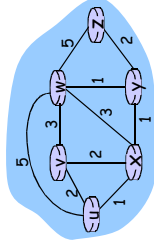


Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
 - RIP
 - OSPF
 - BGP
- 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 Routing algorithms
 - Link state
 - Distance Vector
- 4.6 Routing in the Internet
- 4.7 Broadcast and multicast routing

Network Layer 4-1

Graph abstraction



Graph: $G = (N, E)$

N = set of routers (nodes) = $\{u, v, w, x, y, z\}$

E = set of links (edges) = $\{(u,v), (u,x), (v,w), (x,w), (x,y), (w,z), (y,z)\}$

Remark: Graph abstraction is useful in other network contexts

Example: P2P, where N is set of peers and E is set of TCP connections

Network Layer 4-3

Routing Algorithm classification

Global or local (decentralized) information?

- Global:
 - all routers have complete topology, link cost info
 - "link state" algorithms
- Local/Decentralized:
 - router knows physically-connected neighbors, link costs to neighbors
 - iterative process of computation, exchange of info with neighbors
 - "distance vector" algorithms

Static or dynamic?

- Static:
 - routes change slowly over time
- Dynamic:
 - routes change more quickly
 - periodic update
 - in response to link cost changes

Network Layer 4-5

A Link-State (Global) Routing Algorithm

Dijkstra's algorithm

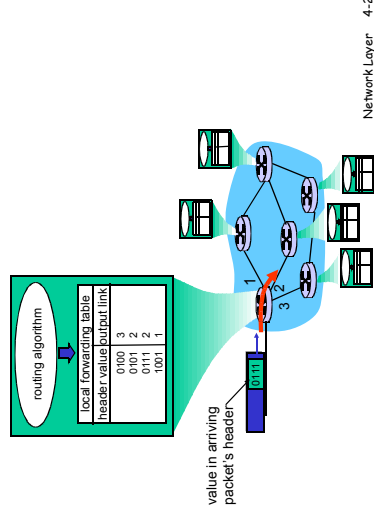
- net topology, link costs known to all nodes
 - accomplished via "link state broadcast"
 - all nodes have same info
- computes least cost paths from one node ("source") to all other nodes
 - gives forwarding table for that node
- iterative: after k iterations, know least cost path to k dest.'s

Notation:

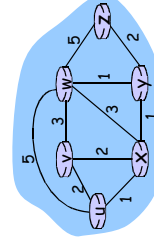
- $c(x,y)$: link cost from node x to y ; $= \infty$ if not direct neighbors
- $D(v)$: current value of cost of path from source to dest. v
- $p(v)$: predecessor node along path from source to v
- N' : set of nodes whose least cost path definitively known

Network Layer 4-7

Interplay between routing, forwarding



Graph abstraction: costs



$c(x,x') = \text{cost of link } (x,x')$

- e.g., $c(w,z) = 5$

"cost" could always be 1, or inversely related to bandwidth, or inversely related to congestion

Cost of path $(x_1, x_2, x_3, \dots, x_p) = c(x_1, x_2) + c(x_2, x_3) + \dots + c(x_{p-1}, x_p)$

Question: What's the least-cost path between u and z ?

Routing algorithm: algorithm that finds **least-cost path**

Network Layer 4-4

Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
 - 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 Routing algorithms
 - Link state
 - Distance Vector
- 4.6 Routing in the Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast and multicast routing

Network Layer 4-6

Dijkstra's Algorithm: greedy

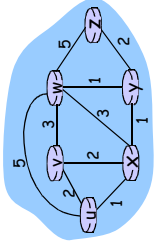
1 Initialization:

- 2 $N' = \{u\}$ // starting from source node
- 3 for all nodes v
- 4 if v adjacent to u
- 5 then $D(v) = c(u,v)$ // neighbor connected
- 6 else $D(v) = \infty$ // not connected
- 7
- 8 Loop
- 9 find w not in N' such that $D(w)$ is a **minimum** // greedy!
- 10 add w to N'
- 11 update $D(v)$ for all v adjacent to w and not in N' :
 $D(v) = \min(D(v), D(w) + c(w,v))$
- 12 J^* new cost to v is either old cost to v or known
- 13 **shortest path cost to w plus cost from w to v**
- 14 **until all nodes in N'**
- 15

Network Layer 4-8

Dijkstra's algorithm: example

Step		N'	D(v),p(v)	D(w),p(w)	D(x),p(x)	D(y),p(y)	D(z),p(z)
0		u	2,u	5,u	1,u	∞	∞
1		ux	2,u	4,x	2,x	∞	∞
2		uxy	2,u	3,y	2,x	4,y	∞
3		uxyv	2,u	3,y	2,x	4,y	4,y
4		uxyvw	2,u	3,y	2,x	4,y	4,y
5		uxyvwz	2,u	3,y	2,x	4,y	4,y



Network Layer 4-9

Dijkstra's algorithm, discussion

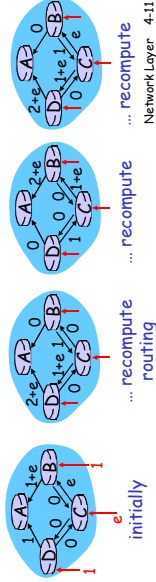
Algorithm complexity: n nodes

- each iteration: need to check all nodes, w, not in N
- $n(n+1)/2$ comparisons: $O(n^2)$
- more efficient implementations possible (using a heap): $O(n \log n)$

Oscillations possible (dynamically changing cost):

- e.g., link cost = amount of carried traffic

- Link costs are not symmetric**



Network Layer 4-11

Distance Vector Algorithm

Bellman-Ford Equation (dynamic programming)

Define

$d_x(y)$:= cost of least-cost path from x to y

Then

$$d_x(y) = \min_v \{c(x,v) + d_v(y)\}$$

where min is taken over all neighbors v of x

Network Layer 4-13

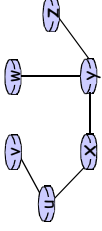
Distance Vector Algorithm

- $D_x(y)$ = estimate of least cost from x to y
- Node x knows cost to each neighbor v: $c(x,v)$
- Node x maintains distance vector $D_x = [D_x(y): y \in N]$
- Node x also maintains its neighbors' distance vectors
 - For each neighbor v, x maintains $D_v = [D_v(y): y \in N]$

Network Layer 4-15

Dijkstra's algorithm: example (2)

Resulting shortest-path tree from u:



Resulting forwarding table in u:

destination	link output port
v	(u,v)
x	(u,x)
y	(u,x)
w	(u,x)
z	(u,x)

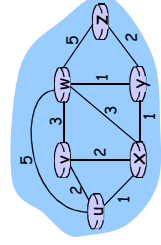
Network Layer 4-10

Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 Routing algorithms
 - Link state
 - Distance Vector
- 4.6 Routing in the Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast and multicast routing

Network Layer 4-12

Bellman-Ford example



B-F equation says:

$$d_u(z) = \min \{ c(u,v) + d_v(z), c(u,x) + d_x(z), c(u,w) + d_w(z) \}$$

$$= \min \{ 2 + 5, 1 + 3, 5 + 3 \} = 4$$

Clearly, $d_v(z) = 5$, $d_x(z) = 3$, $d_w(z) = 3$

Node that achieves minimum is next hop in shortest path \rightarrow forwarding table

Network Layer 4-14

Distance vector algorithm (4)

Basic idea:

- From time-to-time, each node sends its own distance vector estimate to neighbors
- Asynchronous 异步的
- When a node x receives new DV estimate from neighbor, it updates its own DV using B-F equation: $D_x(y) \leftarrow \min_v \{c(x,v) + D_v(y)\}$ for each node $y \in N$
- Under minor, natural conditions, the estimate $D_x(y)$ converge to the actual least cost $d_x(y)$
 - 最终会逼近实际距离

Network Layer 4-16

Distance Vector Algorithm (5)

Iterative, asynchronous:
each local iteration caused by:

- local link cost change
- DV update message from neighbor

Distributed:

- 邻居信息, 无需全局信息
- each node notifies neighbors *only* when its DV changes
 - neighbors then notify their neighbors if necessary

Each node:

wait for (change in local link cost or msg from neighbor)

