# 5

# THE NETWORK LAYER

| Email | FTP | . . . |
|-------|-----|-------|
| TCP | | |
| IP | | |
| ATM | | |
| Data link | | |
| Physical | | |

Fig. 5-1. Running TCP/IP over an ATM subnet.

| Issue | Datagram subnet | VC subnet |
|---|---|---|
| Circuit setup | Not needed | Required |
| Addressing | Each packet contains the full source and destination address | Each packet contains a short VC number |
| State information | Subnet does not hold state information | Each VC requires subnet table space |
| Routing | Each packet is routed independently | Route chosen when VC is set up; all packets follow this route |
| Effect of router failures | None, except for packets lost during the crash | All VCs that passed through the failed router are terminated |
| Congestion control | Difficult | Easy if enough buffers can be allocated in advance for each VC |

Fig. 5-2. Comparison of datagram and virtual circuit subnets.

| Upper layer | Type of subnet | |
| --- | --- | --- |
| | Datagram | Virtual circuit |
| Connectionless | UDP over IP | UDP over IP over ATM |
| Connection-oriented | TCP over IP | ATM AAL1 over ATM |

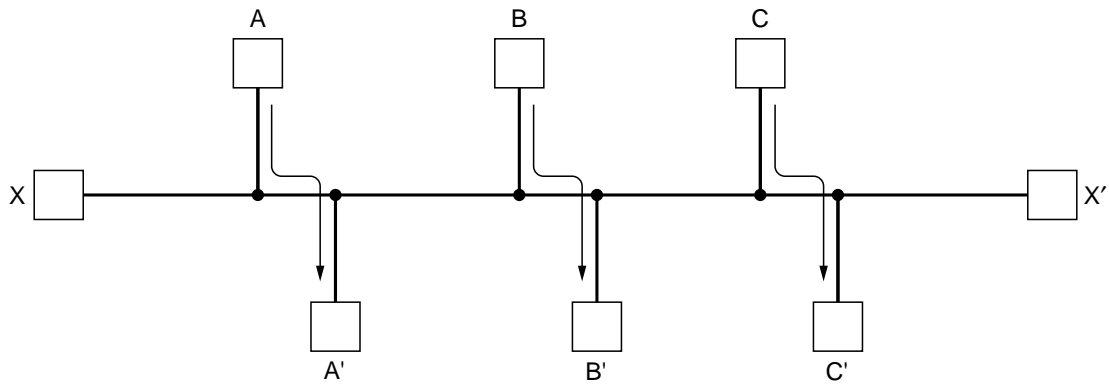Fig. 5-3. Examples of different combinations of service and subnet structure.

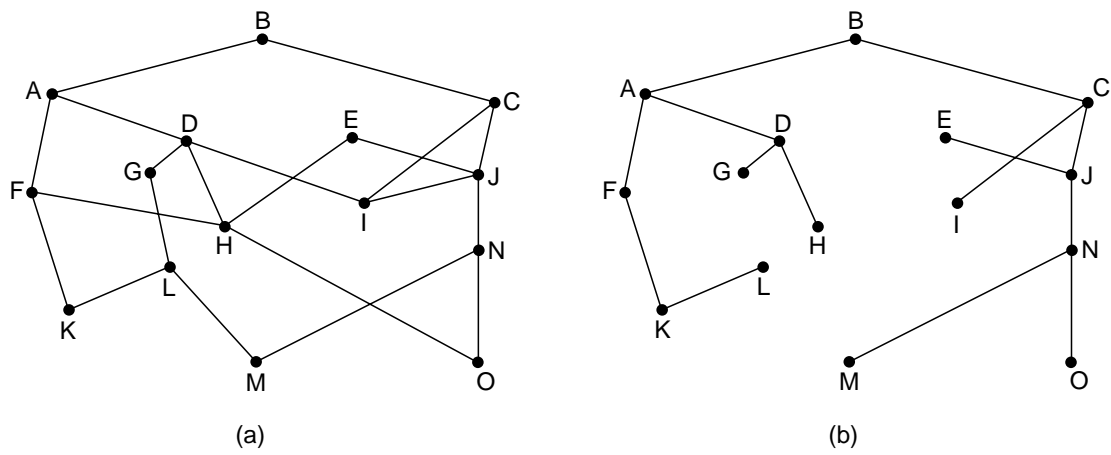Fig. 5-4. Conflict between fairness and optimality.

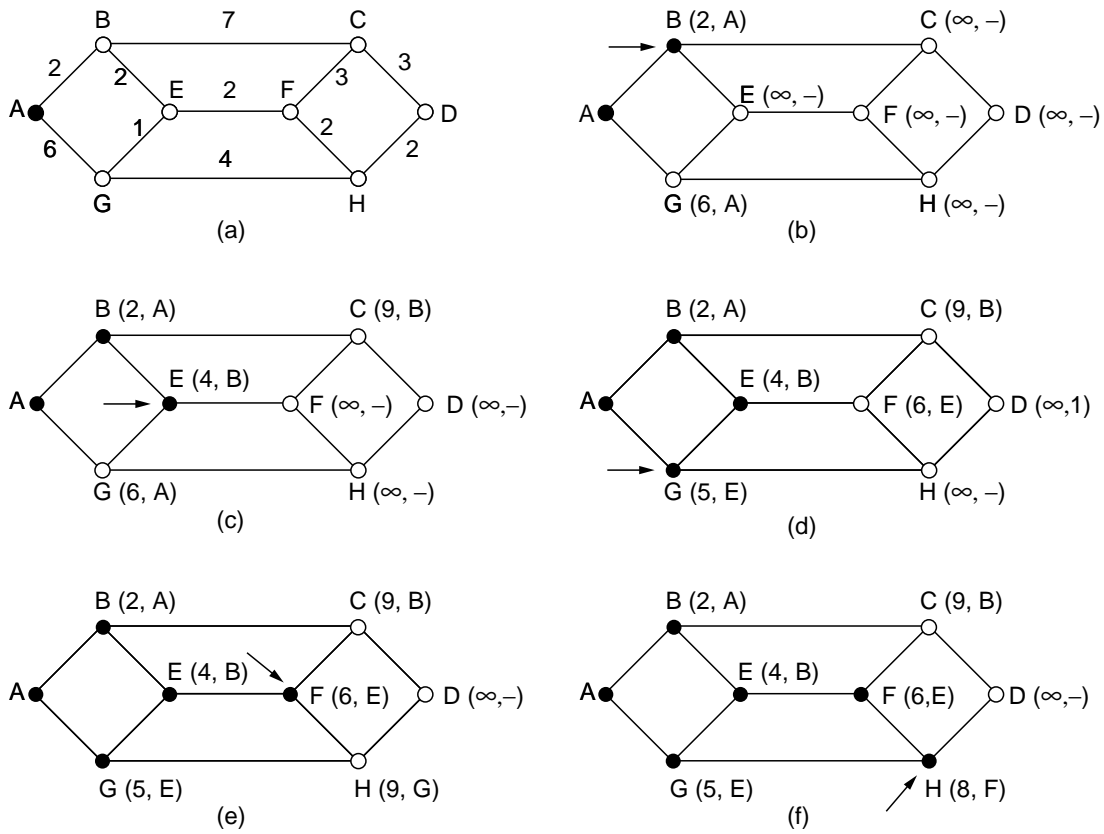Fig. 5-5. (a) A subnet. (b) A sink tree for router *B*.

Fig. 5-6. The first five steps used in computing the shortest path from *A* to *D*. The arrows indicate the working node.

```c
#define MAX_NODES 1024                      /* maximum number of nodes */
#define INFINITY 1000000000                 /* a number larger than every maximum path */
int n, dist[MAX_NODES][MAX_NODES];          /* dist[i][j] is the distance from i to j */

void shortest_path(int s, int t, int path[])
{ struct state {                            /* the path being worked on */
      int predecessor;                      /* previous node */
      int length;                           /* length from source to this node */
      enum {permanent, tentative} label;    /* label state */
 } state[MAX_NODES];

 int i, k, min;
 struct state *
             p;
 for (p = &state[0]; p < &state[n]; p++) {  /* initialize state */
     p->predecessor = −1;
     p->length = INFINITY;
     p->label = tentative;
 }
 state[t].length = 0;  state[t].label = permanent;
 k = t;                                     /* k is the initial working node */
 do {                                       /* Is there a better path from k? */
     for (i = 0; i < n; i++)                /* this graph has n nodes */
         if (dist[k][i] != 0 && state[i].label == tentative) {
             if (state[k].length + dist[k][i] < state[i].length) {
                 state[i].predecessor = k;
                 state[i].length = state[k].length + dist[k][i];
             }
         }

     /* Find the tentatively labeled node with the smallest label. */
     k = 0; min = INFINITY;
     for (i = 0; i < n; i++)
         if (state[i].label == tentative && state[i].length < min) {
             min = state[i].length;
             k = i;
         }
     state[k].label = permanent;
 } while (k != s);

 /* Copy the path into the output array. */
 i = 0;  k = s;
 do {path[i++] = k; k = state[k].predecessor; } while (k >= 0);
}
```
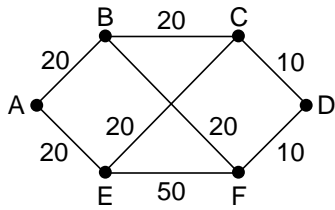
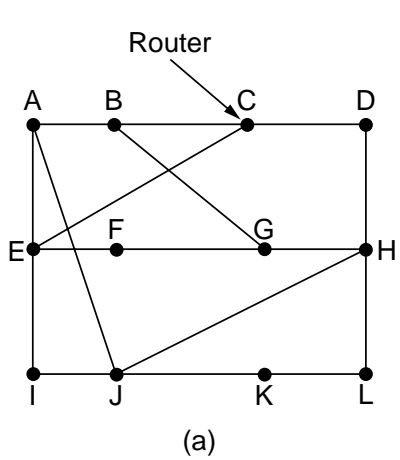Fig. 5-7. Dijkstra's algorithm to compute the shortest path.

