

Cloud access and cloud interconnection networks

Unquestionably, communication is at the heart of cloud computing. The decades-long evolution of microprocessor and storage technologies, computer architecture and software systems, parallel algorithms and distributed control strategies enabled cloud computing yet, only *interconnectivity supported by a continually evolving Internet made cloud computing feasible*.

This chapter presents basic concepts necessary for understanding the intricacies of the communication infrastructure used for cloud computing. Communication and computing are inseparable concepts as discussed in Chapter 3, where we have seen that, even at the scale of a single core, efficient communication between the CPU, the memory, and the I/O sub-system are critical for optimal performance.

Communication, in the case of a single-multicore processor or on a system on a chip (SoC), is very efficient because it involves only hardware busses or networks on a chip and router-based packet switching networks between SoC modules. Large-scale systems such as supercomputers are built around complex and expensive interconnection networks controlled mostly by hardware thus support lower-latency and higher-bandwidth communication than cloud interconnects.

Distributed systems and computer clouds communicate over communication networks or the Internet. The Internet is a network of networks interconnected by switches operating under the control of routing algorithms and it is based on intricate software-based communication protocols that increase the communication latency and limit the bandwidth.

The designers of a cloud computing infrastructure are acutely aware that the farther data travels from the CPU the lower the bandwidth and the larger the communication latency. The limits of cloud interconnection networks bandwidth and latency are tested by the demands of systems with millions of servers; some cloud workloads are more affected by these limits than others.

Cloud workloads fall into four broad categories based on their dominant resource needs: CPU-intensive, memory-intensive, I/O-intensive, and storage-intensive. While the first two benefit from, but do not require, high-performing networking, the latter two do. Networking performance directly impacts the performance of I/O- and storage-intensive workloads.

The costs of the networking infrastructure continue to rise at a time when the cost of other components of cloud infrastructure continue to decrease. Moreover, applications in science and engineering are data- and network-intensive and require more expensive networks with higher bandwidth and lower latency than today's cloud interconnects.

This chapter starts with an overview of the Internet and the World Wide Web in Sections 6.1, 6.2, and 6.3, covering basic concepts on network architecture. The Internet is in fact *content-centric*, i.e., it is used primarily to access content, and Section 6.4 presents research ideas for a content-centric Internet. Software-defined networks, covered in Section 6.5, support an alternative to network management that enables dynamic, programmatically efficient network configuration to improve network performance and monitoring.

The focus then changes to the communication fabric used by the cloud infrastructure and to analysis of interconnection networks architecture and algorithms. After an overview of the interconnection networks in Section 6.6, the chapter covers multistage networks, Infiniband, and storage area networks in Sections 6.7, 6.8, and 6.9, respectively.

A scalable data-center architecture and network resource management are the topics of Sections 6.10 and 6.11. Then, the focus changes again, this time to content-delivery networks and vehicular networks in Sections 6.12 and 6.13. Further readings, historical notes, and a set of exercises and problems conclude the chapter.

6.1 Packet-switched networks and the Internet

Computer clouds are accessed through the Internet, a network of packet-switched networks. A packet-switched network transports data units called *packets* through a maze of *switches* where packets are queued and routed towards their destination. Packets are subject to random delays, loss, and may arrive at their final destination out of order.

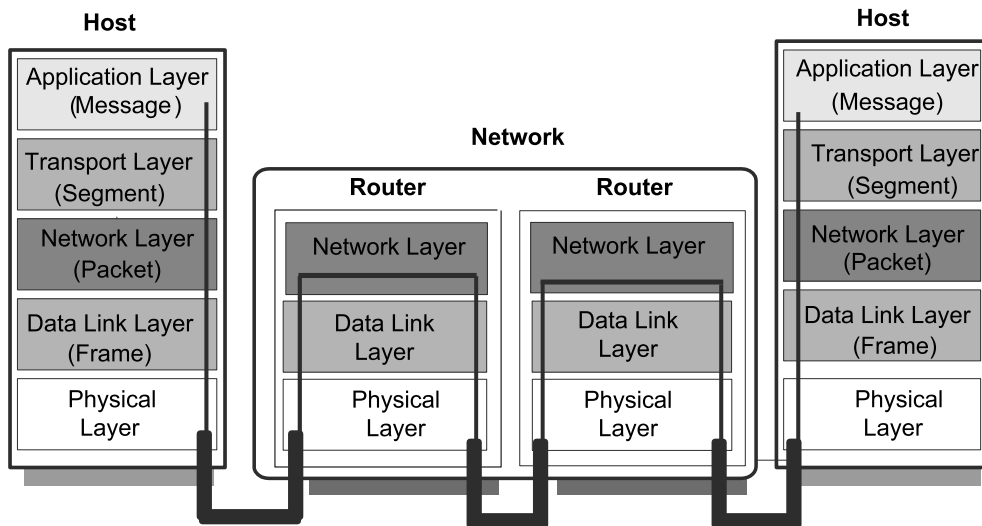
A few basic concepts are defined next. A *datagram* is the transfer unit in a packet-switched network. In addition to its *payload*, a datagram has a *header* containing control information necessary for its transport through the network. A *network architecture* describes the protocol stack used for communication. A *protocol* is a set of rules on how to communicate, specifying the actions taken by the sender and the receiver of a data unit. A network *host* identifies a system located at the network edge capable to initiate and to receive communication, a computer, a mobile device, e.g., a phone, or a sensor.

Network architecture and protocols. A packet-switched network has a *network core* consisting of routers and control systems interconnected by very-high-bandwidth communication channels and a *network edge* where the end-user systems reside.

A packet-switched network is a complex system consisting of a large number of autonomous components subject to complex and, sometimes contradictory requirements. Basic strategies for implementing a complex system are *layering* and *modularization*. Layering means decomposing a complex function into elements interacting through well-defined channels; a layer can only communicate with its adjacent layers. *Modularization* means dividing a system into interchangeable modules to create a flexible system while reducing the number of unique building blocks.