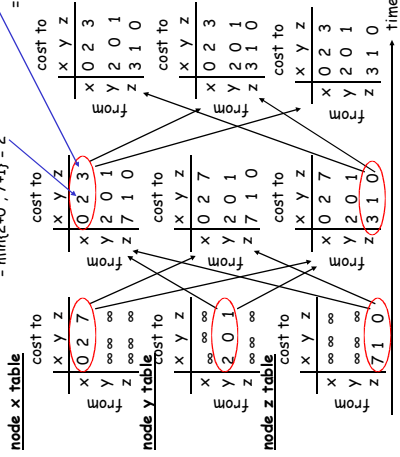
recompute estimates

if DV to any dest has changed, *notify* neighbors

Network Layer 4-17

$$D_x(y) = \min(c(x,y) + D_y(y), c(x,z) + D_z(y)) = \min(2+0, 7+1) = 2$$

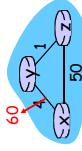
$$D_x(z) = \min(c(x,y) + D_y(z), c(x,z) + D_z(z)) = \min(2+1, 7+0) = 3$$



Network Layer 4-19

Distance Vector: link cost changes

- Link cost changes:
 - good news travels fast
 - bad news travels slow
 - if link cost increases, all nodes must propagate the new cost
 - if link cost decreases, only nodes on the path to the destination need to propagate the new cost
- Routing reverse:
 - If Z routes through Y to get to X:
 - If Z links Y its (C(z,y) distance to X is infinite, so Y won't route to X via Z)
 - all links completely route cost to all early precursors



4->60, but z does not know;
z->x is still 5
Then, y->x becomes 5, via z: y tells z that d(y,x)=5
Then, z->x becomes 6, via y: z tells y that d(z,x)=6
This repeats until the path cost larger than 50

Network Layer 4-21

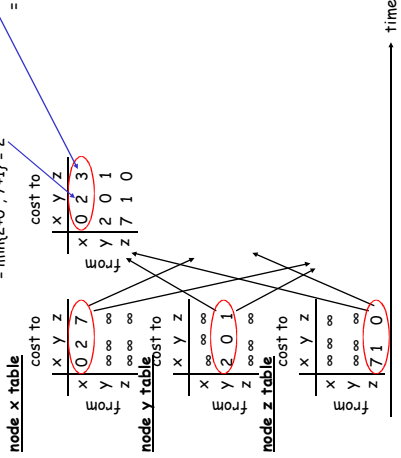
Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
 - 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 Routing algorithms
 - Link state
 - Distance Vector
- 4.6 Routing in the Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast and multicast routing

Network Layer 4-23

$$D_x(y) = \min(c(x,y) + D_y(y), c(x,z) + D_z(y)) = \min(2+0, 7+1) = 2$$

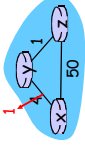
$$D_x(z) = \min(c(x,y) + D_y(z), c(x,z) + D_z(z)) = \min(2+1, 7+0) = 3$$



Network Layer 4-18

Distance Vector: link cost changes

- Link cost changes:
 - node detects local link cost change
 - updates routing info, recalculates distance vector
 - if DV changes, notify neighbors



"good news travels fast"

At time t_0 , y detects the link-cost change, updates its DV, and informs its neighbors.
At time t_1 , z receives the update from y and updates its DV. It computes a new least cost to x and sends its neighbors its DV.
At time t_2 , y receives z's update and updates its distance table. y's least costs do not change and hence y does not send any message to z.

Network Layer 4-20

Comparison of LS and DV algorithms

Message complexity
Robustness: what happens if router malfunctions?

- LS: with n nodes, E links, $O(nE)$ msgs sent
- DV: exchange between neighbors only
- convergence time varies
- Speed of Convergence
 - LS: $O(n^2)$ algorithm requires $O(nE)$ msgs
 - may have oscillations
 - DV: convergence time varies
 - may be routing loops
 - count-to-infinity problem
- LS: node can advertise incorrect *path* cost
- DV: each node's table used by others
 - error propagate thru network
- LS: node can advertise incorrect *link* cost
- DV: each node computes only its own table

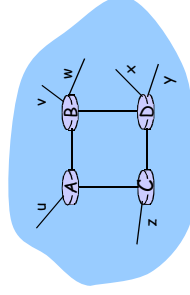
Network Layer 4-22

RIP (Routing Information Protocol)

- (Local, decentralized) distance vector (DV) algorithm
- distance metric: # of hops (max = 15 hops) 距离

From router A to subnets:

destination	hops
subnet u	1
subnet v	2
subnet w	2
subnet x	3
subnet y	3
subnet z	2



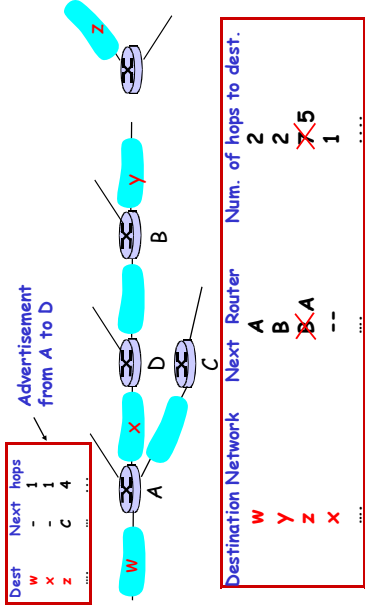
Network Layer 4-24

RIP advertisements

- *distance vectors*: exchanged among neighbors every 30 sec via Response Message (also called **advertisement**)
- 每 30 秒更新一次，邻居路由器之间的信息交换
- each advertisement: list of up to 25 destination subnets

Network Layer 4-25

RIP: Example



Routing/Forwarding table in D

Network Layer 4-27

OSPF "advanced" features (not in RIP)

- **security**: all OSPF messages authenticated (to prevent malicious intrusion, **net layer auth. using MD5 hash**)
- multiple same-cost **paths** allowed (only one path in RIP)
- Integrated uni- and **multicast** support:
 - Multicast OSPF (MOSPF) uses same topology data base as OSPF
- **hierarchical** OSPF in large domains.

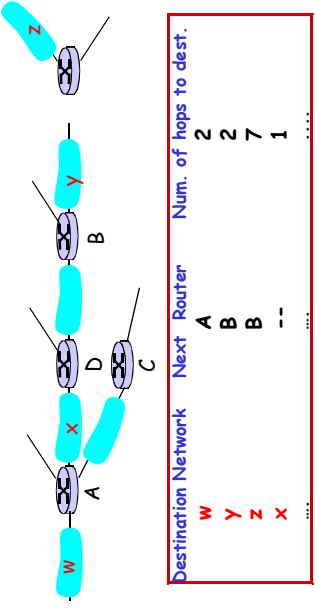
Network Layer 4-29

Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
 - 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 Routing algorithms
 - Link state
 - Distance Vector
 - Hierarchical routing
- 4.6 Routing in the Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast and multicast routing

Network Layer 4-31

RIP: Example



Routing/Forwarding table in D

Network Layer 4-26

OSPF (Open Shortest Path First)

- "open": publicly available
- uses **Link State algorithm**
 - LS packet dissemination
 - topology map at each node
 - route computation using **Dijkstra's** algorithm
- OSPF advertisement carries one entry per neighbor router
- advertisements disseminated to **entire AS** via **flooding 开销太大**
 - carried in OSPF messages directly over IP (rather than TCP or UDP)

Network Layer 4-28

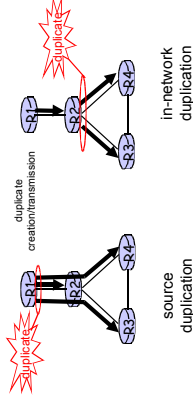
Internet inter-AS routing: BGP

- **BGP (Border Gateway Protocol)**: *the de facto standard 事实上在用的标准*
- BGP provides each AS a means to:
 1. Obtain subnet reachability information from neighboring ASs.
 2. Propagate reachability information to all AS-internal routers.
 3. Determine "good" routes to subnets based on reachability information and policy.
- allows subnet to advertise its existence to rest of Internet: *"I am here"*

Network Layer 4-30

Broadcast (广播) Routing

- deliver packets from source to all other nodes
- source duplication is inefficient:



- source duplication: how does source determine recipient addresses?

