, W14

L4: Routing & Control

10

Intradomain Routing

Distributing link state packets

- * - Trickiest part *
- All routers must get packets quickly and reliably
- If routers have different versions of topology odd behaviour will result
- Use flooding to distribute link state packets *← a type of routing (must route to route)*
- Every incoming packet is sent out on every outgoing line
 - Except the one it arrived on
 - How do you keep from swamping the network? *← control spread?*
- Packets have sequence numbers (Seq)
 - Each time you (the source) send out a new packet increment its sequence number, k
 - Routers keep track of (source router, $k_{largest}$) pairs *← if get packet with smaller k discard it*
 - If incoming packet has $k < k_{largest}$ discard it
 - 32-bit k 's would take 137 to loop at 1 packet per second
 - Age field is decremented in case router goes down

keep propagating the packets

A
Seq.
Age

don't worry about running out of k

in order
for each new link state packet u originate (not forward)
recall Age is given by link state packet itself

router re-boots with $k=0$? Age field used,

CSE 3213, W14

L4: Routing & Control

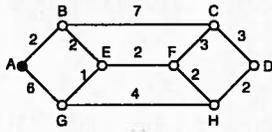
11

(source router, $k_{largest}$, Age) set to 1 min. (or whatever latest arrival was) & decrement every sec. if no update in that time, delete records (i.e. will accept $k=0$)

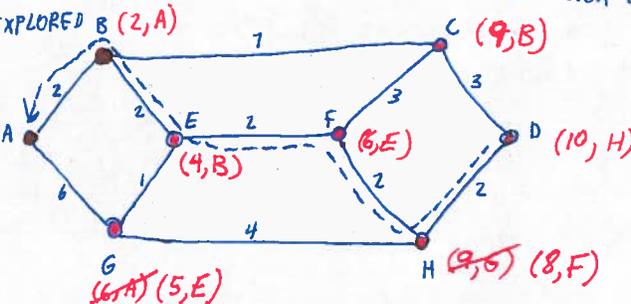
Age also decremented by each router (i.e. of link state packet itself) prevents packets from lingering in network

Shortest-Path Finding

Now apply Dijkstra's algorithm to find shortest path



next steps
from start point: explore all possible paths & choose best link - record them
from best link: explore all possible next steps - record them
until reach destination - from record find best possible link & remove from record



routing table entry

$D \rightarrow H \rightarrow F \rightarrow E \rightarrow B \rightarrow A$
Destination Next hop
D B

- (x,y) in FOUND
- (x,y) in EXPLORED
- ① put all possible hops in FOUND
- ② put cheapest from FOUND in EXPLORED
- ③ put all possible hops from cheapest in FOUND
- ④ if hop is already in FOUND update it if new hop is cheaper
- ⑤ back to ② repeat until destination put in EXPLORED

CSE 3213, W14

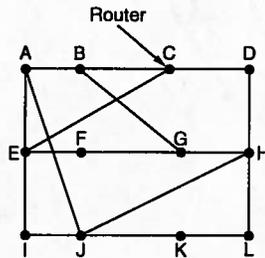
L4: Routing & Control

12

Interdomain Routing

- **BGP: Border Gateway Protocol**
 - Path vector routing (a form of distance vector routing)
 - Exchanging data with neighbours to incrementally form a global view of the network

table "vector" holding best known distance to each destination
i.e. each router has this a global view



CSE 3213, W14

L4: Routing & Control

[©Pearson, Tanenbaum]

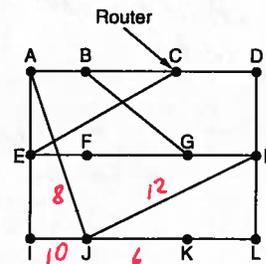
13

Vector Tables

- Each router maintains a vector table (e.g. J)
 - For each router (destination) in the network keeps track of...
 - The next hop it should take
 - The total (estimated) distance to the destination

Destination	Next Hop	Total Distance
A	A	8
B	A	15
C	H	12
D	A	17
E	I	22

- It can start a basic network table by talking to its neighbours
 - JA=8, JI=10, JH=12, JK=6



CSE 3213, W14

L4: Routing & Control

14

100ms or so?

Vector Table Updates

- Every period T each router shares its vector table with each of its neighbours, X, and their distance to destination Y

- Looks at their distance, estimates XY
 - Neighbour X's estimate to destination Y
 - If my estimate to Y, $R_Y > X_Y + R_X$...
 - ... update my row for Y to
 - next hop: X
 - distance: $X_Y + R_X$

To	A	I	H	K
A	0	24	20	21
B	12	36	31	28
C	25	18	19	36
D	40	27	8	24
E	14	7	30	22
F	23	20	19	40
G	18	31	6	31
H	17	20	0	19
I	21	0	14	22
J	9	11	7	10
K	24	22	22	0
L	29	33	9	9

J receives these vectors from neighbours

cost · next hop

8	A
20	A
28	I
20	H
17	I

min. (RX + XY)

old routing table not used

- For example what happens to J's vector table with...
 - $JA=8, JI=10, JH=12, JK=6$
- This technique is called the distributed Bellman-Ford algorithm