The Internet protocol stack, based on the TCP/IP network architecture, is shown in Fig. 6.1. Data flows down the protocol stack of the sending host from the application layer to the transport layer, then to the network layer, and to the data link layer. The physical layer pushes the streams of bits through a physical communication link encoded either as electrical, optical, or electromagnetic signals. At the receiving host, the packets flow up from the physical, to the data link, then to the network and transport layers and are finally delivered to the application layer. The corresponding data units for the five-layer architecture are: messages, segments, packets, frames, and encoded bits, respectively.

The *transport layer* is responsible for end-to-end communication, from an application running on the sending host to its peer running on the destination host using either TCP or UDP protocols. The network layer decides where the packet should be sent, either to another router, or to a destination host connected to a local area network connected to the router. IP, the *network layer* protocol, guides packets through the packet-switched network from the point of entry to the place where a packet exits

**FIGURE 6.1**

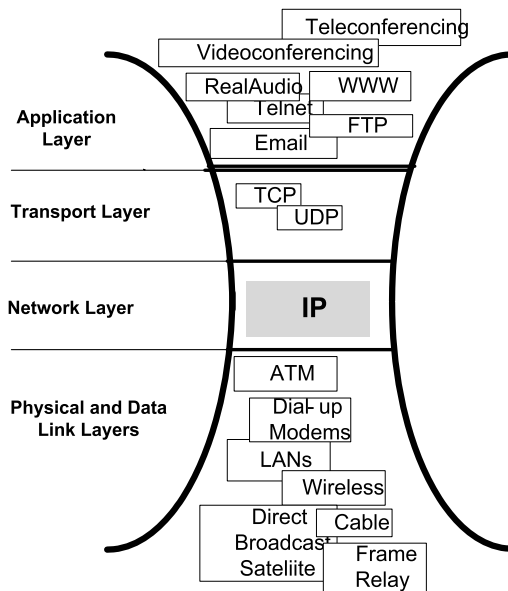
Internet protocol stack. Applications running on hosts at the edge of the network communicate using application-layer protocols. The transport layer deals with end-to-end delivery. The network layer is responsible for routing a packet through the network. The data link layer ensures reliable communication between adjacent nodes of the network, and the physical layer transports streams of bits encoded as electrical, optical, or electromagnetic signals (the thick lines represent such bit pipes).

the network. The *data link layer* encapsulates the packet for the communication link to the next hop. Once a packet reaches a router, the bits are passed to the data link and then to the network layer.

A protocol on one system communicates with its *peer* on another system. For example, the transport protocol on the sender, host A, communicates with the transport protocol on the receiver, host B. On the sending side, A, the transport protocol encapsulates the data from the application layer and adds control information as headers that can only be understood by its peer, the transport layer on host B. When the peer receives the data unit, it carries out a decapsulation, retrieves the control information, removes the headers, and then passes the payload to the next layer up, the application layer on host B.

The payload for the data link layer at the sending site includes the network header and the payload at the network layer. In turn, the network layer payload includes transport layer header and its payload consisting of the application layer header and application data.

The Internet. The Internet is a network of networks, a collection of separate, autonomous, and distinct networks. All networks adhere to a common framework and use: (i) globally unique IP addresses; (ii) the IP (Internet Protocol) routing protocol; and (iii) the Border Gateway Routing (BGP) protocol. BGP is a path vector reachability protocol that makes core routing decisions. BGP maintains a table of IP networks designating network reachability among autonomous systems. BGP makes routing decisions based on path, network policies, and/or rule sets.

**FIGURE 6.2**

The hourglass network architecture of the Internet. Regardless of the application, the transport protocol, and the physical network, all packets are routed from the source to the destination using the IP protocol and the IP address of the destination.

An *IP address* is a string of integers uniquely identifying every host connected to the Internet. An IP address allows the network to identify first the destination network and then the host in that network where a datagram should be delivered. A host may have multiple IP addresses, and it may be connected to more than one network. A host could be a supercomputer, a workstation, a laptop, a mobile phone, a network printer, or any other physical device with a network interface.

The Internet is based on a hourglass network architecture, as shown in Fig. 6.2. *The hourglass architecture is partially responsible for the explosive growth of the Internet*, it allowed the lower layers of the architecture to evolve independently from the upper layers. The communication technology drives dramatic changes of the lower layers of the Internet architecture, including the increase of the communication bandwidth and the widespread use of wireless networks and satellite communication. The software and the applications are the engines of progress for the upper layers of the architecture.

The hourglass model reflects the *end-to-end* architectural design principle. The model captures the fact that all packets transported through the Internet use IP to reach their destination. IP provides only best-effort delivery because any router along the path from the source to the destination may drop a packet when it is overloaded.

Another important architectural design principle of the Internet is the *separation between the routing and the forwarding planes*. This separation has allowed the forwarding plane to function while the routing evolved. The *forwarding plane* decides what to do with packets arriving on an inbound interface of a router. This plane uses a table to lookup the destination address of an incoming packet, then

retrieves the information to determine the path from the receiving interface to the proper outgoing interface(s) through the internal forwarding fabric of the router. The *routing plane* is responsible for building the routing table in each router.

In addition to the IP or logical address, each network interface, the hardware connecting a host with a network, has a unique *physical* or *MAC address*. While the MAC address is permanently assigned to a network interface of the device, the IP address may be dynamically assigned. The IP address of a mobile device changes depending on the device location and the network it is connected to.