Network Layer 4-32

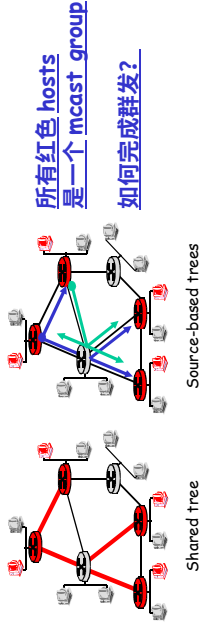
In-network duplication

- flooding: when node receives brdcst pkt, sends copy to all neighbors (**infinite cascading**)
 - Problems: cycles & broadcast storm
- controlled flooding: node only brdcsts pkt if it hasn't brdcst same packet before (只发给没收到的)
 - Node keeps track of pkt ids already brdcsted
 - Or reverse path forwarding (RPF): only forward pkt if it arrived on shortest path between node and source
- spanning tree (最优解)
 - No redundant packets received by any node

Network Layer 4-33

Multicast Routing: Problem Statement

- Goal:** find a tree (or trees) connecting routers having local mcast group members
 - tree:** not all paths between routers used
 - source-based:** different tree from each sender to rcvrs
 - shared-tree:** same tree used by all group members



Source-based trees

Shared tree

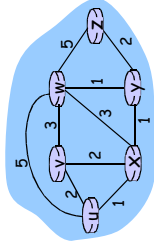
Network Layer 4-3

Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
 - 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 **Routing algorithms**
 - Link state
 - Distance Vector
- 4.6 Routing in the Internet
 - RIP
 - OSPF
 - BGP
- 4.7 Broadcast and multicast routing

Network Layer 4-1

Graph abstraction



Graph: $G = (N, E)$

N = set of routers (nodes) = $\{u, v, w, x, y, z\}$

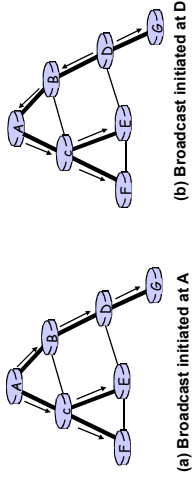
E = set of links (edges) = $\{(u,v), (u,x), (v,x), (v,w), (x,w), (x,y), (w,y), (w,z), (y,z)\}$

Remark: Graph abstraction is useful in other network contexts
Example: P2P, where N is set of peers and E is set of TCP connections

Network Layer 4-3

Spanning Tree

- First construct a spanning tree
- Nodes forward copies only along spanning tree 沿着 tree 发送以避免重复



(a) Broadcast initiated at A

(b) Broadcast initiated at D

B 和 C 之间的 broadcast pkt 无需再转发

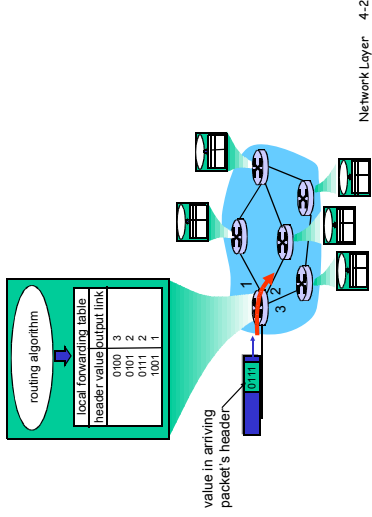
Network Layer 4-34

Chapter 4: summary

- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 **What's inside a router**
 - 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 Routing algorithms
 - Link state
 - Distance Vector
- 4.6 Routing in the Internet
 - RIP
 - OSPF
 - BGP
- 4.7 **Broadcast and multicast routing**

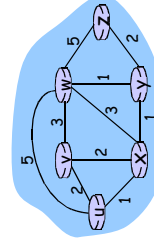
Network Layer 4-36

Interplay between routing, forwarding



Network Layer 4-2

Graph abstraction: costs



$c(x, x') = \text{cost of link } (x, x')$

- e.g., $c(w, z) = 5$

cost could always be 1, or inversely related to bandwidth, or inversely related to congestion

Cost of path $(x_1, x_2, x_3, \dots, x_p) = c(x_1, x_2) + c(x_2, x_3) + \dots + c(x_{p-1}, x_p)$

Question: What's the least-cost path between u and z?

Routing algorithm: algorithm that finds **least-cost path**

Network Layer 4-4

Routing Algorithm classification

Global or local (decentralized) information?

Global:

- all routers have complete topology, link cost info
- "link state" algorithms
- Local/Decentralized:
 - router knows physically-connected neighbors, link costs to neighbors
 - iterative process of computation, exchange of info with neighbors
 - "distance vector" algorithms

Static or dynamic?

- Static:
- routes change slowly over time
- Dynamic:
- routes change more quickly
 - periodic update
 - in response to link cost changes

Network Layer 4-5

A Link-State (Global) Routing Algorithm

Dijkstra's algorithm

- net topology, link costs known to all nodes
 - accomplished via "link state broadcast"
 - all nodes have same info
- computes least cost paths from one node ("source") to all other nodes
 - gives **forwarding table** for that node
- iterative: after k iterations, know least cost path to k dest.'s

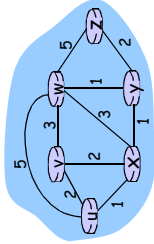
Notation:

- $c(x,y)$: link cost from node x to y; $= \infty$ if not direct neighbors
- $D(v)$: current value of cost of path from source to dest. v
- $p(v)$: predecessor node along path from source to v
- N' : set of nodes whose least cost path definitively known

Network Layer 4-7

Dijkstra's algorithm: example

Step		N'	D(v),p(v)	D(w),p(w)	D(x),p(x)	D(y),p(y)	D(z),p(z)
0		U	2,u	5,u	1,u	∞	∞
1	UX		2,u	4,x	2,x	4,y	∞
2	UXV		2,u	3,v	2,x	4,y	4,y
3	UXVW		2,u	3,v	2,x	4,y	4,y
4	UXVWZ		2,u	3,v	2,x	4,y	4,y
5	UXVWZ		2,u	3,v	2,x	4,y	4,y



Network Layer 4-9

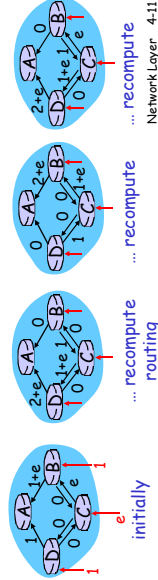
Dijkstra's algorithm, discussion

Algorithm complexity: n nodes

- each iteration: need to check all nodes, w, not in N
- $n(n+1)/2$ comparisons: $O(n^2)$
- more efficient implementations possible (using a heap): $O(n \log n)$

Oscillations possible (dynamically changing cost):

- e.g. link cost = amount of carried traffic
- Link costs are not symmetric**



Network Layer 4-11

Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.7 Broadcast and multicast routing
- 4.5 Routing algorithms
 - Link state
 - Distance Vector
- 4.6 Routing in the Internet
 - RIP
 - OSPF
 - BGP

Network Layer 4-6

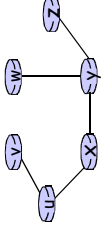
Dijkstra's Algorithm: greedy

- Initialization:**
- $N' = \{u\}$ // starting from source node
- for all nodes v
- if v adjacent to u
- then $D(v) = c(u,v)$ // neighbor connected
- else $D(v) = \infty$ // not connected
-
- Loop**
- find w not in N' such that $D(w)$ is a **minimum** // greedy!
- add w to N'
- update $D(v)$ for all v adjacent to w and not in N' :
 - $D(v) = \min(D(v), D(w) + c(w,v))$
 - // new cost to v is either old cost to v or known
 - shortest path cost to w plus cost from w to v**
- until all nodes in N'**

Network Layer 4-8

Dijkstra's algorithm: example (2)

Resulting shortest-path tree from u:



Resulting forwarding table in u:

destination	link output port
v	(u,v)
x	(u,x)
y	(u,x)
w	(u,x)
z	(u,x)

Network Layer 4-10

Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
- 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.7 Broadcast and multicast routing
- 4.5 Routing algorithms
 - Link state
 - Distance Vector
- 4.6 Routing in the Internet
 - RIP
 - OSPF
 - BGP