Destination

|   | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| A |  | 9 AB | 4 ABC | 1 ABFD | 7 AE | 4 AEF |
| B | 9 BA |  | 8 BC | 3 BFD | 2 BFE | 4 BF |
| C | 4 CBA | 8 CB |  | 3 CD | 3 CE | 2 CEF |
| D | 1 DFBA | 3 DFB | 3 DC |  | 3 DCE | 4 DF |
| E | 7 EA | 2 EFB | 3 EC | 3 ECD |  | 5 EF |
| F | 4 FEA | 4 FB | 2 FEC | 4 FD | 5 FE |  |

Source

(a)

(b)

Fig. 5-8. (a) A subnet with line capacities shown in kbps. (b) The traffic in packets/sec and the routing matrix.

| i | Line | $\lambda_i$ (pkts/sec) | $C_i$ (kbps) | $\mu C_i$ (pkts/sec) | $T_i$ (msec) | Weight |
|---|------|------------------------|--------------|----------------------|--------------|--------|
| 1 | AB | 14 | 20 | 25 | 91 | 0.171 |
| 2 | BC | 12 | 20 | 25 | 77 | 0.146 |
| 3 | CD | 6 | 10 | 12.5 | 154 | 0.073 |
| 4 | AE | 11 | 20 | 25 | 71 | 0.134 |
| 5 | EF | 13 | 50 | 62.5 | 20 | 0.159 |
| 6 | FD | 8 | 10 | 12.5 | 222 | 0.098 |
| 7 | BF | 10 | 20 | 25 | 67 | 0.122 |
| 8 | EC | 8 | 20 | 25 | 59 | 0.098 |

Fig. 5-9. Analysis of the subnet of Fig. 5-0 using a mean packet size of 800 bits. The reverse traffic (*BA*, *CB*, etc.) is the same as the forward traffic.

Router

A   B   C   D
E   F   G   H
I   J   K   L

(a)

| To | A | | I | | H | | K | | ↓ | Line |
|---|---|---|---|---|---|---|---|---|---|---|
| A | 0 | | 24 | | 20 | | 21 | | 8 | A |
| B | 12 | | 36 | | 31 | | 28 | | 20 | A |
| C | 25 | | 18 | | 19 | | 36 | | 28 | I |
| D | 40 | | 27 | | 8 | | 24 | | 20 | H |
| E | 14 | | 7 | | 30 | | 22 | | 17 | I |
| F | 23 | | 20 | | 19 | | 40 | | 30 | I |
| G | 18 | | 31 | | 6 | | 31 | | 18 | H |
| H | 17 | | 20 | | 0 | | 19 | | 12 | H |
| I | 21 | | 0 | | 14 | | 22 | | 10 | I |
| J | 9 | | 11 | | 7 | | 10 | | 0 | – |
| K | 24 | | 22 | | 22 | | 0 | | 6 | K |
| L | 29 | | 33 | | 9 | | 9 | | 15 | K |

JA delay is 8   JI delay is 10   JH delay is 12   JK delay is 6

New routing table for J

Vectors received from J's four neighbors

(b)

Fig. 5-10. (a) A subnet. (b) Input from *A*, *I*, *H*, *K*, and the new routing table for *J*.

(a)

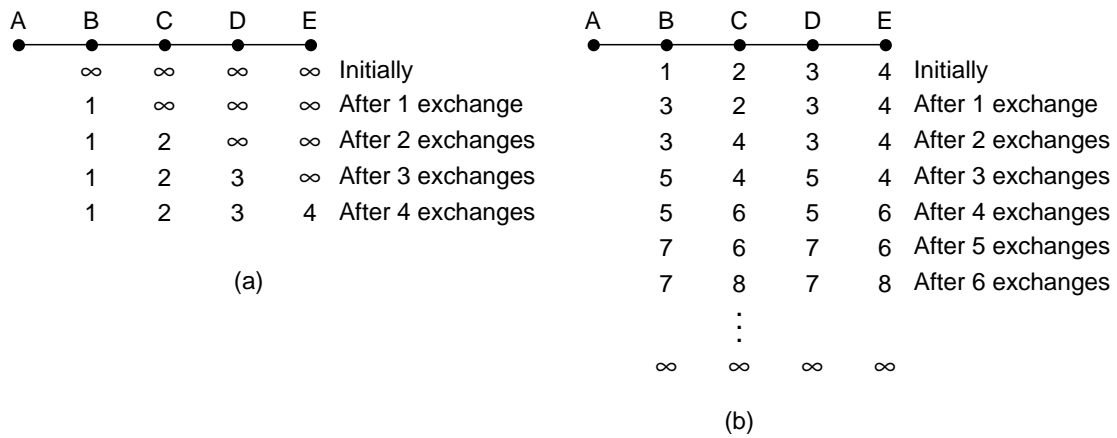| A | B | C | D | E | |
|---|---|---|---|---|---|
| • | ∞ | ∞ | ∞ | ∞ | Initially |
| | 1 | ∞ | ∞ | ∞ | After 1 exchange |
| | 1 | 2 | ∞ | ∞ | After 2 exchanges |
| | 1 | 2 | 3 | ∞ | After 3 exchanges |
| | 1 | 2 | 3 | 4 | After 4 exchanges |

(b)

| A | B | C | D | E | |
|---|---|---|---|---|---|
| • | 1 | 2 | 3 | 4 | Initially |
| | 3 | 2 | 3 | 4 | After 1 exchange |
| | 3 | 4 | 3 | 4 | After 2 exchanges |
| | 5 | 4 | 5 | 4 | After 3 exchanges |
| | 5 | 6 | 5 | 6 | After 4 exchanges |
| | 7 | 6 | 7 | 6 | After 5 exchanges |
| | 7 | 8 | 7 | 8 | After 6 exchanges |
| | ⋮ | | | | |
| | ∞ | ∞ | ∞ | ∞ | |

Fig. 5-11. The count-to-infinity problem.

Fig. 5-12. An example where split horizon fails.

Fig. 5-13. (a) Nine routers and a LAN. (b) A graph model of (a).

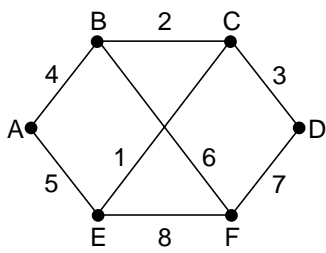Fig. 5-14. A subnet in which the East and West parts are connected by two lines.

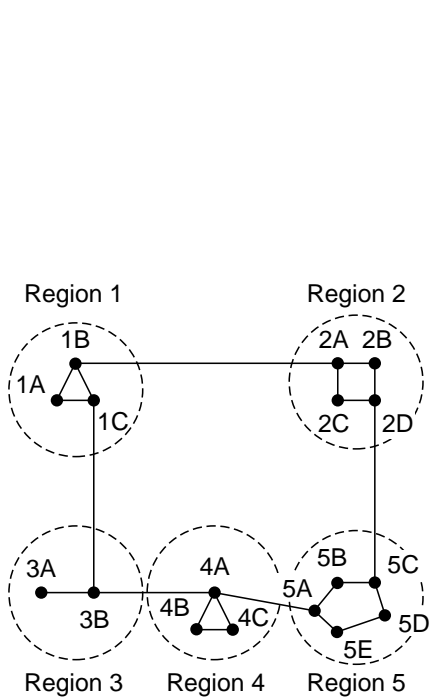Fig. 5-15. (a) A subnet. (b) The link state packets for this subnet.

| Source | Seq. | Age | Send flags | | | ACK flags | | | Data |
|---|---|---|---|---|---|---|---|---|---|
| | | | A | C | F | A | C | F | |
| A | 21 | 60 | 0 | 1 | 1 | 1 | 0 | 0 | |
| F | 21 | 60 | 1 | 1 | 0 | 0 | 0 | 1 | |
| E | 21 | 59 | 0 | 1 | 0 | 1 | 0 | 1 | |
| C | 20 | 60 | 1 | 0 | 1 | 0 | 1 | 0 | |
| D | 21 | 59 | 1 | 0 | 0 | 0 | 1 | 1 | |

Fig. 5-16. The packet buffer for router *B* in Fig. 5-15.

Full table for 1A

| 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 table for 1A

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

(a)

(b)

(c)

Fig. 5-17. Hierarchical routing.