The Dynamic Host Configuration Protocol (DHCP) is an automatic configuration protocol. DHCP assigns an IP address to a client system. A DHCP server has three methods of allocating IP addresses:

1. **Dynamic allocation**—a network administrator assigns a range of IP addresses to DHCP. During network initialization, each client computer on the LAN is configured to request an IP address from the DHCP server. The request-and-grant process uses a lease concept with a controllable time period, allowing the DHCP server to reclaim (and then reallocate) IP addresses that are not renewed.
2. **Automatic allocation**—the DHCP server permanently assigns a free IP address to a client from the range defined by the administrator.
3. **Static allocation**—the DHCP server allocates an IP address based on a manually filled in table with (MAC address–IP address) pairs. Only a client with a MAC address listed in this table is allocated an IP address.

Once a packet reaches the destination host, it is delivered to the proper transport protocol daemon which, in turn, delivers it to the application that listens to a *port*, an abstraction of the end-point of a logical communication channel, Fig. 6.3. The processes or threads running an application use an abstraction called *socket* to send and receive data through the network. A socket manages a queue of incoming messages and one for outgoing messages.

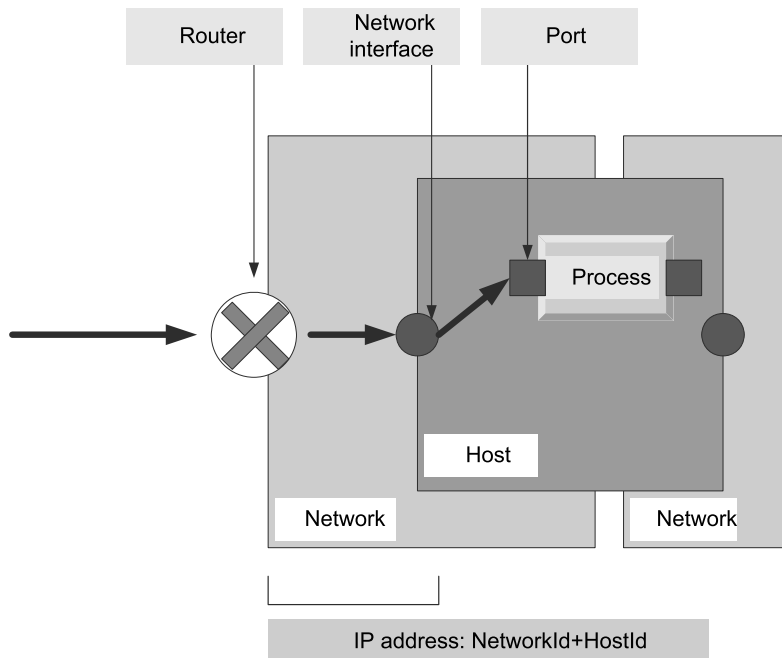
Internet transport protocols. Internet uses two transport protocols, a connectionless datagram protocol, UDP (User Datagram Protocol), and a connection-oriented protocol, TCP (Transport Control Protocol). The header of a datagram contains information sufficient for routing through the network from the source to the destination. The arrival time and the datagram delivery order are not guaranteed.

To ensure efficient communication, UDP assumes that error checking and error correction are either not necessary or performed by the application. Datagrams may arrive out of order, duplicated, or may not arrive at all. Applications using UDP include: DNS (Domain Name System), VoIP (Voice over IP), TFTP (Trivial File Transfer Protocol), streaming media applications such as IPTV, and online games.

TCP provides reliable, ordered delivery of a stream of bytes from an application on one system to its peer on the destination system; an application sends/receives data units called *segments* to/from a specific port, an abstraction of an endpoint of a logical communication link. TCP is the transport protocol used by the World Wide Web, email, file transfer, remote administration, and many other important applications.

TCP uses an end-to-end *flow-control mechanism* based on a sliding window, a range of packets the sender can send before receiving an acknowledgment from the receiver. These mechanisms enable the receiver to control the rate of segments sent and to process them reliably.

A network has a finite capacity to transport data, and when its load approaches capacity, we witness undesirable effects, such as routers start dropping packets and the delays and the jitter increase. An obvious analogy is a highway where the time to travel from point A to point B increases dramatically

**FIGURE 6.3**

Packet delivery to processes and threads; a packet is first routed by the IP protocol to the destination network and then to the host specified by the IP address. Applications listen to *ports*, abstractions of the endpoint of a communication channel.

in the case of congestion. A solution for traffic management is to introduce traffic lights limiting the rate at which new traffic is allowed to enter the highway, and this is precisely what the TCP emulates.

TCP uses several mechanisms for *congestion control* discussed in Section 6.3. These mechanisms control the rate of the data entering the network, keeping the data flow below a rate that would lead to a network collapse, and enforcing a fair allocation among flows. Acknowledgments coupled with timers are used to infer network conditions between the sender and receiver. TCP congestion control policies are based on four algorithms, *slow-start*, *congestion avoidance*, *fast retransmit*, and *fast recovery*. These algorithms use local information, such as the RTO (retransmission timeout) based on the estimated RTT (round-trip time) between the sender and receiver, as well as the variance in this round trip time to implement the congestion control policies. UDP is a connectionless protocol, thus there are no means to control the UDP traffic.

The review of basic networking concepts in this section shows why process-to-process communication incurs a significant overhead. While the raw speed of fiber-optic channels can reach Tbps, the actual transmission rate for end-to-end communication over a wide area network can only be on the order of Gbps, and the latency is on the order of milliseconds. The term “speed” is used informally to describe the maximum data-transmission rate, or the capacity of a communication channel; this ca-

capacity is determined by the physical bandwidth of the channel, and this explains why the term channel “bandwidth” is also used to measure the channel capacity, or the maximum data rate.

6.2 Internet evolution

The Internet is continually evolving under the pressure of its own success and the need to accommodate new applications and a larger number of users. Initially conceived as a data network, a network designed to transport data files, the Internet has morphed into today’s network supporting data streaming and applications with real-time constraints such as the Lambda service offered by the AWS. The discussion in this section is restricted to the aspects of the Internet evolution relevant to cloud computing.