Network Layer 4-12

Distance Vector Algorithm

Bellman-Ford Equation (dynamic programming)

Define

$d_x(y) :=$ cost of least-cost path from x to y

Then

$$d_x(y) = \min \{c(x,v) + d_v(y)\}$$

where min is taken over all neighbors v of x

Network Layer 4-13

Distance Vector Algorithm

- $D_x(y)$ = estimate of least cost from x to y
- Node x knows cost to each neighbor v : $c(x,v)$
- Node x maintains distance vector $D_x = [D_x(y): y \in N]$
- Node x also maintains its neighbors' distance vectors
 - For each neighbor v , x maintains $D_v = [D_v(y): y \in N]$

Network Layer 4-15

Distance Vector Algorithm (5)

Iterative, asynchronous:

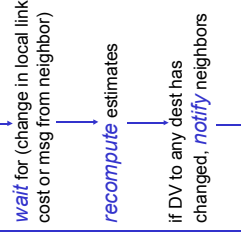
each local iteration caused by:

- local link cost change
- DV update message from neighbor

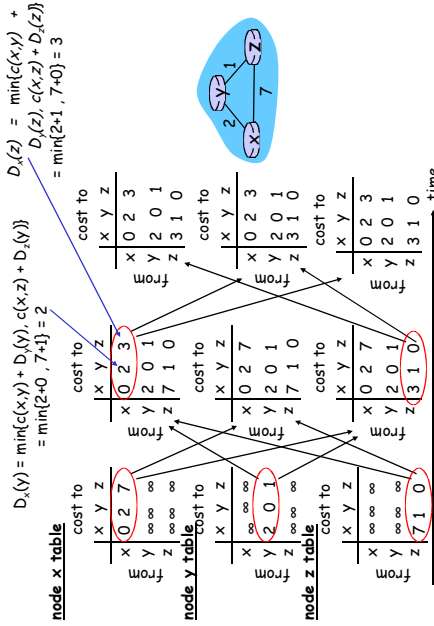
Distributed:

- 邻居信息，无需全局信息
- each node notifies neighbors *only* when its DV changes
 - neighbors then notify their neighbors if necessary

Each node:

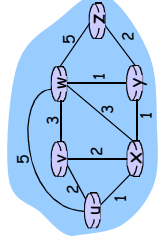


Network Layer 4-17



Network Layer 4-19

Bellman-Ford example



B-F equation says:

$$d_u(z) = \min \{ c(u,v) + d_v(z), c(u,x) + d_x(z), c(u,w) + d_w(z) \}$$

$$= \min \{ 2 + 5, 1 + 3, 5 + 3 \} = 4$$

Clearly, $d_u(z) = 5$, $d_x(z) = 3$, $d_w(z) = 3$

Node that achieves minimum is next hop in shortest path → forwarding table

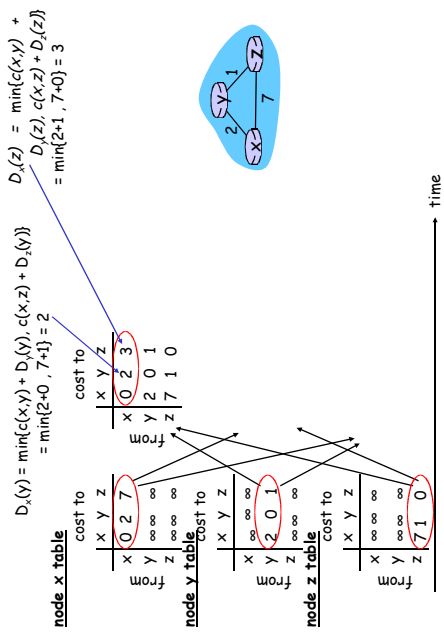
Network Layer 4-14

Distance vector algorithm (4)

Basic idea:

- From time-to-time, each node sends its own distance vector estimate to neighbors
- Asynchronous 异步的
- When a node x receives new DV estimate from neighbor, it updates its own DV using B-F equation: $D_x(y) \leftarrow \min_v \{c(x,v) + D_v(y)\}$ for each node $y \in N$
- Under minor, natural conditions, the estimate $D_x(y)$ converge to the actual least cost $d_x(y)$
 - 最终会逼近实际距离

Network Layer 4-16

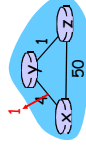


Network Layer 4-18

Distance Vector: link cost changes

Link cost changes:

- node detects local link cost change
- updates routing info, recalculates distance vector
- if DV changes, notify neighbors

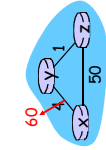


"good news travels fast"

At time t_0 , y detects the link-cost change, updates its DV, and informs its neighbors.
 At time t_1 , z receives the update from y and updates its table. It computes a new least cost to x and sends its neighbors its DV.
 At time t_2 , x receives z 's update and updates its distance table. x 's least costs do not change and hence x does *not* send any message to z .

Network Layer 4-20

Distance Vector: link cost changes



- Link cost changes:
 - good news travels fast
 - bad news travels slow - **count to infinity problem**
 - 24 iterations before algorithm stabilizes once first
- Poisoned reverse:
 - If Z routes through Y to get to X, Y will advertise infinite cost to X via Z
 - will this completely solve count to infinity problem?

4→60, but z does not know;
z→x is still 5
Then, y→x becomes 5, via z: y
tells z that $d(y,x)=5$
Then, z→x becomes 6, via y: z
tells y that $d(z,x)=6$
This repeats until the path cost
larger than 50

Network Layer 4-21

Chapter 4: Network Layer

- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
 - 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
 - 4.5 Routing algorithms
 - Link state
 - Distance Vector
 - 4.6 Routing in the Internet
 - RIP
 - OSPF
 - BGP
 - 4.7 Broadcast and multicast routing

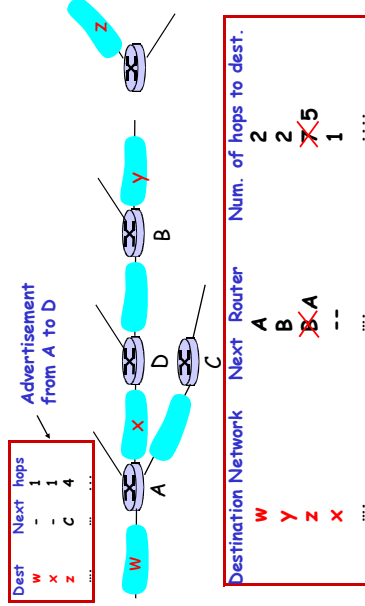
Network Layer 4-23

RIP advertisements

- distance vectors: exchanged among neighbors every 30 sec via Response Message (also called **advertisement**)
- 每 30 秒更新一次，邻居路由器之间的信息交换
- each advertisement: list of up to 25 destination subnets

Network Layer 4-25

RIP: Example



Routing/Forwarding table in D

Network Layer 4-27

Comparison of LS and DV algorithms

- Message complexity
 - LS: with n nodes, E links, $O(nE)$ msgs sent
 - DV: exchange between neighbors only
 - convergence time varies
- Speed of Convergence
 - LS: $O(n^2)$ algorithm requires $O(nE)$ msgs
 - may have oscillations
 - DV: convergence time varies
 - may be routing loops
 - count-to-infinity problem
- Robustness: what happens if router malfunctions?
 - LS:
 - node can advertise incorrect link cost
 - each node computes only its own table
 - DV:
 - DV node can advertise incorrect path cost
 - each node's table used by others
 - error propagate thru network

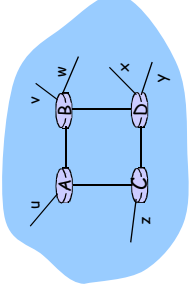
Network Layer 4-22

RIP (Routing Information Protocol)

- (Local, decentralized) distance vector (DV) algorithm
- distance metric: # of hops (max = 15 hops) 距离

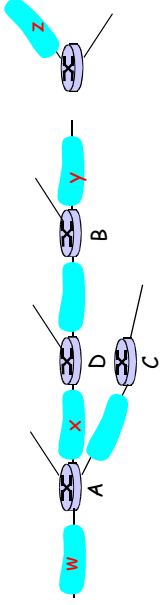
From router A to subnets:

destination	hops
subnet u	1
subnet v	2
subnet w	2
subnet x	3
subnet y	3
subnet z	2



Network Layer 4-24

RIP: Example



Destination Network	Next Router	Num. of hops to dest.
W	A	2
Y	B	2
Z	B	7
X	--	1

Routing/Forwarding table in D

Network Layer 4-26

OSPF (Open Shortest Path First)

- "open": publicly available
- uses Link State algorithm
 - L5 packet dissemination
 - topology map at each node
 - route computation using Dijkstra's algorithm
- OSPF advertisement carries one entry per neighbor router
- advertisements disseminated to entire AS via flooding 开销太大
 - carried in OSPF messages directly over IP (rather than TCP or UDP)

Network Layer 4-28

OSPF "advanced" features (not in RIP)

- **security**: all OSPF messages authenticated (to prevent malicious intrusion, net layer auth. using MD5 hash)
- **multiple same-cost paths** allowed (only one path in RIP)
- Integrated uni- and **multicast** support:
 - Multicast OSPF (MOSPF) uses same topology data base as OSPF
- **hierarchical** OSPF in large domains.