Fig. 5-18. A WAN to which LANs, MANs, and wireless cells are attached.
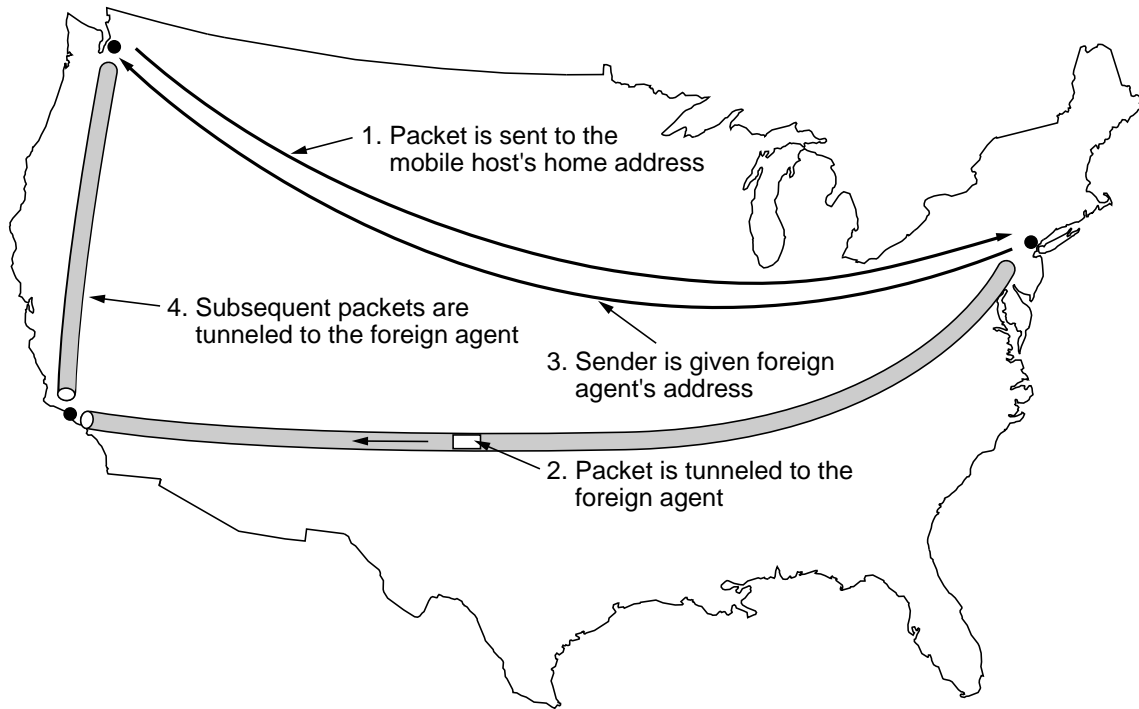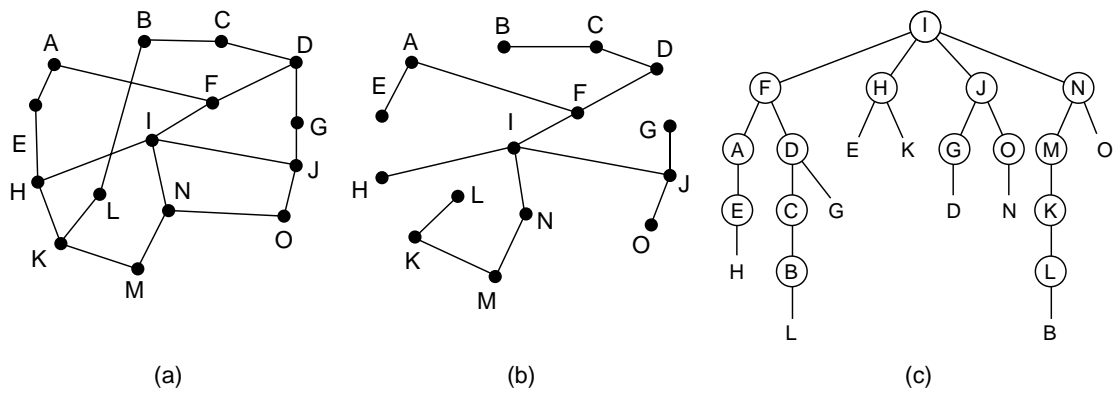
Fig. 5-19. Packet routing for mobile users.

Fig. 5-20. Reverse path forwarding. (a) A subnet. (b) A spanning tree. (c) The tree built by reverse path forwarding.
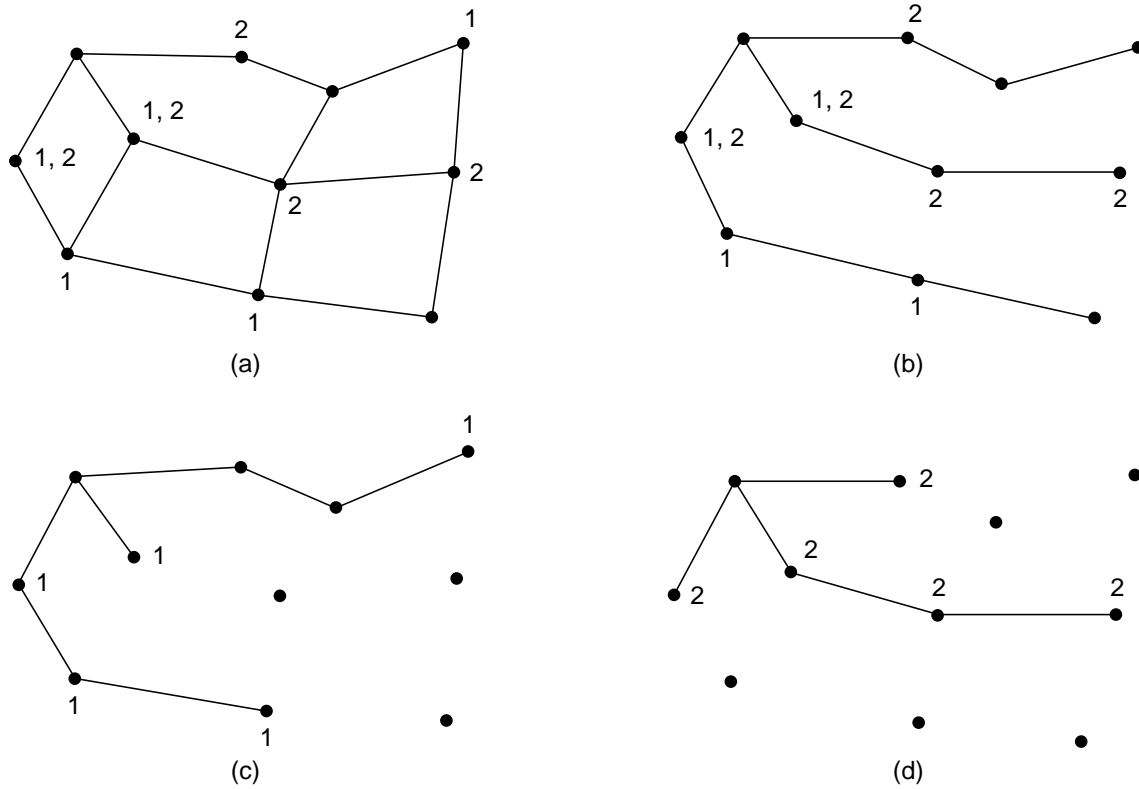
Fig. 5-21. (a) A subnet. (b) A spanning tree for the leftmost router. (c) A multicast tree for group 1. (d) A multicast tree for group 2.
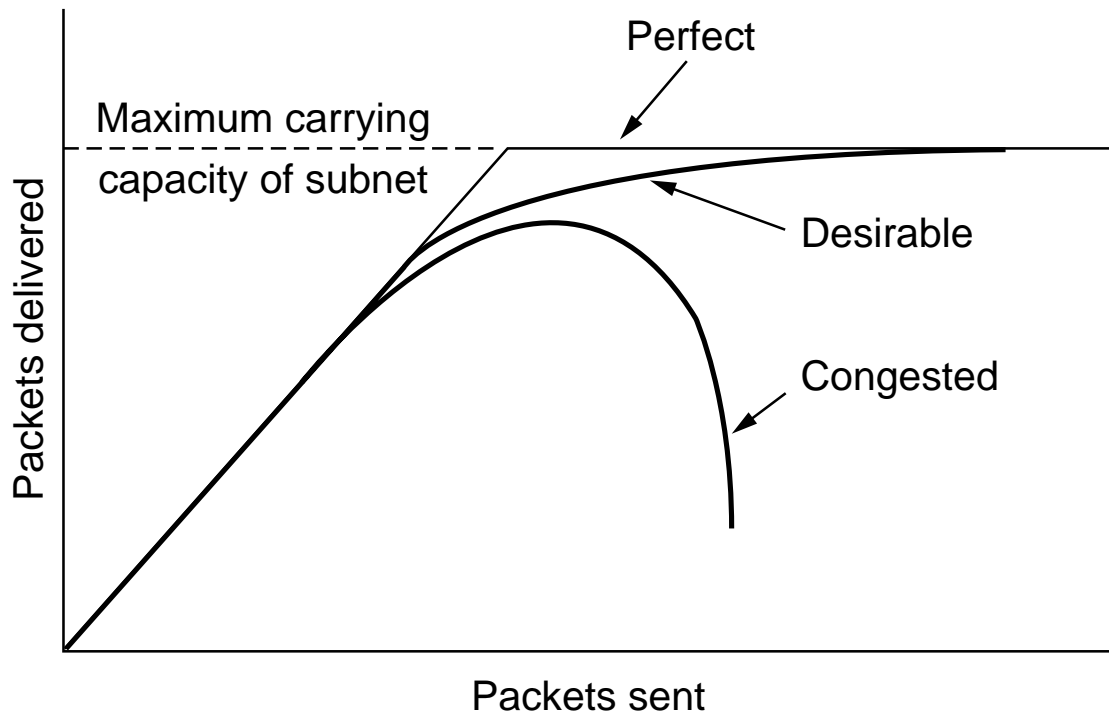
Fig. 5-22. When too much traffic is offered, congestion sets in and performance degrades sharply.

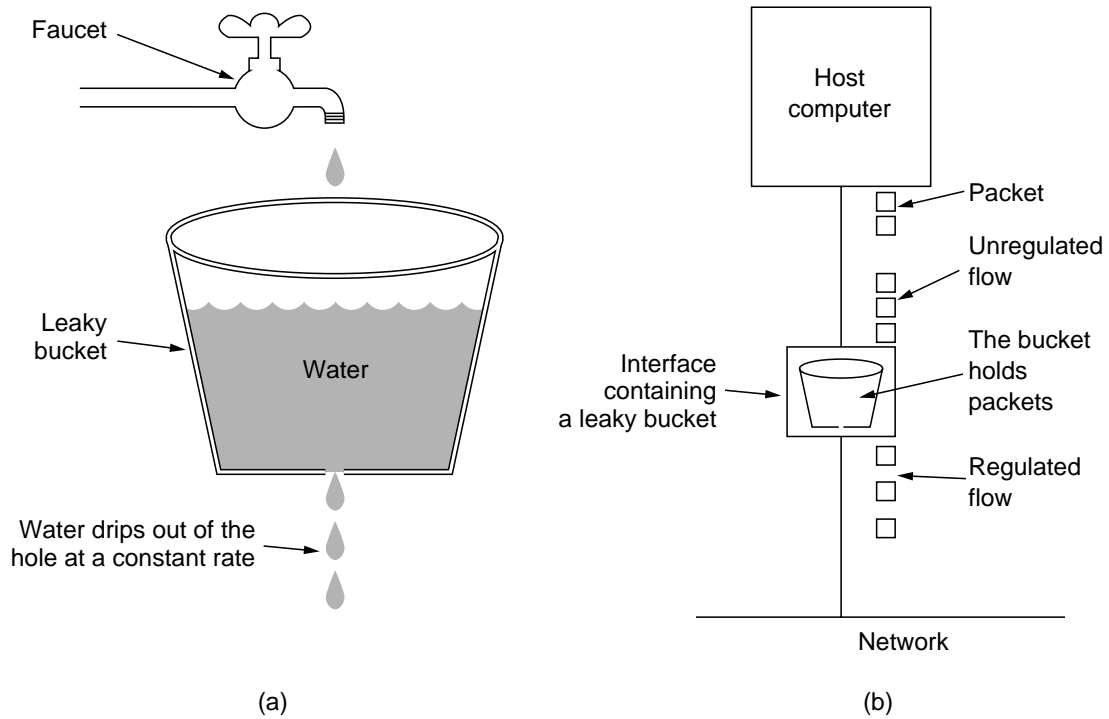| Layer | Policies |
|-------|----------|
| Transport | • Retransmission policy<br>• Out-of-order caching policy<br>• Acknowledgement policy<br>• Flow control policy<br>• Timeout determination |
| Network | • Virtual circuits versus datagram inside the subnet<br>• Packet queueing and service policy<br>• Packet discard policy<br>• Routing algorithm<br>• Packet lifetime management |
| Data link | • Retransmission policy<br>• Out-of-order caching policy<br>• Acknowledgement policy<br>• Flow control policy |

Fig. 5-23. Policies that affect congestion.