J to... 8 10 12 16

CSE 3213, W14

L4: Routing & Control

15

Distance Vector Convergence

- Good news travels fast
 - Delay metric is number of hops
 - A is down initially
 - But when it comes up...
 - ... each exchange propagates the news in a linear fashion
- Bad news travels slow
 - A suddenly goes down
 - B does not hear from A...
 - ...but C thinks it is 2 hops away
 - B thinks it can get to A from C
 - But B & D think they are 3 away
 - So C updates to 4, etc., etc.
 - Distance = 1 + min(neighbour)
 - Slow count to infinity

metrics for node distance to A

	A	B	C	D	E
Initially	0	∞	∞	∞	∞
After 1 exchange	0	1	∞	∞	∞
After 2 exchanges	0	1	2	∞	∞
After 3 exchanges	0	1	2	3	∞
After 4 exchanges	0	1	2	3	4

most optimistic distance

	A	B	C	D	E
Initially	0	∞	∞	∞	∞
After 1 exchange	0	3	2	3	4
After 2 exchanges	0	3	4	3	4
After 3 exchanges	0	5	4	5	4
After 4 exchanges	0	5	6	5	6
After 5 exchanges	0	7	6	7	6
After 6 exchanges	0	7	8	7	8

@1 @B @C
A C B B D
A 2 3 1 3

@2 @B @C
A C B B D
A 2 3 3 3

good news propagate naturally through network

but when link totally goes down network is fooled by its local view instead of announce a big loss it propagates updates one step at a time

choosing most optimistic distance & incrementing by 1 hop

hears B gets from A + updates → C gets B's info about A & updates, etc.

CSE 3213, W14

L4: Routing & Control

16

skew towards optimism → natural thing to do, but slow in responding to catastrophic failure

Routing Summary

- **Intradomain**
 - OSPF
 - link-state
 - • global communication
 - local computation
 - Runs in network layer
 - acknowledged IP
 - Tends to have high memory requirements
 - Keeping track of each router's link state
 - More computation in implementing graph search
- **Interdomain**
 - BGP
 - distance vector
 - local communication
 - global computation
 - Runs in application layer
 - utilizes TCP
 - Slow at pruning out bad links
 - Bad news travels slowly (count-to-infinity problem)

share neighbour data with everyone else locally compute global network routes

exchanges info only with neighbours
or at least distributed. every router participates in a part of a global computation

CSE 3213, W14

L4: Routing & Control

17

Error Detection

- Physical links, router and host hardware can corrupt messages
- Routers should have the ability to detect errors
- Many approaches are used at many levels
 - Line coding at the physical layer (will discuss layers in L4)
 - Turbo codes
 - Reed-Solomon codes
 - LDPC codes
- At the data link and network layers
 - Use header checksums
 - A parity scheme

CSE 3213, W14

L4: Routing & Control

18

Retransmission of Erroneous Information

- What do you do when you detect an error?
- You can just forget about the erroneous packet
 - System doesn't breakdown, but information is lost
- Or you can request that the same packet be retransmitted
 - Foundation of a reliable delivery service
 - Such a scheme is called: Automatic Repeat reQuest (ARQ)
 - Units send messages and expect acknowledgments
 - A number of strategies are employed to make this approach reliable and efficient
 - Present in both the data link and transport layers
 - In the Internet this is typically carried out by the hosts not the routers

CSE 3213, W14

L4: Routing & Control

19

Congestion Control

- What happens if many hosts send packets through one link?
- Router buffer is overwhelmed and must start discarding packets
- Hosts don't get acknowledgments and thus slow down the rate at which they send information

(so you are typically controlling this with the endpoints)

CSE 3213, W14

L4: Routing & Control

20

Flow Control

- A fast transmitter can overwhelm a slow receiver
- In this case receiver indicates (in acknowledgments) to the transmitter how much buffer space it has remaining
- Transmitter doesn't send unless it is aware of enough buffer space in receiver
 - Implemented within ARQ

↑ automatic repeat request
→ here there is explicit comms between hosts about state of exchange
→ in congestion control we react to circumstances (i.e. no ack)

} at least this is a basic thing

CSE 3213, W14

L4: Routing & Control

21

Congestion & Flow Control Summary

- | | |
|---|--|
| <ul style="list-style-type: none">• Congestion Control<ul style="list-style-type: none">– Prevents overflowing the router buffers– Concerned with network internals | <ul style="list-style-type: none">• Flow Control<ul style="list-style-type: none">– Prevents overflowing the destination buffer– Concerned with end-to-end operation |
|---|--|

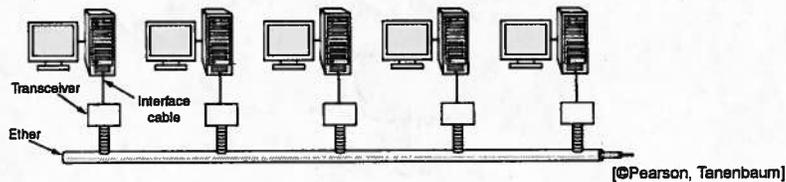
CSE 3213, W14

L4: Routing & Control

22

Medium Access Control (MAC)

- Multiple hosts try to communicate over one medium
 - one wire
 - one radio channel
- How do the units organize their behaviour in order to achieve useful communication?
 - This is the job of the MAC
- This is more the job of LAN and less a specific Internet function



CSE 3213, W14

L4: Routing & Control

23