Network Layer 4-29

Chapter 4: Network Layer

- 4.1 Introduction
 - 4.5 Routing algorithms
 - Link state
 - Distance Vector
 - Hierarchical routing
- 4.2 Virtual circuit and datagram networks
- 4.3 What's inside a router
 - 4.6 Routing in the Internet
 - RIP
 - OSPF
 - BGP
- 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- **4.7 Broadcast and multicast routing**

Network Layer 4-31

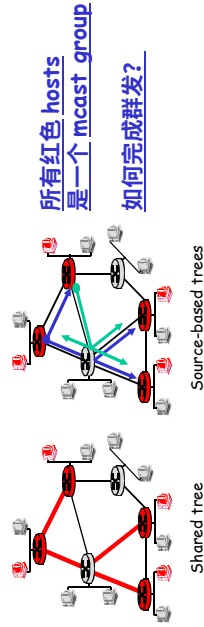
In-network duplication

- flooding: when node receives brdcst pkt, sends copy to all neighbors (Infinite cascading)
 - Problems: cycles & broadcast storm
- controlled flooding: node only brdcsts pkt if it hasn't brdcst same packet before (**只发给没收到的**)
 - Node keeps track of pkt ids already brdcsted
 - Or reverse path forwarding (RPF): only forward pkt if it arrived on shortest path between node and source
- spanning tree (**最优解**)
 - No redundant packets received by any node

Network Layer 4-33

Multicast Routing: Problem Statement

- **Goal**: find a tree (or trees) connecting routers having local mcast group members
 - **Tree**: not all paths between routers used
 - **source-based**: different tree from each sender to rcvrs
 - **shared-tree**: same tree used by all group members



Source-based trees

Shared tree

Network Layer 4-36

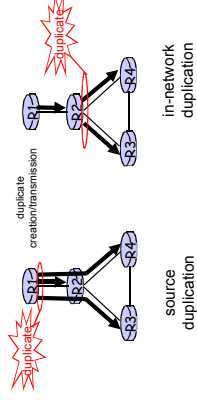
Internet inter-AS routing: BGP

- **BGP (Border Gateway Protocol)**: the de facto standard **事实上在用的标准**
- BGP provides each AS a means to:
 1. Obtain subnet reachability information from neighboring ASs.
 2. Propagate reachability information to all AS-internal routers.
 3. Determine "good" routes to subnets based on reachability information and policy.
- allows subnet to advertise its existence to rest of Internet: **"I am here"**

Network Layer 4-30

Broadcast (广播) Routing

- deliver packets from source to all other nodes
- source duplication is inefficient:

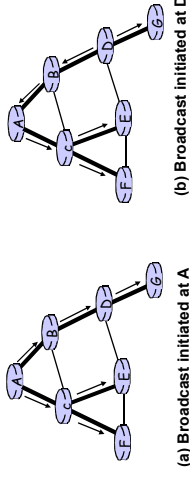


- source duplication: how does source determine recipient addresses?

Network Layer 4-32

Spanning Tree

- First construct a spanning tree
- Nodes forward copies only along spanning tree **沿着 tree 发送以避免重复**



B和C之间的broadcast pkt 无需再转发

Network Layer 4-34

Chapter 4: summary

- 4.1 Introduction
- 4.2 Virtual circuit and datagram networks
- 4.3 **What's inside a router**
 - 4.4 IP: Internet Protocol
 - Datagram format
 - IPv4 addressing
 - ICMP
 - IPv6
- 4.5 Routing algorithms
 - **Link state**
 - **Distance Vector**
- 4.6 Routing in the Internet
 - RIP
 - OSPF
 - BGP
- **4.7 Broadcast and multicast routing**

File into TCP segments

- ❑ **rdt: a pkt stream; each pkt has a seq. #**
- ❑ **TCP: a byte stream; each byte has seq. #**

Q: how a file (500,000 bytes) divided into TCP segments? MSS = 1000 bytes.

- ❑ 500,000 bytes can create 500 segments.
- ❑ 1st byte seq. # is 0, the last byte seq. # is 499,999
- ❑ Segment's seq. # is its 1st byte's seq. #
 - Ex. 1st segment seq # = 0; 2nd segment seq # = 1000

Transport Layer 3-19

TCP Round Trip Time and Timeout

Q: how to set TCP

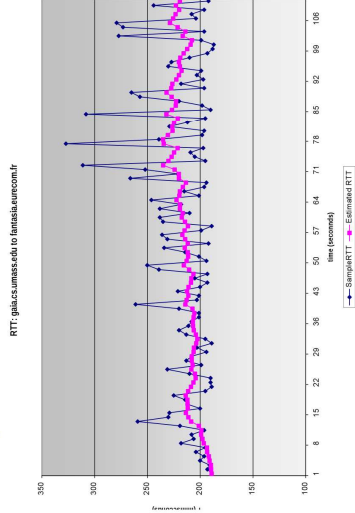
- ❑ timeout value?
 - longer than RTT
 - but RTT varies
- ❑ too short: premature timeout
 - unnecessary retransmissions
- ❑ too long: slow reaction to segment loss

Q: how to estimate RTT?

- ❑ **sampleRTT**: measured time from segment transmission until ACK receipt
 - ignore retransmissions
- ❑ **sampleRTT** will vary, want estimated RTT "smoother"
 - **average several recent measurements**, not just current **sampleRTT**

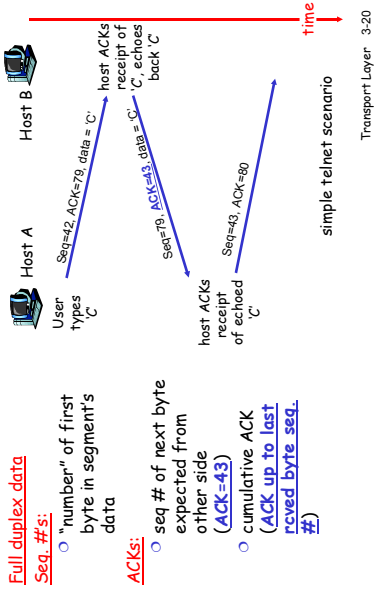
Transport Layer 3-21

Example RTT estimation:



Est. RTT can roughly reflect the real RTT value. Transport Layer 3-23

Telnet: TCP seq. #'s and ACKs



TCP Round Trip Time and Timeout

$$\text{EstimatedRTT} = (1 - \alpha) * \text{EstimatedRTT} + \alpha * \text{sampleRTT}$$

- ❑ Exponential weighted moving average
- ❑ influence of past sample decreases exponentially fast
- ❑ typical value: $\alpha = 0.125$
- ❑ **Heuristic: May not be the best**

Transport Layer 3-22

TCP Round Trip Time and Timeout

Setting the timeout

- ❑ **EstimatedRTT** plus "safety margin"
 - large variation in **EstimatedRTT** → larger safety margin
- ❑ first estimate of how much **sampleRTT** deviates from **EstimatedRTT**:

$$\text{DevRTT} = (1 - \beta) * \text{DevRTT} + \beta * |\text{sampleRTT} - \text{EstimatedRTT}|$$