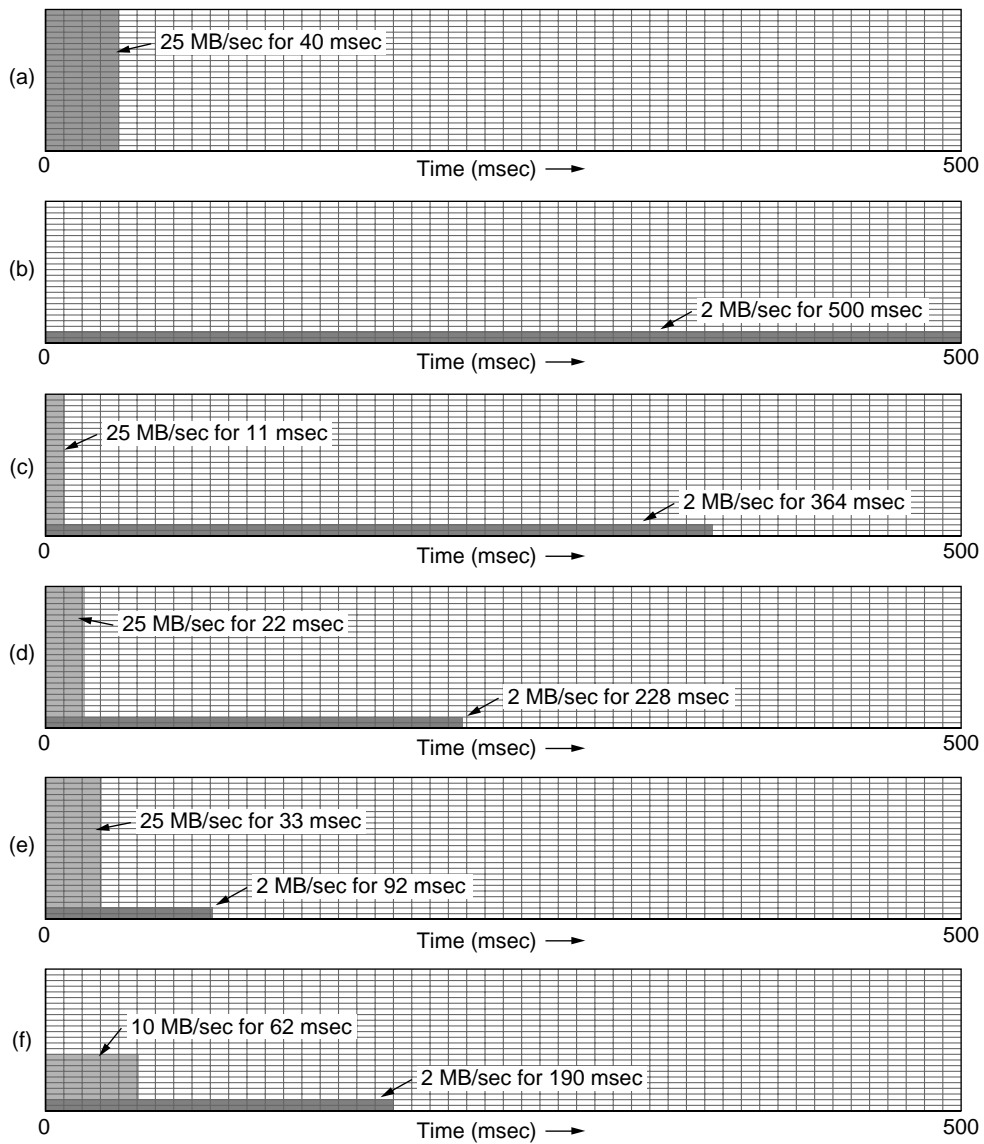Fig. 5-24. (a) A leaky bucket with water. (b) A leaky bucket with packets.

Fig. 5-25. (a) Input to a leaky bucket. (b) Output from a leaky bucket. (c) - (e) Output from a token bucket with capacities of 250KB, 500KB, and 750KB. (f) Output from a 500KB token bucket feeding a 10 MB/sec leaky bucket.
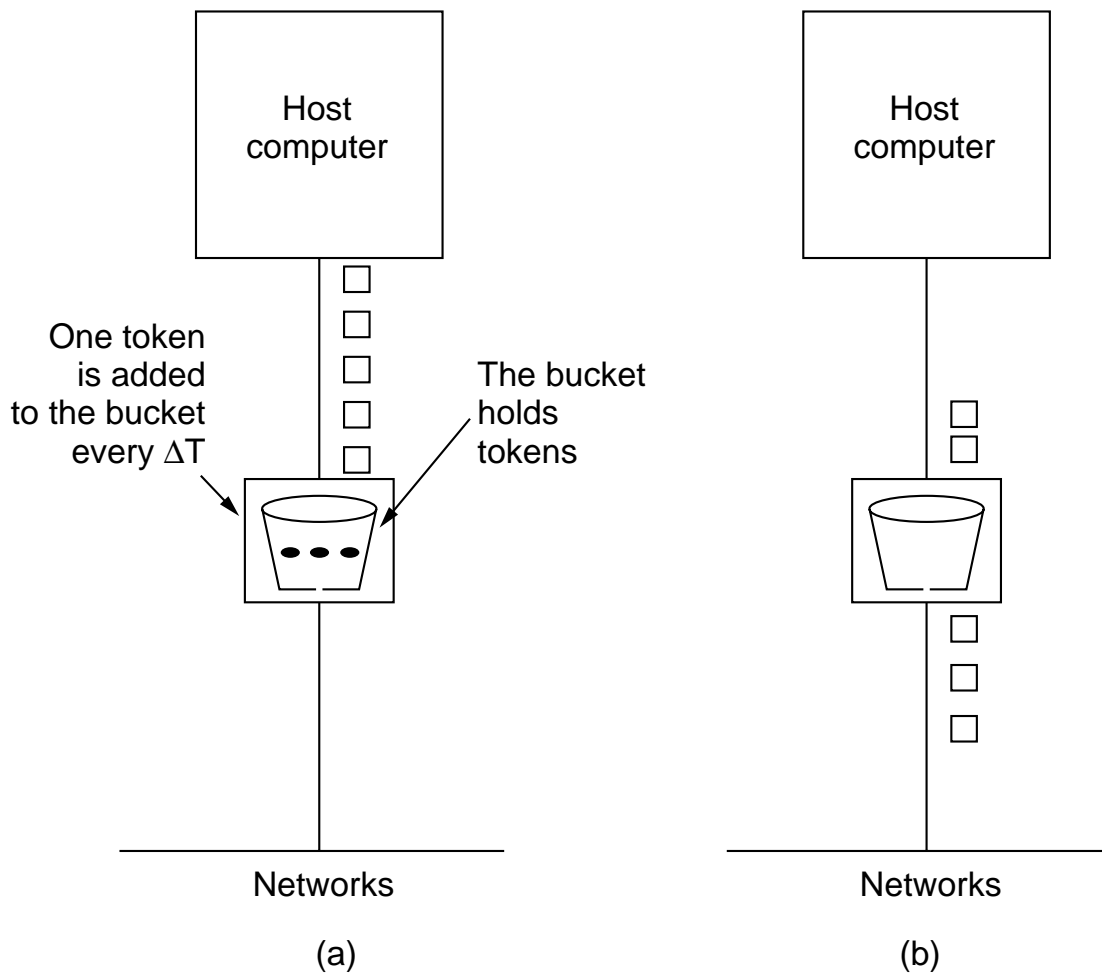
Fig. 5-26. The token bucket algorithm. (a) Before. (b) After.

| Characteristics of the Input | Service Desired |
| --- | --- |
| Maximum packet size (bytes) | Loss sensitivity (bytes) |
| Token bucket rate (bytes/sec) | Loss interval ($\mu$sec) |
| Token bucket size (bytes) | Burst loss sensitivity (packets) |
| Maximum transmission rate (bytes/sec) | Minimum delay noticed ($\mu$sec) |
| | Maximum delay variation ($\mu$sec) |
| | Quality of guarantee |