Tier 1, 2, and 3 networks. To understand the architectural consequences of Internet evolution, we discuss first the relationship between two networks. *Peering* means that two networks exchange traffic between each other’s customers freely. *Transit* requires a network to pay another one for accessing the Internet. The term *customer* means that a network is receiving money to allow Internet access.

Based on these relationships, the networks are commonly classified as Tier 1, 2, and 3. A *Tier 1 network* can reach every other network on the Internet without purchasing IP transit or paying settlements; examples of Tier 1 networks are Verizon, ATT, NTT, and Deutsche Telecom; see Fig. 6.4.

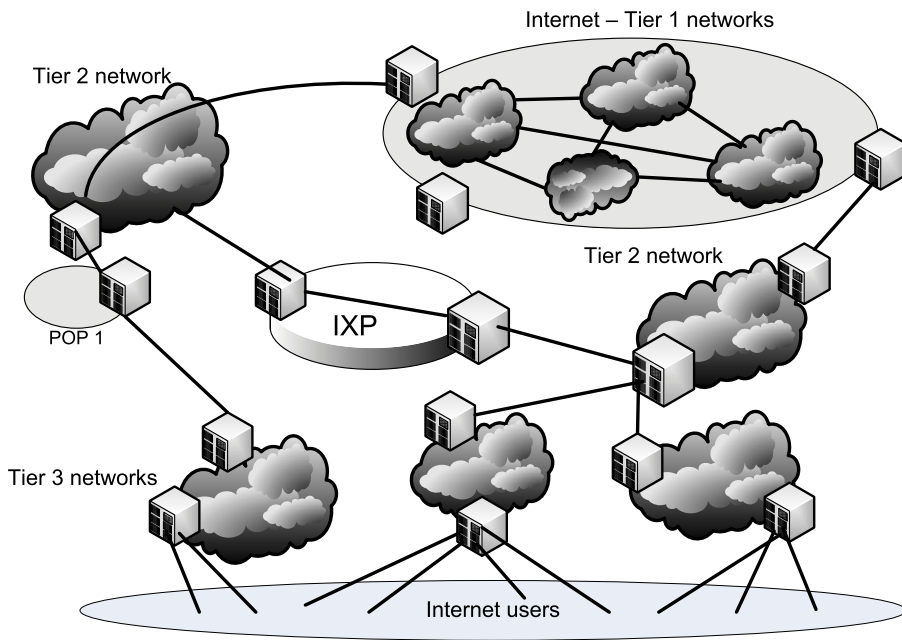
A *Tier 2 network* is an Internet service provider that engages in the practice of peering with other networks but still purchases IP transit to reach some portion of the Internet; Tier 2 providers are the most common providers on the Internet. A *Tier 3 network* purchases transit rights from other networks (typically Tier 2 networks) to reach the Internet. A *point-of-presence (POP)* is an access point from one place to the rest of the Internet.

An *Internet exchange point (IXP)* is a physical infrastructure enabling *Internet Service Providers (ISPs)* to exchange Internet traffic. IXPs interconnect networks directly via the exchange, rather than through one or more third-party networks. The advantages of the direct interconnection are numerous, but the primary reasons to implement an IXP are cost, latency, and bandwidth. Traffic passing through an exchange is typically not billed by any party, whereas traffic to an ISP’s upstream provider is.

IXPs reduce the portion of an ISP’s traffic that must be delivered via their upstream transit providers, thereby reducing the average per-bit delivery cost of their service. Furthermore, the increased number of paths found through the IXP improves routing efficiency and fault tolerance. A typical IXP consists of one or more network switches, to which each of the participating ISPs connects.

New technologies, such as web applications, cloud computing, and content-delivery networks, are reshaping the definition of a network. The World Wide Web, gaming, and entertainment are merging, and more computer applications are moving to the cloud. Data streaming consumes an increasingly larger fraction of the available bandwidth as high definition TV sets become less expensive and content providers such as Netflix and Hulu offer customers services that require a significant increase in the network bandwidth.

Does the network infrastructure keep up with the demand for bandwidth? A natural question to ask is: Where is the actual bottleneck limiting the bandwidth available to a typical Internet broadband user? The answer is: the “last mile,” the link connecting home to the ISP network. Recognizing that the broadband access infrastructure ensures continual growth of the economy and enables people to

**FIGURE 6.4**

There are three classes of networks, Tier 1, 2, and 3; an IXP is a physical infrastructure allowing ISPs to exchange Internet traffic.

work from any site, Google initiated in early 2010 a project to provide 1 Gbps access to individual households through FTTH.¹ ATT FTTH expects to have some seven million subscribers by 2022.

Migration to IPv6. Internet Protocol, Version 4 (IPv4), provides an addressing capability of 2^{32} , or approximately 4.3 billion IP addresses, a number that proved to be insufficient. Indeed, the Internet Assigned Numbers Authority (IANA) assigned the last batch of *five* address blocks to the Regional Internet Registries in February 2011, officially depleting the global pool of completely fresh blocks of addresses; each of the address blocks represents approximately 16.7 million possible addresses. Network Address Translation (NAT) provides a partial solution to this problem because it allows a single public IP address to support hundreds or even thousands of private IP addresses.

Internet Protocol Version 6 (IPv6) provides an addressing capability of 2^{128} , or 3.4×10^{38} addresses. There are other major differences between IPv4 and IPv6:

¹ The fiber-to-the-home (FTTH) is a broadband network architecture that uses optical fiber to replace the copper-based local loop used for the last-mile network access to homes.