(typically, $\beta = 0.25$) **统计学概念: 指数加权移动平均 EWMA**
(**Exponentially Weighted Moving Average**)

Then set timeout interval: **Weighted approach**

$$\text{TimeoutInterval} = \text{EstimatedRTT} + 4 * \text{DevRTT}$$

Transport Layer 3-24

rdt = reliable data transfer

- Reliable**
- Data + **reply** (ACK/NAK)
- Data pkt, or ACK may be
 - Corrupted (Bit error) → checksum
 - Lost → seq. # to prevent duplicate receptions
 - Delayed → count down timer

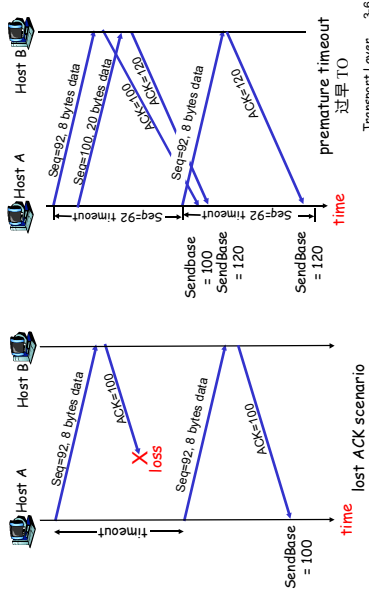
Transport Layer 3-1

TCP reliable data transfer

- Retransmissions are triggered by:
 - timeout events
 - duplicate acks
- Initially consider **simplified/ simplest** TCP sender:
 - ignore duplicate acks
 - ignore flow control, congestion control
- Pipelined segments
- Cumulative acks
 - Ack up to the last received byte
- TCP uses single retransmission **timer**

Transport Layer 3-3

TCP: retransmission scenarios



Transport Layer 3-6

TCP ACK generation [RFC 1122, RFC 2581]

Event at Receiver	TCP Receiver action
Arrival of in-order segment with expected seq #. All data up to expected seq # already ACKed	Delayed ACK. Wait up to 500ms for next segment. If no next segment, send ACK (网络通畅下一个也马上来)
Arrival of in-order segment with expected seq #. One other segment has ACK pending	Immediately send single cumulative ACK. ACKing both in-order segments (正要发ACK, 来了一个新segment, 一次两个)
Arrival of out-of-order segment higher-than-expected seq. #. Gap detected	Immediately send duplicate ACK , indicating seq. # of next expected byte (顺序错误, 中间缺一段)
Arrival of segment that partially or completely fills gap	Immediate send ACK, provided that segment starts at lower end of gap (顺序错误, 中间缺一段, 从底部开始补)

Transport Layer 3-8

Chapter 3 outline

- 3.1 Transport-layer services
- 3.2 Connectionless transport: UDP
- 3.3 Connection-oriented transport: TCP
 - segment structure
 - TCP's reliable data transfer
 - flow control
 - connection management
- 3.4 Principles of congestion control
- 3.5 TCP congestion control

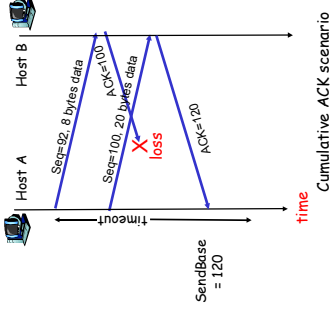
Transport Layer 3-2

TCP sender events:

- data rcvd from app:**
- Create segment with seq #
- seq # is byte-stream number of first data byte in segment
- start timer if not already running (think of **timer as for oldest un-acked segment**)
- expiration interval: TimeoutInterval
- timeout:**
- retransmit segment that caused timeout
- restart timer
- Ack rcvd:**
- If Rx acknowledges previously unacked segments
 - update what is known to be acked
 - **start a new timer** for a new outstanding segments

Transport Layer 3-4

TCP retransmission scenarios (more)



Transport Layer 3-7

Fast Retransmit=Re-Tx b4 TO

- **Re-transmit before timeout**
- If segment is lost, there will likely be many duplicate ACKs.
- Trigger event for fast re-Tx: If sender receives 3 ACKs for the same data, it supposes that segment after ACKed data was lost:
 - **fast retransmit: resend segment before timer expires**
- Time-out period often relatively long:
 - long delay before resending lost packet
- Detect lost segments via duplicate ACKs.
 - Sender often sends many segments back-to-back (sender 一个接一个, 背靠背地发送 segment)

Transport Layer 3-9

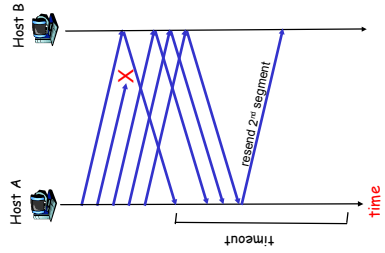


Figure 3.37 Resending a segment after triple duplicate ACK Layer 3-10

Fast retransmit algorithm:

```

event: ACK received, with ACK field value of y
if (y > SendBase) {
    SendBase = y
    if (there are currently not-yet-acknowledged segments)
        start timer
}
else {
    increment count of dup ACKs received for y
    if (count of dup ACKs received for y = 3) {
        resend segment with sequence number y
    }
}

```

a duplicate ACK for already ACKed segment

fast retransmit

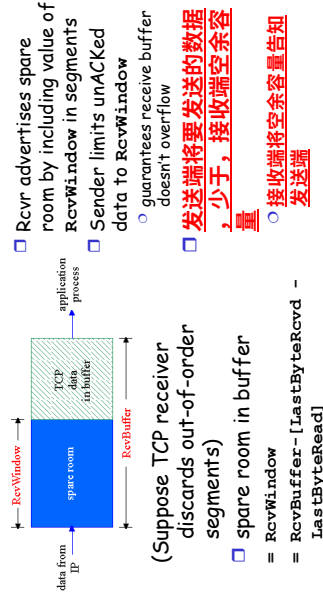
Transport Layer 3-11

Chapter 3 outline

- 3.1 Transport-layer services
- 3.2 Multiplexing and demultiplexing
- 3.3 Connectionless transport: UDP
- 3.4 Principles of reliable data transfer
- 3.5 Connection-oriented transport: TCP
 - segment structure
 - reliable data transfer
 - flow control
 - connection management
- 3.6 Principles of congestion control
- 3.7 TCP congestion control

Transport Layer 3-12

TCP Flow control: how it works



Transport Layer 3-14

TCP Connection Management

Three way handshake:

- Recall:** TCP sender, receiver establish "connection" before exchanging data segments
- initialize TCP variables:
 - seq. #s
 - buffers, flow control info (e.g. RcvrWindow)
- client: connection initiator


```

Socket clientSocket = new Socket("hostname", "port number");

```
- server: contacted by client


```

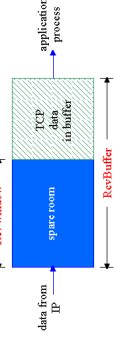
Socket connectionSocket = weIcomeSocket.accept();

```
- Step 1:** client host sends TCP SYN segment to server
 - specifies initial seq #
 - no data
- Step 2:** server host receives SYN, replies with SYNACK segment
 - server allocates buffers
 - specifies server initial seq. #
- Step 3:** client receives SYNACK, replies with ACK segment, which may contain data

Transport Layer 3-16

TCP Flow Control

- receive side of TCP connection has a receive buffer:



- app process may be slow at reading from buffer

- flow control: sender won't overflow receiver's buffer by transmitting too much, too fast
- speed-matching service: matching the send rate to the receiving app's drain rate
- 使发送端的发送速度, 尽量匹配, 接收端的读取速度

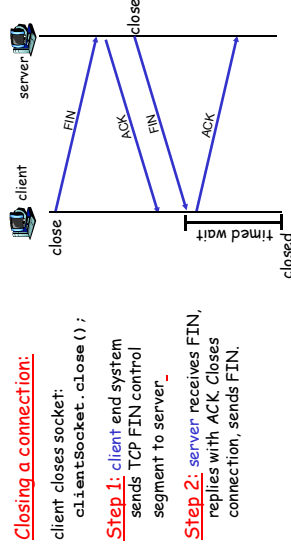
Transport Layer 3-13

Chapter 3 outline

- 3.1 Transport-layer services
- 3.2 Multiplexing and demultiplexing
- 3.3 Connectionless transport: UDP
- 3.4 Principles of reliable data transfer
- 3.5 Connection-oriented transport: TCP
 - segment structure
 - reliable data transfer
 - flow control
 - connection management
- 3.6 Principles of congestion control
- 3.7 TCP congestion control

Transport Layer 3-15

TCP Connection Management (cont.)