Fig. 5-27. An example flow specification.
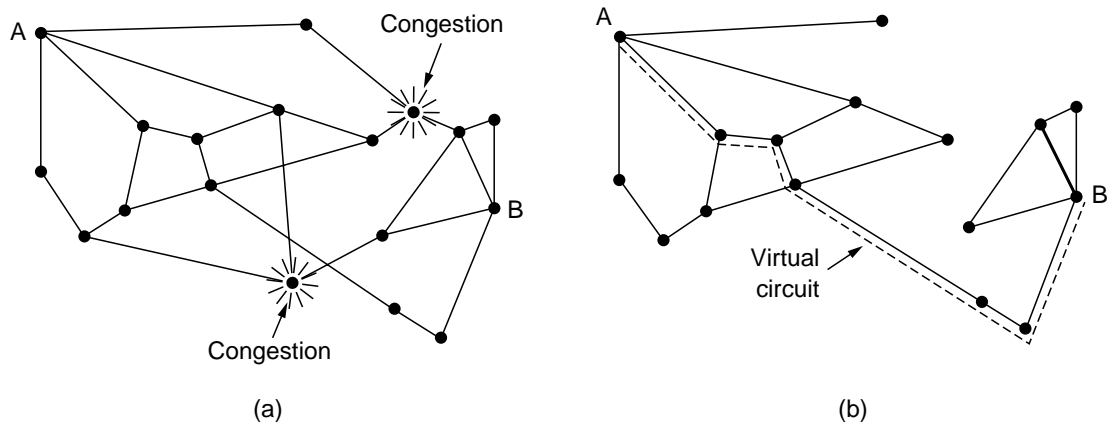
Fig. 5-28. (a) A congested subnet. (b) A redrawn subnet that eliminates the congestion and a virtual circuit from *A* to *B*.

| | | | | | |
|---|---|---|---|---|---|
| A | 1 | 6 | 11 | 15 | 19 | 20 |
| B | 2 | 7 | 12 | 16 |
| C | 3 | 8 |
| D | 4 | 9 | 13 | 17 |
| E | 5 | 10 | 14 | 18 |

O

(a)

| Packet | Finishing time |
|---|---|
| C | 8 |
| B | 16 |
| D | 17 |
| E | 18 |
| A | 20 |

(b)

Fig. 5-29. (a) A router with five packets queued for line *O*. (b) Finishing times for the five packets.

Fig. 5-30. (a) A choke packet that affects only the source. (b) A choke packet that affects each hop it passes through.

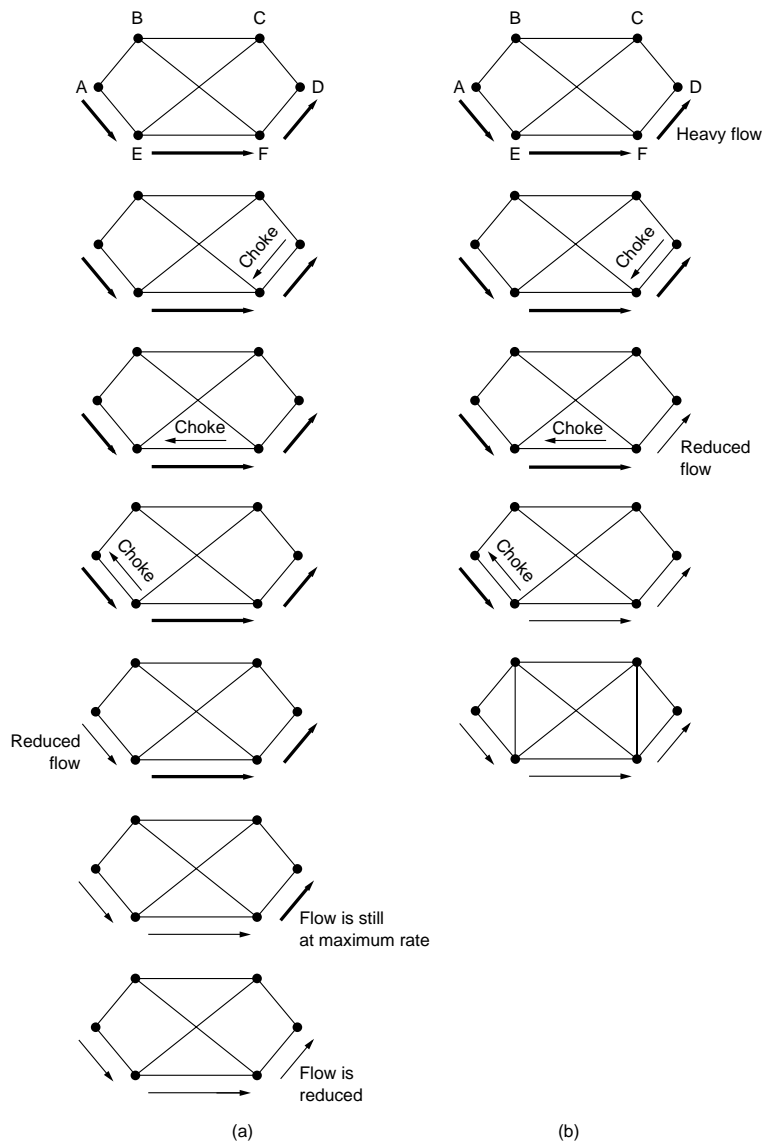Fig. 5-31. (a) A network. (b) The multicast spanning tree for host 1. (c) The multicast spanning tree for host 2.

Fig. 5-32. (a) Host 3 requests a channel to host 1. (b) Host 3 then requests a second channel, to host 2. (c) Host 5 requests a channel to host 1.

Fig. 5-33. Network interconnection.

Network 1
packets here

Network 2
packets here

Network 1   G   Network 2

(a)

Full
gateway

Network 2
packets here

G   Network 2

Network 1

Network 1
packets here

(b)

Half-gateway   Neutral packets here

Network 1   Network 2

(c)

Fig. 5-34. (a) A full gateway between two WANs. (b) A full gate-
way between a LAN and a WAN. (c) Two half-gateways.

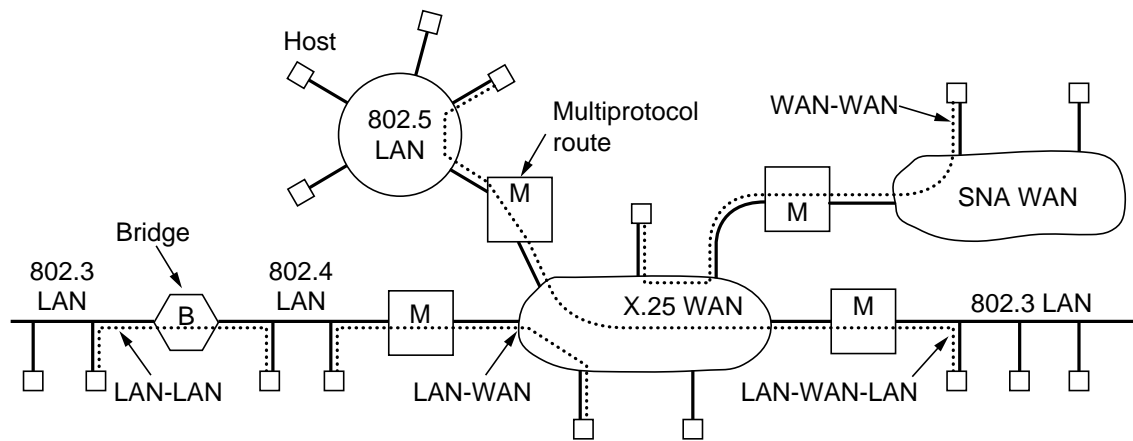| Item | Some Possibilities |
|---|---|
| Service offered | Connection-oriented versus connectionless |
| Protocols | IP, IPX, CLNP, AppleTalk, DECnet, etc. |
| Addressing | Flat (802) versus hierarchical (IP) |
| Multicasting | Present or absent (also broadcasting) |
| Packet size | Every network has its own maximum |
| Quality of service | May be present or absent; many different kinds |
| Error handling | Reliable, ordered, and unordered delivery |
| Flow control | Sliding window, rate control, other, or none |
| Congestion control | Leaky bucket, choke packets, etc. |
| Security | Privacy rules, encryption, etc. |
| Parameters | Different timeouts, flow specifications, etc. |
| Accounting | By connect time, by packet, by byte, or not at all |

Fig. 5-35. Some of the many ways networks can differ.

Fig. 5-36. Internetworking using concatenated virtual circuits.
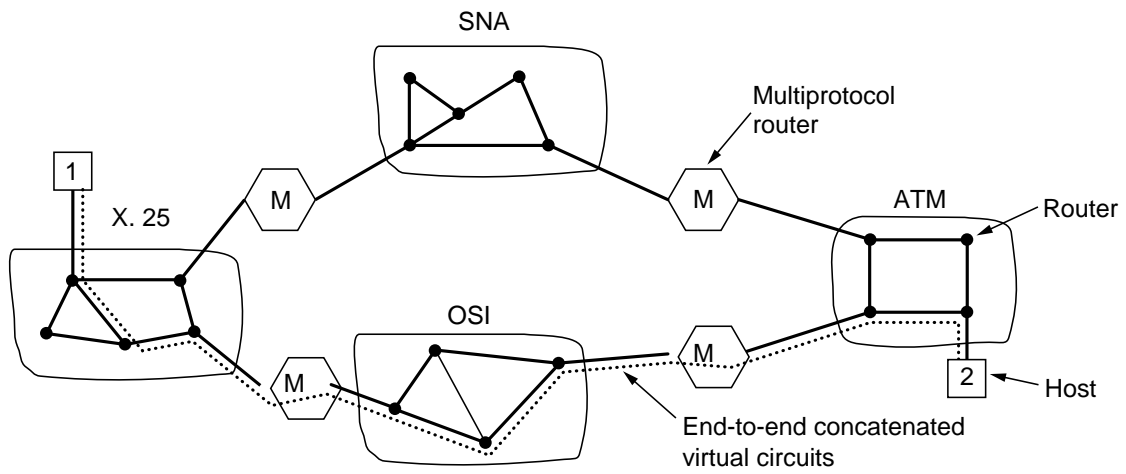
Fig. 5-37. A connectionless internet.

Fig. 5-38. Tunneling a packet from Paris to London.

Fig. 5-39. Tunneling a car from France to England.

Fig. 5-40. (a) An internetwork. (b) A graph of the internetwork.

Fig. 5-41. (a) Transparent fragmentation. (b) Nontransparent fragmentation.