1. *Multicasting.* IPv6 does not implement traditional IP broadcast, i.e., the transmission of a packet to all hosts on the attached link using a special broadcast address and, therefore, does not define broadcast addresses. IPv6 supports new multicast solutions, including embedding rendezvous point addresses in an IPv6 multicast group address. This solution simplifies the deployment of inter-domain solutions.
2. *Stateless address auto-configuration (SLAAC).* IPv6 hosts can configure themselves automatically using router-discovery messages when connected to a routed IPv6 network using ICMPv6 (Internet Control Message Protocol v6). When first connected to a network, a host sends a link-local router solicitation multicast request for its configuration parameters. If suitably configured, routers respond to such a request with a router advertisement packet that contains network-layer configuration parameters.
3. *Mandatory support for network security.* Internet Network Security (IPsec) is an integral part of the base protocol suite in IPv6, while it is optional for IPv4. IPsec is a protocol suite² operating at the IP layer. Each IP packet is authenticated and encrypted. Other security protocols, e.g., Secure Sockets Layer (SSL), Transport Layer Security (TLS), and Secure Shell (SSH) operate at the upper layers of the TCP/IP suite.

Migration to IPv6 is a very challenging and costly proposition because it involves upgrading applications, hosts, routers, and DNS infrastructure. Moving to IPv6 requires backward compatibility, and any organization migrating to IPv6 should maintain a complete IPv4 infrastructure [113].

6.3 TCP congestion control

Network congestion occurs when a router or a communication channel is forced to carry out more traffic than it can handle. Congestion leads to increased delays, packet loss, and retransmissions, blocking new connections, and ultimately to lower network capacity. Increased round-trip time of a *connection-oriented protocol* signals network congestion and could trigger congestion control mechanisms. *Connectionless protocols* like UDP cannot detect network congestion and cannot contribute to congestion control or avoidance.

TCP is a connection-oriented Internet transport protocol. TCP uses two mechanisms, flow control and congestion control, to achieve high channel utilization, avoid congestion, and, at the same time, ensure a fair sharing of the network bandwidth. The *flow-control* mechanism throttles the sender feedback from the receiver and forces the sender to transmit only the amount of data the receiver is able to buffer and then process. TCP uses a sliding-window flow-control protocol. If W is the window size, then the sender data rate S is:

$$S = \frac{W \times MSS}{RTT} \text{ bps} = \frac{W}{RTT} \text{ packets/second} \quad (6.1)$$

² IPsec uses several protocols: (1) Authentication Headers (AH) supports connectionless integrity, data origin authentication for IP datagrams, and protection against replay attacks; (2) Encapsulating Security Payloads (ESP) supports confidentiality, data-origin authentication, connectionless integrity, an anti-replay service, and limited traffic-flow confidentiality; (3) Security Associations (SA) provide the parameters necessary to operate the AH and/or ESP operations.

where MSS and RTT denote the Maximum Segment Size and the Round-Trip Time, respectively, assuming that MSS is one packet. If S is too small, the transmission rate is smaller than the channel capacity, while a large S leads to congestion. Channel capacity depends on the network load because the physical channels along the path of a flow are shared with many other Internet flows.

The actual window size W is affected by two factors: (a) the ability of the receiver to accept new data and (b) the sender's estimation of the available network capacity. The receiver specifies the amount of additional data it is willing to accept in the *receive window* field of every frame. The receiver's window shifts when the receiver receives and acknowledges a new segment of data. When a receiver advertises a window size of zero, the sender stops sending data and starts the *persist timer*. This timer is used to avoid deadlock when a subsequent window-size update from the receiver is lost.

When the persist timer expires, the sender sends a small packet, and the receiver responds by sending another acknowledgment containing the new window size. In addition to the flow control provided by the receiver, the sender attempts to infer the available network capacity and to avoid overloading the network. The source uses the losses and the delay to determine the level of congestion. If $awnd$ denotes the receiver window and $cwnd$ the congestion window set by the sender, the actual window should be:

$$W = \min(cwnd, awnd). \quad (6.2)$$

Several algorithms are used to calculate $cwnd$, including Tahoe and Reno, developed by Van Jacobson in 1988 and 1990. Tahoe was based on slow start (SS), congestion avoidance (CA), and fast retransmit (FR). The sender probes the network for spare capacity and detects congestion based on loss. *Slow start* means that the sender starts with a window of two times MSS, $init_cwnd = 1$. For every packet acknowledged, the congestion window increases by one MSS so that the congestion window effectively doubles for every RTT.

When the congestion window exceeds the threshold, *i.e.*, $cwnd \geq ssthresh$, the algorithm enters the CA state. In the CA state, on each successful acknowledgment $cwnd \leftarrow cwnd + 1/cwnd$ and on each RTT $cwnd \leftarrow cwnd + 1$. *Fast retransmit* is motivated by the fact that the timeout is too long thus, a sender retransmits immediately after three duplicate acknowledgments without waiting for a timeout. Two adjustments are made in this case:

$$flightsize = \min(awnd, cwnd) \quad \text{and} \quad ssthresh \leftarrow \max(flightsize/2, 2), \quad (6.3)$$

and the system enters in the slow start state, $cwnd = 1$. The pseudocode describing the Tahoe algorithm is:

```

for every ACK {
    if (W < ssthresh) then W++      (SS)
    else    W += 1/W                (CA)
}
for every loss {
    ssthresh = W/2
    W = 1
}

```

While the majority of the Internet traffic is due to long-lived, bulk-data transfers, *e.g.*, video and audio streaming, the majority of transactions, *e.g.*, web requests, are short-lived. So, a major challenge is to ensure some fairness for short-lived transactions.

To overcome the limitations of the slow start, application strategies have been developed to reduce the time needed to download data over the Internet. For example, two browsers, Firefox 3 and Google Chrome, open up to six TCP connections per domain to increase the parallelism and to boost start-up performance when downloading a web page. Internet Explorer 8 opens 180 connections. Clearly, these strategies circumvent the mechanisms for congestion control and incur a considerable overhead. A better solution is to increase the initial congestion window of TCP for several reasons [155]:

1. TCP latency is dominated by the number of RTT's during the slow start phase. Increasing the *init_cwnd* parameter enables the data transfer to be completed with fewer RTTs.
2. A single TCP connection requires multiple RTTs to download a single page because the average page size is 384 KB.
3. This will ensure fairness between short-lived transactions, which are a majority of Internet transfers, and long-lived transactions, which transfer very large amounts of data.
4. This will allow faster recovery after losses through Fast Retransmission.