Closing a connection:

- client closes socket:

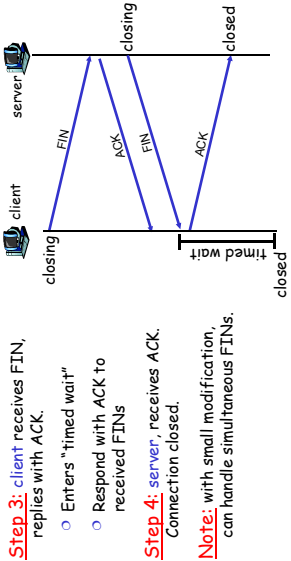

```

clientSocket.close();

```
- Step 1:** client end system sends TCP FIN control segment to server
- Step 2:** server receives FIN, replies with ACK. Closes connection, sends FIN.

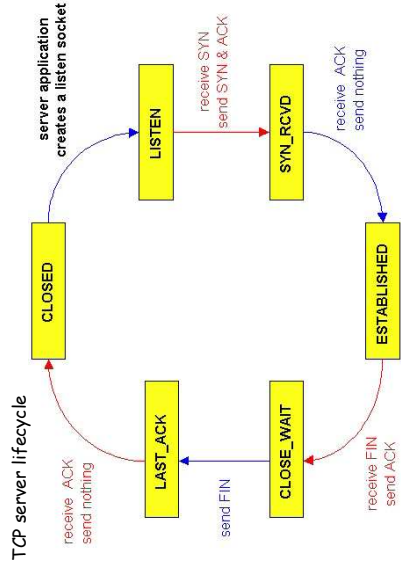
Transport Layer 3-17

TCP Connection Management (cont.)



Transport Layer 3-18

TCP Connection Management (cont.)



Principles of Congestion Control

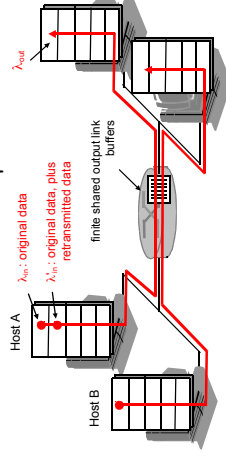
Congestion:

- informally: "too many sources sending too much data too fast for *network* to handle"
- different from flow control!
 - Flow control**: 接收端能力有限，忙不过来
 - Congestion control**: 网络太拥挤，发不过来
- manifestations:
 - lost packets (buffer overflow at routers)
 - long delays (queueing in router buffers)
 - a top-10 problem!**

Transport Layer 3-22

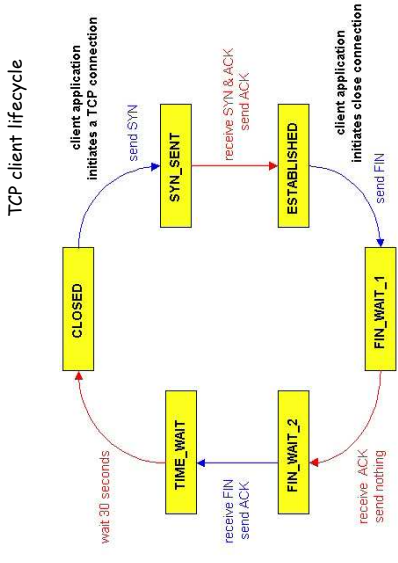
Causes/costs of congestion: scenario 2

- Case 2: one router, *finite* buffers**
- Pkt loss at router
- sender retransmission of lost packet



Transport Layer 3-24

TCP Connection Management (cont.)



Chapter 3 outline

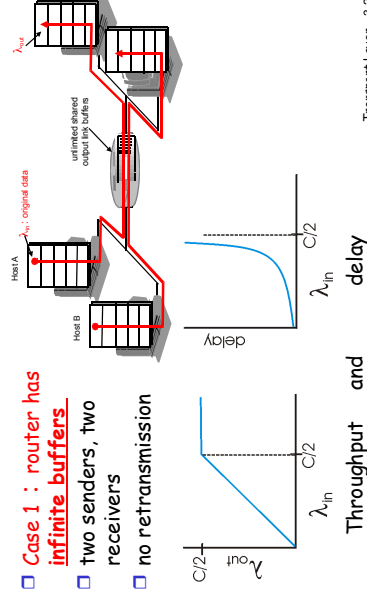
- 3.1 Transport-layer services
- 3.2 Multiplexing and demultiplexing
- 3.3 Connectionless transport: UDP
- 3.4 Principles of reliable data transfer
- 3.5 Connection-oriented transport: TCP
 - segment structure
 - reliable data transfer
 - flow control
 - connection management
- 3.6 Principles of congestion control
- 3.7 TCP congestion control

Transport Layer 3-21

Causes/costs of congestion: scenario 1

Case 1: router has infinite buffers

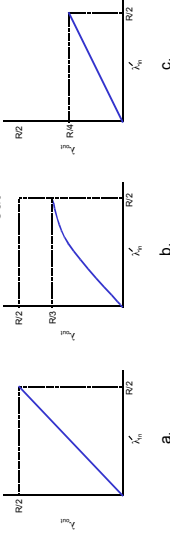
- two senders, two receivers
- no retransmission



Transport Layer 3-23

Causes/costs of congestion: scenario 2

- always: $\lambda_{in} = \lambda_{out}$ (goodput)
- "perfect" retransmission only when loss: $\lambda'_{in} > \lambda_{out}$, λ'_{in} larger
- retransmission of delayed (not lost) packet makes λ'_{in} larger (than perfect case) for same λ_{out}



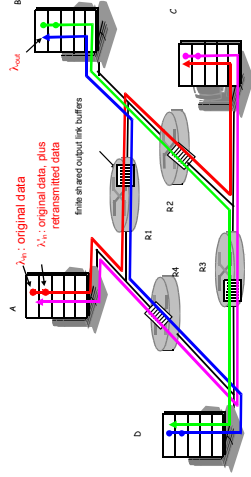
- "costs" of congestion:
 - more work (retrans) for given "goodput"
 - unneded retransmissions: link carries multiple copies of pkt

Transport Layer 3-25

Causes/costs of congestion: scenario 3

- four senders
- multi-hop paths
- timeout/retransmit

Q: what happens as λ_{in} and λ'_{in} increase?



Flow B-D outcompetes flow A-C at R2!!!

Transport Layer 3-26

Approaches towards congestion control

Two broad approaches towards congestion control:

Network-assisted congestion control:

- 网络辅助拥塞控制
- 网络路由转发信息辅助
- routers provide feedback to end systems
 - single bit indicating congestion
 - explicit rate sender should send at

End-end congestion control:

- 端到端拥塞控制
- 接收端 ACK 辅助，发送端猜测
- no explicit feedback from network
- congestion inferred from end-system observed loss, delay
- approach taken by TCP

Transport Layer 3-28

TCP

- Rdt
 - ACK
 - Checksum
 - Seq. #
 - Timer
 - Flow control
- TCP
 - Flow control
 - TCP Congestion control
 - Additive increase, multiplicative decrease (AIMD 线性增加、成倍递减)
 - Congestion window (CongWin)
 - Slow start
 - Collision avoidance
 - Threshold
 - Fast re-tx and fast recovery

Transport Layer 3-30

TCP Congestion Control: details

- Window: sender limits transmission such as $\text{LastByteSent} - \text{LastByteAcks} \leq \text{CongWin}$
- Roughly,

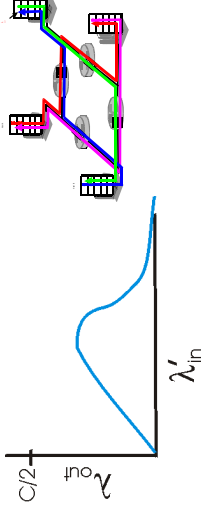
$\text{rate} = \frac{\text{CongWin}}{\text{RTT}} \text{ Bytes/sec}$
- CongWin is dynamic, function of perceived network congestion

How does sender perceive congestion/loss?