Number of the first elementary fragment in this packet

Packet number | End of packet bit | 1 byte

| 27 | 0 | 1 | A | B | C | D | E | F | G | H | I | J |

Header

(a)

| 27 | 0 | 0 | A | B | C | D | E | F | G | H |   | 27 | 8 | 1 | I | J |

Header                                              Header

(b)

| 27 | 0 | 0 | A | B | C | D | E |   | 27 | 5 | 0 | F | G | H |   | 27 | 8 | 1 | I | J |

Header                          Header                       Header

(c)

Fig. 5-42. Fragmentation when the elementary data size is 1 byte. (a) Original packet, containing 10 data bytes. (b) Fragments after passing through a network with maximum packet size of 8 bytes. (c) Fragments after passing through a size 5 gateway.

Packet
filtering
router

Application
gateway

Packet
filtering
router

Connections
to outside
networks

Backbone

Corporate
network

Security
perimeter

Inside
LAN

Outside
LAN

Firewall

Fig. 5-43. A firewall consisting of two packet filters and an application gateway.

Fig. 5-44. The Internet is an interconnected collection of many networks.

Fig. 5-45. The IP (Internet Protocol) header.

| Option | Description |
| --- | --- |
| Security | Specifies how secret the datagram is |
| Strict source routing | Gives the complete path to be followed |
| Loose source routing | Gives a list of routers not to be missed |
| Record route | Makes each router append its IP address |
| Timestamp | Makes each router append its address and timestamp |

Fig. 5-46. IP options.

Fig. 5-47. IP address formats.

| | |
|---|---|
| 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 | This host |
| 0 0    . . .    0 0     Host | A host on this network |
| 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 | Broadcast on the local network |
| Network    1 1 1 1    . . .    1 1 1 1 | Broadcast on a distant network |
| 127     (Anything) | Loopback |

Fig. 5-48. Special IP addresses.

32 Bits

| 10 | Network | Subnet | Host |

Subnet mask

1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0

Fig. 5-49. One of the ways to subnet a class B network.

| Message type | Description |
| --- | --- |
| Destination unreachable | Packet could not be delivered |
| Time exceeded | Time to live field hit 0 |
| Parameter problem | Invalid header field |
| Source quench | Choke packet |
| Redirect | Teach a router about geography |
| Echo request | Ask a machine if it is alive |
| Echo reply | Yes, I am alive |
| Timestamp request | Same as Echo request, but with timestamp |
| Timestamp reply | Same as Echo reply, but with timestamp |

Fig. 5-50. The principal ICMP message types.

Router has
2 IP addresses
192.31.60.4
192.31.65.1

Router has
2 IP addresses
192.31.60.7
192.31.63.3

To WAN

192.31.65.7      192.31.65.5

192.31.63.8

F2

| 1 |      | 2 |

F1        F3

| 3 |      | 4 |

E1        E2        E3        E4        E5        E6

Ethernet
addresses

CS Ethernet
192.31.65.0

Campus
FDDI ring
192.31.60.0

EE Ethernet
192.31.63.0

Fig. 5-51. Three interconnected class C networks: two Ethernets
and an FDDI ring.

WAN 1

A  B  C  D  E  WAN 2  F  I  J

LAN 1

H  LAN 2

G

WAN 3

(a)

W1

10  12

A  B  C  4  D  6  8  W2  F

4  6  16

2  3  13  J

2  H  3  4

1  L2  2

17  12

G  W3

L1

(b)

Fig. 5-52. (a) An autonomous system. (b) A graph representation of (a).

Fig. 5-53. The relation between ASes, backbones, and areas in OSPF.

| Message type | Description |
| --- | --- |
| Hello | Used to discover who the neighbors are |
| Link state update | Provides the sender's costs to its neighbors |
| Link state ack | Acknowledges link state update |
| Database description | Announces which updates the sender has |
| Link state request | Requests information from the partner |

Fig. 5-54. The five types of OSPF messages.

Information F receives
from its neighbors about D

From B: "I use BCD"
From G: "I use GCD"
From I:  "I use IFGCD"
From E: "I use EFGCD"

(a)                                                    (b)

Fig. 5-55. (a) A set of BGP routers. (b) Information sent to *F*.

Fig. 5-56. The IPv6 fixed header (required).

| Prefix (binary) | Usage | Fraction |
| --- | --- | --- |
| 0000 0000 | Reserved (including IPv4) | 1/256 |
| 0000 0001 | Unassigned | 1/256 |
| 0000 001 | OSI NSAP addresses | 1/128 |
| 0000 010 | Novell NetWare IPX addresses | 1/128 |
| 0000 011 | Unassigned | 1/128 |
| 0000 1 | Unassigned | 1/32 |
| 0001 | Unassigned | 1/16 |
| 001 | Unassigned | 1/8 |
| 010 | Provider-based addresses | 1/8 |
| 011 | Unassigned | 1/8 |
| 100 | Geographic-based addresses | 1/8 |
| 101 | Unassigned | 1/8 |
| 110 | Unassigned | 1/8 |
| 1110 | Unassigned | 1/16 |
| 1111 0 | Unassigned | 1/32 |
| 1111 10 | Unassigned | 1/64 |
| 1111 110 | Unassigned | 1/128 |
| 1111 1110 0 | Unassigned | 1/512 |
| 1111 1110 10 | Link local use addresses | 1/1024 |
| 1111 1110 11 | Site local use addresses | 1/1024 |
| 1111 1111 | Multicast | 1/256 |

Fig. 5-57. IPv6 addresses

| Extension header | Description |
| --- | --- |
| Hop-by-hop options | Miscellaneous information for routers |
| Routing | Full or partial route to follow |
| Fragmentation | Management of datagram fragments |
| Authentication | Verification of the sender's identity |
| Encrypted security payload | Information about the encrypted contents |
| Destination options | Additional information for the destination |

Fig. 5-58. IPv6 extension headers.

| Next header | 0 | 194 | 0 |
|:---:|:---:|:---:|:---:|
| Jumbo payload length | | | |

Fig. 5-59. The hop-by-hop extension header for large datagrams (jumbograms).

| Next header | 0 | Number of addresses | Next address |
|---|---|---|---|
| | Bit map | | |

| 1 – 24 Adresses |
|---|

Fig. 5-60. The extension header for routing.

Fig. 5-61. A transmission path can hold multiple virtual paths, each of which can hold multiple virtual circuits.

|  |  |  |  |  |  |  |
|------|------|------|------|------|------|------|
| GFC | VP I | VC I | PTI | C L P | HEC | |

(a)

|  |  |  |  |  |
|------|------|------|------|------|
| VP I | VC I | PTI | C L P | HEC |

GFC: General Flow Control         PTI: Payload Type
VPI: Virtual Path Identifier         CLP: Cell Loss Priority
VCI: Virtual Channel Identification     VCI: Header Error Check

(b)

Fig. 5-62. (a) The ATM layer header at the UNI. (b) The ATM layer header at the NNI.

| Payload type | Meaning |
|:---:|:---|
| 000 | User data cell, no congestion, cell type 0 |
| 001 | User data cell, no congestion, cell type 1 |
| 010 | User data cell, congestion experienced, cell type 0 |
| 011 | User data cell, congestion experienced, cell type 1 |
| 100 | Maintenance information between adjacent switches |
| 101 | Maintenance information between source and destination switches |
| 110 | Resource Management cell (used for ABR congestion control) |
| 111 | Reserved for future function |

Fig. 5-63. Values of the *PTI* field.

| Message | Meaning when sent by host | Meaning when sent by network |
|---|---|---|
| SETUP | Please establish a circuit | Incoming call |
| CALL PROCEEDING | I saw the incoming call | Your call request will be attempted |
| CONNECT | I accept the incoming call | Your call request was accepted |
| CONNECT ACK | Thanks for accepting | Thanks for making the call |
| RELEASE | Please terminate the call | The other side has had enough |
| RELEASE COMPLETE | Ack for RELEASE | Ack for RELEASE |

Fig. 5-64. Messages used for connection establishment and release.

Source
host                          Switch #1                    Switch #2                  Destination
                                                                                     host

Setup
Call proceeding
                              Setup
                              Call proceeding
                                                           Setup
                                                           Connect
Connect                       Connect                      Connect ack
Connect ack                   Connect ack

(a)