It can be shown that the latency of a transfer completing during the slow start without losses is given by the expression:

$$\left\lceil \log_{\gamma} \left(\frac{L(\gamma - 1)}{\text{init_cwnd}} + 1 \right) \right\rceil \times RTT + \frac{L}{C}, \quad (6.4)$$

with L being the transfer size, C is the bottleneck-link rate, and γ is a constant equal to 1.5 or 2 depending on whether the acknowledgments are delayed or not; $L/\text{init_cwnd} \geq 1$. In the experiments reported in [155], the TCP latency was reduced from about 490 msec when *init_cwnd* = 3 to about 466 msec for *init_cwnd* = 16.

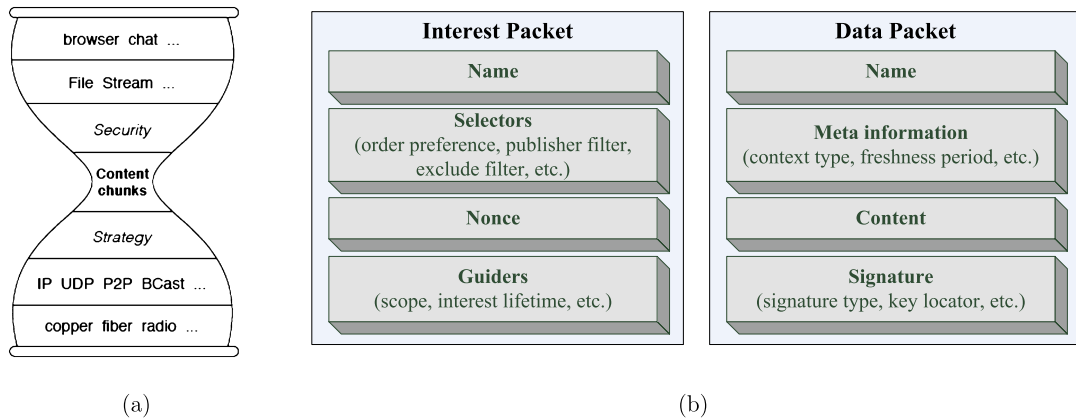
6.4 Content-centric networks; named data networks (R)

Internet data is tied to a particular host, and packets only name the communication end points; this makes data replication and migration difficult. The Internet is mostly used by applications in digital media, electronic commerce, Big Data analytics, etc. as a *distribution network*. In a distribution network, communication is *content-centric*, named objects are requested by an agent, and, once a site holding the object is located, the network transports it to the site that requested the object. The end user in a distribution network is oblivious to the location of the data and is only interested in the content.

The idea of a content-centric network has been around for some time. In 1999, the TRIAD project at Stanford³ [210,479] proposed to use an object name for routing towards the object replica. In this proposal, the Internet Relay Protocol performs name-to-address conversion using routing information maintained by relay nodes. The Name-Based Routing Protocol performs a function similar to the BGP protocol, supporting a mechanism for updating routing information in relay nodes.

In 2006 the Data Oriented Network Architecture (DONA) project at UC Berkeley [149] extended TRIAD incorporating security and persistence as primitives of the new network architecture. DONA architecture exposes two primitives: FIND—allows a client to request a particular piece of data by

³ TRIAD stands for Translating Relaying Internet Architecture Integrating Active Directories.

**FIGURE 6.5**

Named Data Networks. (a) NDN hourglass architecture parallels the one of the Internet; it separates the lower layers of the protocol stack from the upper ones, thus data naming can evolve independently from networking. (b) NDN networking service semantics is to fetch a data chunk identified by name, while Internet semantics is to deliver a packet to a given network address through an end-to-end channel identified by a source and a destination IP address. Interest packets and data packets are used for NDN routing and forwarding.

its name and REGISTER—enables content providers to indicate their intent to serve a particular data object. In 2006 Van Jacobson from UCLA argued that Named Data Networks (NDNs) should be the architecture of the future Internet.

In 2012, the Internet Research Task Force (IRTF) established an information-centric networking (ICN) research working group for investigating architecture designs for NDN. A 2014 survey of ICN research and a succinct presentation of NDNs can be found in [522] and [539], respectively. The important features of the NDN architecture addressed by the current research efforts include namespaces, trust models, in-network storage, data synchronization, and, last but not least, rendezvous, discovery, and bootstrapping.

The hourglass model can be extended, and packets can name objects rather than communication endpoints. NDN's hourglass model separates the lower layers of the protocol stack from the upper ones, thus data naming can evolve independently from networking, as in Fig. 6.5(a). NDN packet delivery is driven by data consumers; communication is initiated by an agent generating an *Interest* packet containing the name of the data as in Fig. 6.5(b).

Once the Interest packet reaches a network host that has a copy of the data item, a *Data* packet containing the name, the contents of the data, and the signature is generated. The signature consists of producer's key. The *data* packet follows the route traced by the *interest* packet, and it is delivered to the data consumer agent. The information necessary to forward a packet is maintained by an NDN router in several places:

1. *Content Store*—local cache for *Data* packets previously crossing the router. When an interest packet arrives, a search of the content store determines if a matching data exists, and, if so, the data is forwarded on the same router interface the interest packet was received. If not, the router uses

a data structure, the Forwarding Information Base (FIB), to forward the packet. More recently, NDN supports more persistent and larger-volume in-network storage, called Repositories, providing services similar to that of the Content Delivery Networks.