- Timeout or 3 duplicate acks
- TCP sender reduces rate (CongWin) after loss event
- three mechanisms:
 - AIMD
 - slow start
 - conservative after timeout events

Transport Layer 3-32

Causes/costs of congestion: scenario 3



Another "cost" of congestion:

- when packet dropped, any "upstream" transmission capacity used for that packet was wasted!

Flow A-C: if pkt loss at R2, then resource forwarding this pkt at R1 is wasted!!!!

Transport Layer 3-27

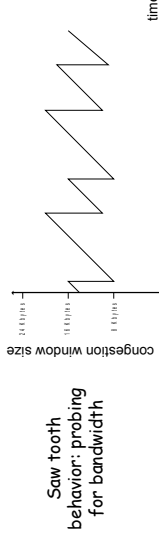
Chapter 3 outline

- 3.1 Transport-layer services
- 3.2 Multiplexing and demultiplexing
- 3.3 Connectionless transport: UDP
- 3.4 Principles of reliable data transfer
- 3.5 Connection-oriented transport: TCP
 - segment structure
 - reliable data transfer
 - flow control
 - connection management
- 3.6 Principles of congestion control
- 3.7 TCP congestion control

Transport Layer 3-29

TCP congestion control: additive increase, multiplicative decrease (AIMD)

- Approach: increase transmission rate (window size), probing for usable bandwidth; until loss occurs
 - AI, additive increase: increase CongWin by 1 MSS every RTT until loss detected
 - MD, multiplicative decrease: cut CongWin in half after loss

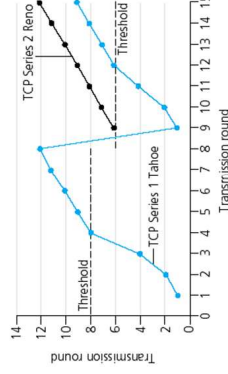


Congestion window (CongWin) = # of bytes that are allowed to be sent in next RTT

Transport Layer 3-31

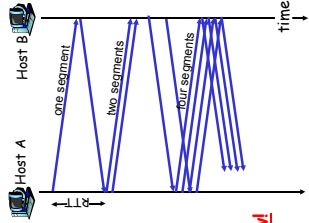
TCP Slow Start

- When connection begins, $\text{CongWin} = 1 \text{ MSS}$
 - If $\text{MSS} = 500 \text{ bytes}$, and $\text{RTT} = 200 \text{ msec}$
 - initial rate = $500 / 0.2 = 2500 \text{ bytes/sec}$
 - = 20 kb/s
- Available bandwidth much greater than MSS/RTT
- Thus quickly increase the CongWin by 1 MSS for each received ACK
 - 1MSS \rightarrow 2MSS \rightarrow 4MSS \dots



TCP Slow Start (more)

- When connection begins, increase rate exponentially
 - by incrementing CongWin for every ACK received
- Summary:** initial rate is slow but ramps up exponentially fast



Transport Layer 3-34

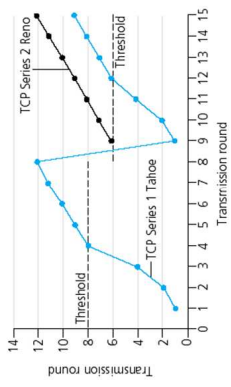
"Slow start" = starting rate is slow!
The increase speed is fast!

Collision avoidance: linear increase

Slow start to linear increase at a threshold

Q: What is threshold value?

A: Half of CongWin before timeout.



Implementation:

- At loss event, Threshold is set to 1/2 of CongWin just before loss event

Transport Layer 3-35

Fast re-tx and fast recovery

- Fast re-tx** and **fast recovery** after 3 dup ACKs (pkt loss):
 - CongWin is cut to half
 - window then grows linearly
- But after Timeout event:
 - CongWin instead set to 1 MSS (to initial start):
 - Slow start: window then grows exponentially
 - Congestions avoidance: Linear growth to a threshold, then grows linearly

Philosophy:

- 3 dup ACKs indicates network capable of delivering some segments
- timeout indicates a "more alarming" congestion scenario

Transport Layer 3-36

TCP sender congestion control

State	Event	TCP Sender Action	Commentary
Slow Start (SS)	ACK receipt for previously unacked data	$\text{CongWin} = \text{CongWin} + \text{MSS}$, If $(\text{CongWin} > \text{Threshold})$ set state to "Congestion Avoidance"	Resulting in a doubling of CongWin every RTT
Congestion Avoidance (CA)	ACK receipt for previously unacked data	$\text{CongWin} = \text{CongWin} + \text{MSS} * (\text{MSS}/\text{CongWin})$	Additive Increase. Adding in increments of 1 MSS every RTT
SS or CA	Loss event detected by triple duplicate ACK	Threshold = $\text{CongWin}/2$, $\text{CongWin} = \text{Threshold}$, Set state to "Slow Start"	Fast recovery. implementing multiplicative decrease. CongWin will not drop below 1 MSS.
SS or CA	Timeout	Threshold = $\text{CongWin}/2$, $\text{CongWin} = 1 \text{ MSS}$, Set state to "Slow Start"	Enter slow start
SS or CA	Duplicate ACK	Increment duplicate ACK count for segment being acked	CongWin and Threshold not changed

Transport Layer 3-38

TCP Futures: TCP over "long, fat pipes"

- Example: 1500 byte segments, 100ms RTT, want 10 Gbps throughput
- Requires window size $W = 83,333$ in-flight segments
- Throughput in terms of loss rate:

$$\frac{1.22 \cdot \text{MSS}}{\text{RTT} \sqrt{L}}$$

- $L = 2 \cdot 10^{-10}$ difficult to achieve
- New versions of TCP for high-speed

Transport Layer 3-40

TCP average throughput

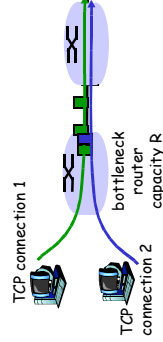
- What's the average throughput of TCP as a function of window size and RTT?
 - Ignore slow start
- Let W be the window size when loss occurs.
 - When window is W , throughput is W/RTT
 - Just after loss, window drops to $W/2$, throughput to $W/2\text{RTT}$.
- Average throughput: $.75 W/\text{RTT}$

Transport Layer 3-39

TCP Fairness

Fairness goal: if K TCP sessions share same bottleneck link of bandwidth R , each should have average rate of R/K

Absolute fairness

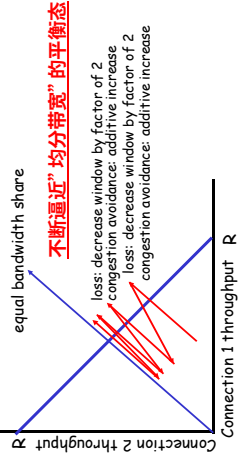


Transport Layer 3-41

Why is TCP fair?

Two competing sessions:

- Additive increase gives slope of 1, as throughput increases
- multiplicative decrease decreases throughput proportionally



Transport Layer 3-42

Chapter 3: Summary

- principles behind transport layer services:
 - multiplexing, demultiplexing (socket and port #)
 - reliable data transfer
 - flow control (发送端与接收端之间协调控制)
 - congestion control(基于对网络拥挤程度猜想的控制)
- instantiation and implementation in the Internet
 - UDP
 - TCP

Next:

 - leaving the network "edge" (application, transport layers)
 - into the network "core"

Transport Layer 3-44

Fairness (more)

Fairness and UDP

- **Multimedia** apps often do not use TCP
 - do not want rate throttled by congestion control

- Instead use UDP:
 - pump audio/video at constant rate, tolerate packet loss

- Research area: TCP fairness
 - 在带宽资源有限时, 但是不能突然中断服务

Fairness and parallel TCP connections

- nothing prevents app from opening parallel connections between 2 hosts.
- **Web browsers** do this
- Example: link of rate R supporting 9 connections:
 - new app asks for 1 TCP, gets rate $R/10$
 - new app asks for 11 TCPs, gets **more than $R/2$** !

Transport Layer 3-43