Release
Release complete
                              Release
                              Release complete
                                                           Release
                                                           Release complete
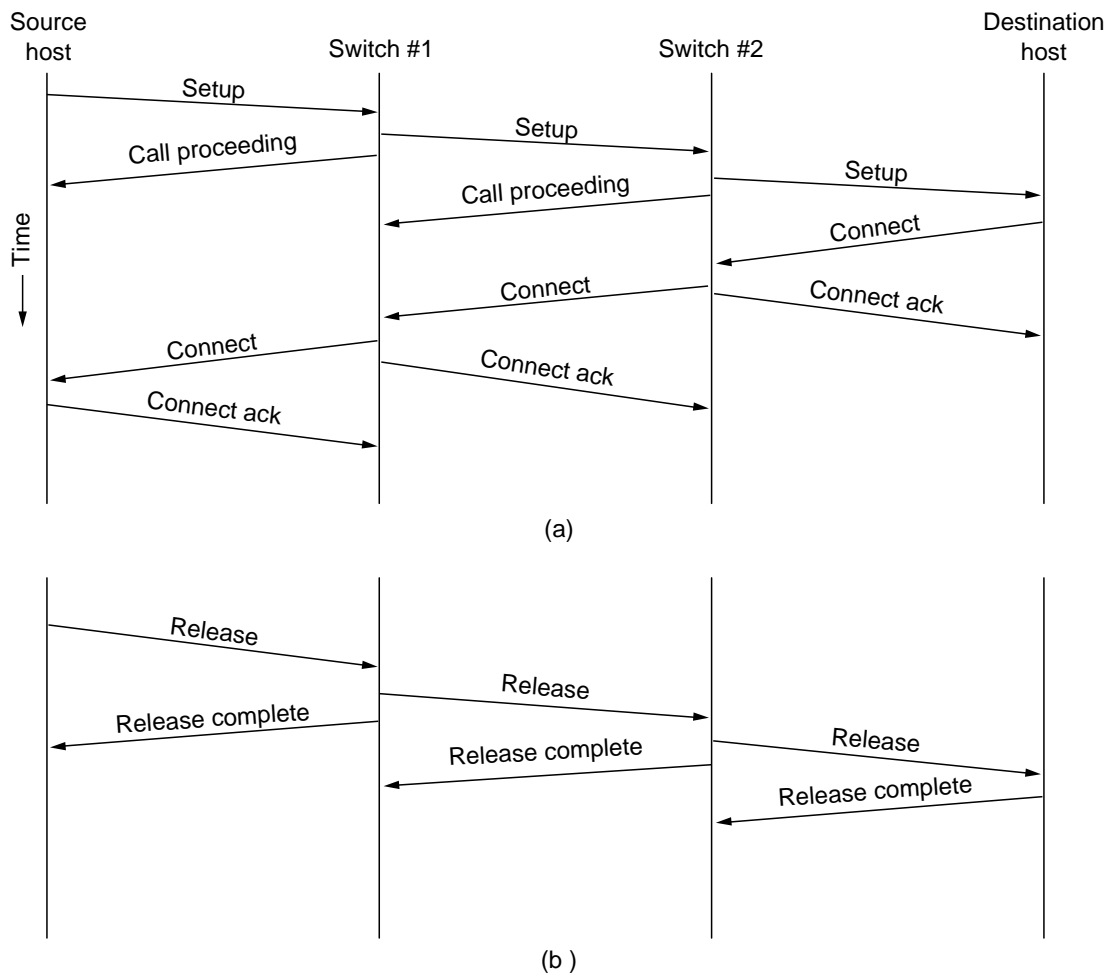
(b )

Fig. 5-65. (a) Connection setup in an ATM network. (b) Connec-
tion release.

Fig. 5-66. Rerouting a virtual path reroutes all of its virtual circuits.

| Source | Incoming line | Incoming VPI | Destina-tion | Outgoing line | Outgoing VPI | Path: |
|--------|---------------|--------------|--------------|---------------|--------------|-------|
| NY | 1 | 1 | SF | 4 | 1 | New |
| NY | 1 | 2 | Denver | 4 | 2 | New |
| LA | 3 | 1 | Minneapolis | 0 | 1 | New |
| DC | 1 | 3 | LA | 3 | 2 | New |
| NY | 1 | 1 | SF | 4 | 1 | Old |
| SF | 4 | 3 | DC | 1 | 4 | New |
| DC | 1 | 5 | SF | 4 | 4 | New |
| NY | 1 | 2 | Denver | 4 | 2 | Old |
| SF | 4 | 5 | Minneapolis | 0 | 2 | New |
| NY | 1 | 1 | SF | 4 | 1 | Old |

Fig. 5-67. Some routes through the Omaha switch of Fig. 5-0.

| Incoming VPI | VPI_table for Minn. Outgoing Line | VPI | VPI_table for DC Outgoing Line | VPI | VPI_table for Dallas Outgoing Line | VPI | VPI_table for LA Outgoing Line | VPI | VPI_table for Denver Outgoing Line | VPI |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | | | | | | | | | | |
| 1 | 3 | 1 | 4 | 1 | | | 0 | 1 | 1 | 1 |
| 2 | 4 | 5 | 4 | 2 | | | 1 | 3 | 1 | 2 |
| 3 | | | 3 | 2 | | | | | 1 | 4 |
| 4 | | | 4 | 3 | | | | | 1 | 5 |
| 5 | | | 4 | 4 | | | | | 0 | 2 |
| 6 | | | | | | | | | | |
| 7 | | | | | | | | | | |
| 8 | | | | | | | | | | |
| 4095 | | | | | | | | | | |
| | Line 0 | | Line 1 | | Line 2 | | Line 3 | | Line 4 | |

Fig. 5-68. The table entries for the routes of Fig. 5-0.

| Class | Description | Example |
|---|---|---|
| CBR | Constant bit rate | T1 circuit |
| RT-VBR | Variable bit rate: real time | Real-time videoconferencing |
| NRT-VBR | Variable bit rate: non-real time | Multimedia email |
| ABR | Available bit rate | Browsing the Web |
| UBR | Unspecified bit rate | Background file transfer |

Fig. 5-69. The ATM service categories.

| Service characteristic | CBR | RT-VBR | NRT-VBR | ABR | UBR |
|---|---|---|---|---|---|
| Bandwidth guarantee | Yes | Yes | Yes | Optional | No |
| Suitable for real-time traffic | Yes | Yes | No | No | No |
| Suitable for bursty traffic | No | No | Yes | Yes | Yes |
| Feedback about congestion | No | No | No | Yes | No |

Fig. 5-70. Characteristics of the ATM service categories.

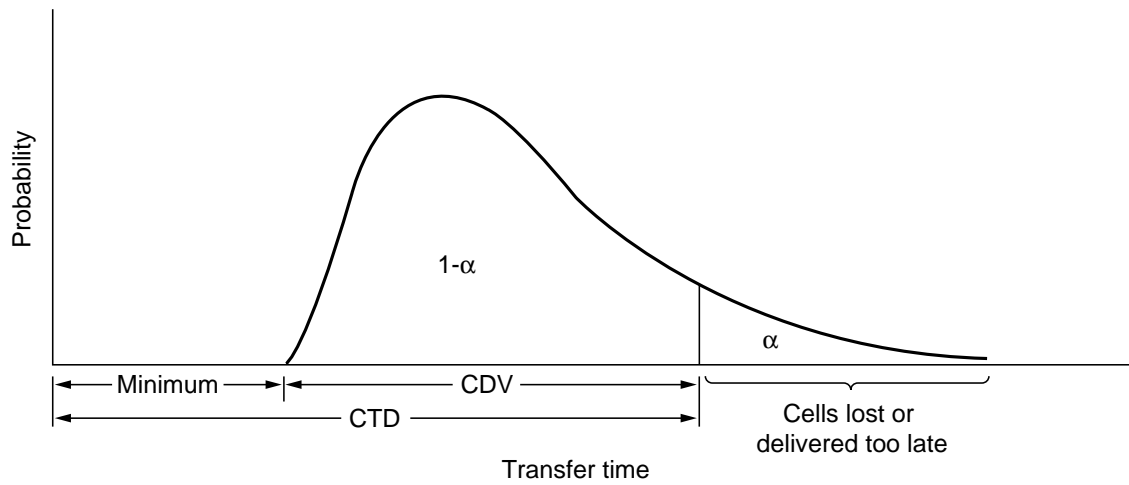| Parameter | Acronym | Meaning |
|---|---|---|
| Peak cell rate | PCR | Maximum rate at which cells will be sent |
| Sustained cell rate | SCR | The long-term average cell rate |
| Minimum cell rate | MCR | The minimum acceptable cell rate |
| Cell delay variation tolerance | CDVT | The maximum acceptable cell jitter |
| Cell loss ratio | CLR | Fraction of cells lost or delivered too late |
| Cell transfer delay | CTD | How long delivery takes (mean and maximum) |
| Cell delay variation | CDV | The variance in cell delivery times |
| Cell error rate | CER | Fraction of cells delivered without error |
| Severely-errored cell block ratio | SECBR | Fraction of blocks garbled |
| Cell misinsertion rate | CMR | Fraction of cells delivered to wrong destination |

Fig. 5-71. Some of the quality of service parameters.

Fig. 5-72. The probability density function for cell arrival times.

Time ⟶

Cell

(a)  1          2   Maximal case.
                    Cell 2 arrives T sec after Cell 1

$t_1$              $t_2$
                                        cell 3 expected
                    |◄——— T ———►|◄  at $t_2$ + T

(b)  1              2   Slow sender.
                       Cell 2 arrives > T sec after cell 1

$t_1$          $t_2$
                                        Cell 3 expected
               |◄——— T ———►|◄  at $t_2$ + T

(c)  1         2   Fast sender.
                   Cell 2 arrives up to L sec early

$t_1$      $t_2$
                                        Cell 3 expected
                   |◄——— T ———►|◄  at $t_1$ + 2T

(d)  1         2   Very fast sender.
                   Cell 2 arrives prior to $t_1$ + T − L.
                   Cell is nonconforming.