2. *Forwarding Information Base*—FIB entries are populated by a name-prefix-based procedure. The Forwarding Strategy retrieves the longest prefix matched entry from the forwarding information base for an interest packet.
3. *Pending Interest Table (PIT)*—stores all interest packets the router has forwarded but have not yet been satisfied. A PIT entry records the data name carried in the Interest, together with its incoming and outgoing router interface(s). When a data packet arrives, the router finds the matching PIT entry and forwards the data to all downstream interfaces listed in that PIT entry; then, it removes the PIT entry and caches the data packet in the Content Store.

Names are essential in NDN, though namespace management is not part of the NDN architecture. The scope and contexts of NDN names are different. Globally accessible data must have globally unique names, while local names require only local routing and must only be locally unique.

There are significant differences between NDN and TCP/IP. NDN namespace cannot be exhausted, while IPv4 address limitation forced migration to IPv6. NDN supports data-centric security; each *Data* packet is cryptographically signed, while TCP/IP security is left to the communication endpoints. An NDN router announces name prefixes for data it is willing to serve, while an IP router announces IP prefixes.

NDN is a universal overlay.⁴ NDN can run over any datagram network and, conversely, any datagram network, including IP, can run over NDN. Classical algorithms such as Open Shortest Path First (OSPF) route data chunks using component-wise, longest-prefix match of the name in an interest packet with the FIB entries of a router.

Wide-area applications can operate over IP tunnels, and islands of NDN nodes could be interconnected by tunneling over non-NDN clouds. Scalable forwarding, along with robust and effective trust management, are critical challenges for the future of NDN networks. Namespace design is at the heart of NDNs because it involves application data, communication, storage models, routing, and security.

NDN could be useful for cloud data centers and, in particular, for those supporting the IaaS cloud delivery model. Data is replicated, and multiple instances could access data concurrently from their nearest storage servers. This would be useful for some MapReduce applications, but it would require major changes in frameworks supporting this paradigm. On the other hand, many applications caching data at the server side would only marginally benefit from NDN support.

6.5 Software-defined networks; SD-WAN

Software-defined Networks (SDNs) allow programmatic control of communication networks extending the basic principles of resource virtualization to networking. SDNs introduce an abstraction layer separating network configuration from the physical communication resources. A network operating system running inside the control layer, sandwiched between the application layer and the infrastructure layer,

⁴ An overlay network is a network built on top of another physical network; its nodes are connected by virtual links, each of which corresponds to a path, possibly through many physical links, in the underlying network.

enables applications to reconfigure dynamically the communication substrate to adapt to their security, scalability, and manageability needs.

Though enthusiastically embraced by networking vendors, there is little agreement either on the architecture, the APIs, or the overlay networks among vendors. A 2012 white paper of the Open Network Foundation (ONF) (<https://www.opennetworking.org/sdn-resources>) promoted OpenFlow-based SDNs. According to ONF, a new network architecture, supporting traffic patterns changes due to extensive use of cloud services and the need for more bandwidth demanded by Big Data applications is necessary. ONF architecture includes several components:

- a. Applications—generate network requirements to controllers.
- b. Controllers—translate application requirements to SDN datapaths and provide an abstract view of the network to applications.
- c. Datapaths—logical network devices that expose control to the physical substrate, consisting of agents and forwarding engines.
- d. Control-to-data-plane—interfaces between SDN controllers and SDN datapaths.
- e. Northbound interfaces—interfaces between applications and controllers.
- f. Interface drivers and agents—driver—agent pairs.
- g. Management and administration.

OpenFlow is an API for programming data plane switches. The datapath and the control paths of an OpenFlow switch consist of a flow table, including an action associated with each flow entry and a controller that programs the flow entry. The controller configures and manages the switch and receives events from the switch.

Software-Defined Wide Area Network (SD-WAN). Informally, SD-WAN is a network where there is a separation between the data service and the rest of the protocol stack. SD-WAN changes WAN from being a hardware-centric network to a software-defined service. An SD-WAN creates a virtual overlay across the underlying data services, separating the upper protocol stack from the network. This separation allows organizations to use the best transport for each application. Transport services used by SD-WANs are:

1. Multiprotocol Label Switching (MPLS)—network routing that directs data from one node to the next based on short path labels rather than network addresses. MPLS avoids complex routing table lookups and speeds traffic flows.
2. Long-Term Evolution/4G/5G. LTE is associated with 4G and 5G wireless networks.
3. Broadband Internet connections.
4. DSL/Cable.

SD-WAN expands network through software rather than hardware and provides a simple interface for WAN management. SD-WAN is a cloud-scale architecture, open and scalable. SD-WAN uses a centralized control function to direct traffic and allows for more consistent, predictable app performance, and improved security. SD-WAN deployment can include existing router/switches, or even virtualized customer-premises equipment.

SD-WAN main advantages: (i) optimized user experience and efficiency for cloud applications; (ii) reduced cost for IT management due to automation and reduced connectivity costs by moving from expensive MPLS to lower-cost connections; (iii) improved organization traffic management; (iv) includes next generation firewalls with malware protection; (v) increased network security with encrypted

network traffic and network segmentation to limit and manage attacks; (vi) automated common tasks and configuration; (vii) redundancy to improve availability and reduce risk/failure; and (viii) easy deployment.

Citrix, Riverbed <https://www.riverbed.com/faq/what-is-sd-wan.html>, CISCO, and Cato Cloud <https://www.catonetworks.com/sd-wan/> are SD-WAN vendors. CISCO offers Meraki <https://meraki.cisco.com/products/security-sd-wan/> as well as Viptela SD-WANs. Citrix <https://www.citrix.com> SD-WAN optimize delivery of virtual, cloud, and SaaS applications on a secure, flexible WAN.