$t_1$
                   Cell 3 expected
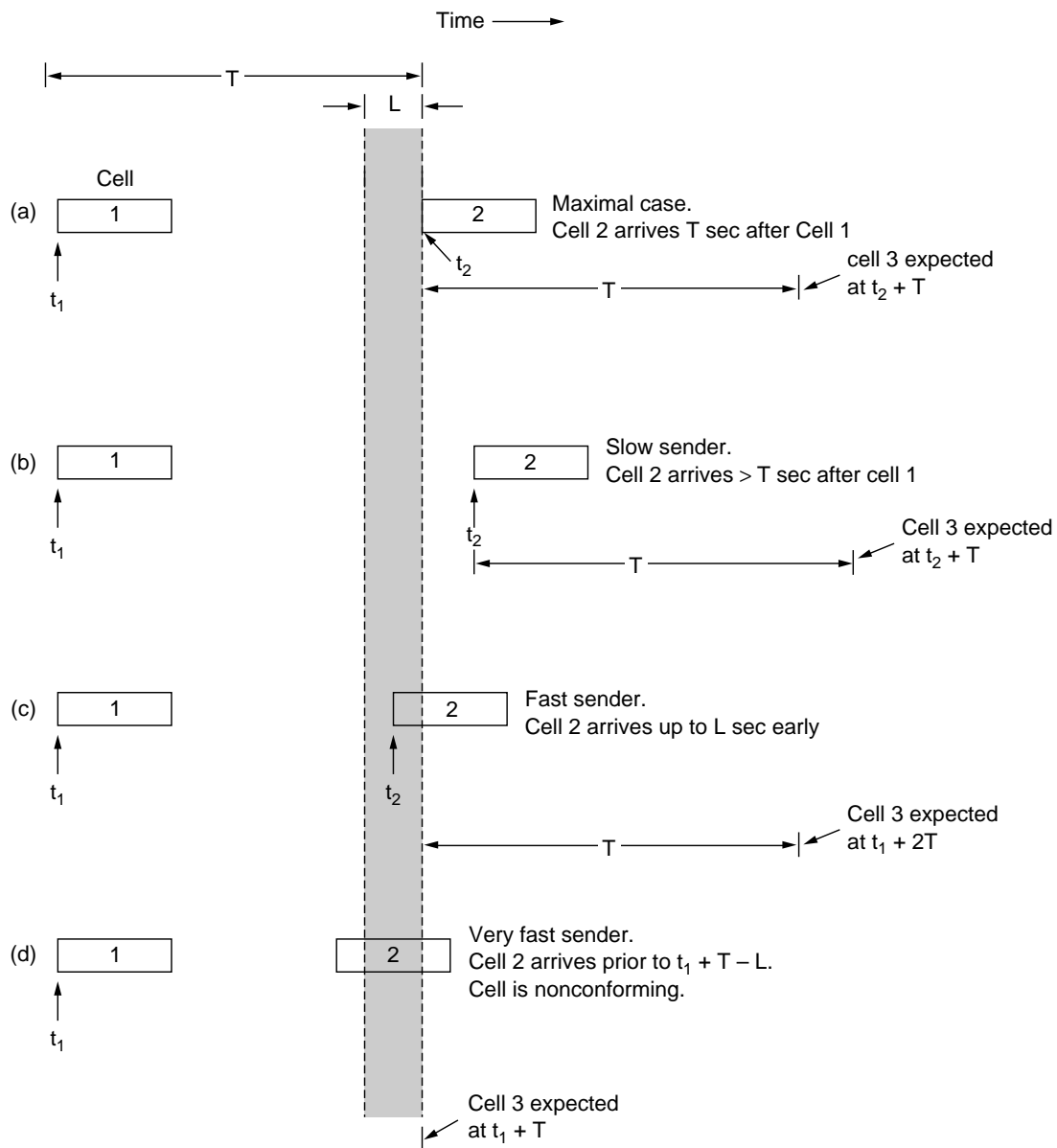                   |◄  at $t_1$ + T

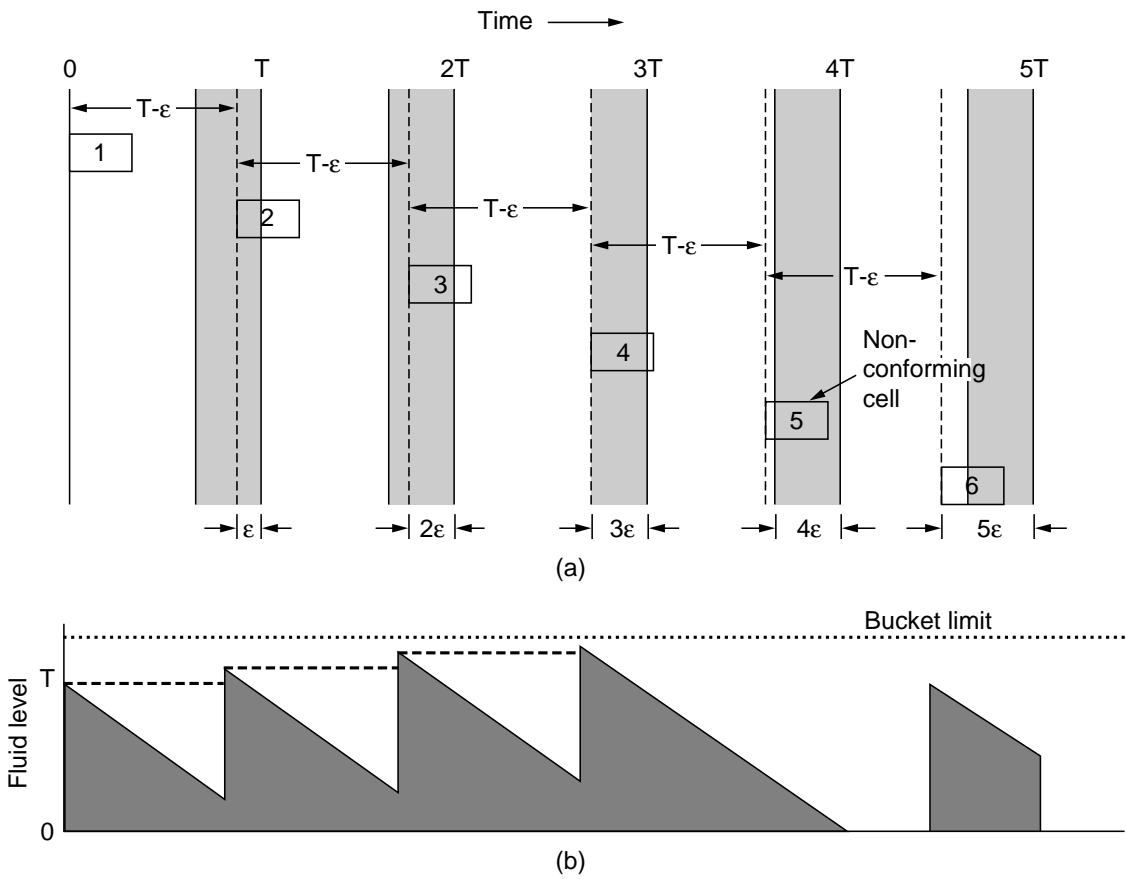Fig. 5-73. The generic cell rate algorithm.

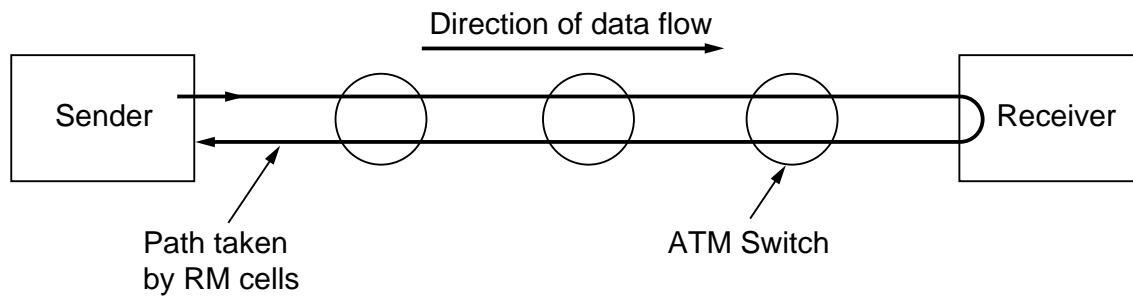Fig. 5-74. (a) A sender trying to cheat. (b) The same cell arrival pattern, but now viewed in terms of a leaky bucket.

Fig. 5-75. The path taken by RM cells.
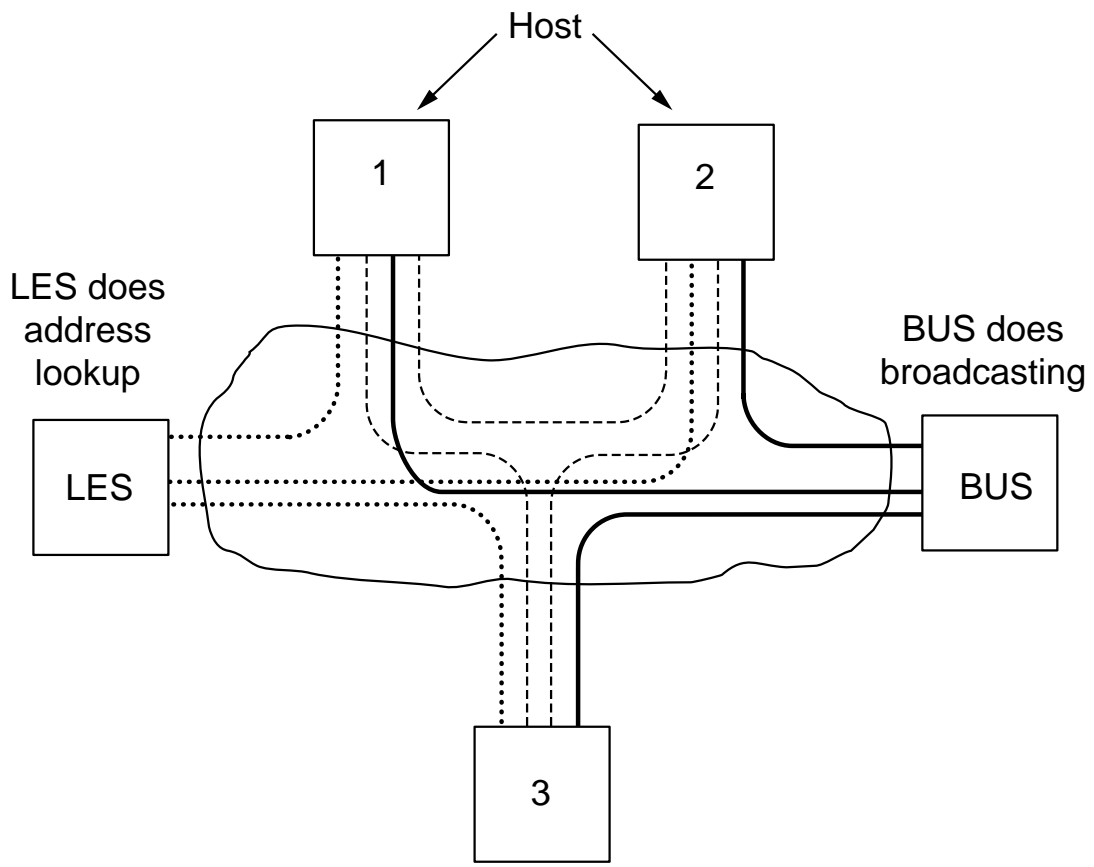
Host

1

2

LES does
address
lookup

LES

BUS does
broadcasting

BUS

3

Fig. 5-76. ATM LAN emulation.