6.6 Interconnection networks for computer clouds

Computing and communication are deeply intertwined, as we have seen in Chapter 3. Interconnection networks discussed in the next sections are critical for the performance of computer clouds and supercomputers. Concepts important for understanding interconnection networks are introduced next.

A network consists of nodes and links or communication channels. The *degree* of a node is the number of links the node is connected with. The *nodes* of an interconnection network could be processors, memory units, or servers. The *network interface* of a node is the hardware connecting it to the network. *Switches* and *communication channels* are the elements of the interconnection fabric. Switches receive data packets, examine each packet to identify the destination IP addresses, and then use the routing tables to forward it to the next hop towards its final destination. An *n-way switch* has *n* ports connected to *n* communication links. An interconnection network can be *non-blocking* if it is possible to connect any permutation of sources and destinations at any time or *blocking* if this requirement is not satisfied.

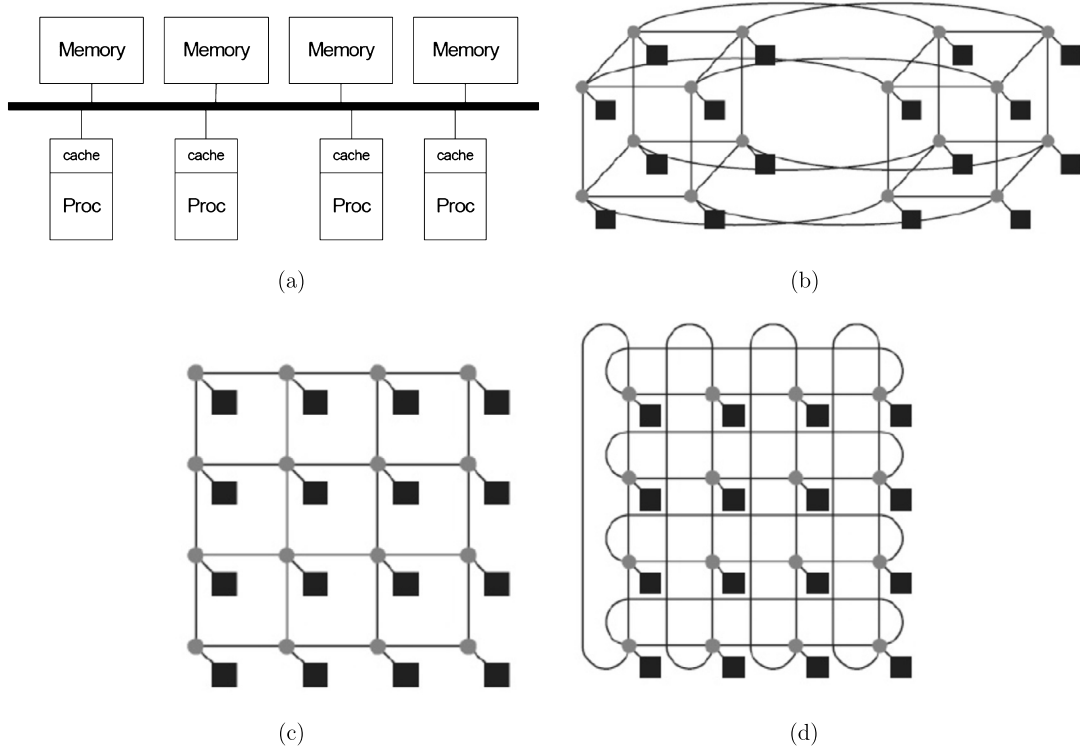
While processor and memory technology have followed Moore's law, interconnection networks have evolved at a slower pace and have become a major factor in determining the overall performance and cost of large-scale systems. For example, from 1997 to 2010, the speed of Ethernet networks has increased from one to 100 Gbps. This increase is slightly slower than the Moore's law for traffic [359] which predicted a 1-Tbps Ethernet by 2013.

Interconnection networks are distinguished by their topology, routing, and flow control. *Network topology* is determined by the way nodes are interconnected, routing decides how a message gets from its source to destination, and the flow control negotiates how the buffer space is allocated. There are two basic types of network topologies:

- Static networks where there are direct connections between servers;
- Switched networks where switches are used to interconnect the servers.

The topology of an interconnection network determines the *network diameter*, the average distance between all pairs of nodes, the *bisection width*, the minimum number of links cut to partition the network into two halves, and the bisection bandwidth, as well as the cost and power consumption [273]. When a network is partitioned into two networks of the same size, the bisection bandwidth measures the communication bandwidth between the two.

The *full bisection bandwidth* allows one half of the network nodes to communicate simultaneously with the other half of the nodes. Assume that half of the nodes inject data into the network at a rate *B* Mbps. When the bisection bandwidth is *B*, then the network has full bisection bandwidth. The most popular topologies with a static interconnect are:

**FIGURE 6.6**

Static networks: (a) Bus; (b) Hypercube; (c) 2D-mesh; (d) 2D-torus.

1. *Bus*, a simple and cost-effective network, see Fig. 6.6(a). It does not scale, but it is easy to implement cache coherence through snooping for distributed memory systems. A bus is often used in shared-memory multiprocessor systems.
2. *Hypercube* of order n ; see Fig. 6.6(b). A hypercube has a good bisection bandwidth; the number of nodes is $N = 2^n$, the degree is $n = \log N$, and the average distance between nodes is $\mathcal{O}(N)$ hops. Example of use: SGI Origin 2000.
3. *2D-mesh*; see Fig. 6.6(c). An $n \times n$ 2D-mesh has many paths to connect nodes, a cost of $\mathcal{O}(n)$, and an average latency of $\mathcal{O}(\sqrt{n})$. A mesh is not symmetric on edges, thus its performance is sensitive to the placement of communicating nodes on edges versus the middle. Example of use: Intel Paragon supercomputer of the 1990s.
4. *Torus*, avoids mesh asymmetry but has a higher cost due to a larger number of components. A torus is effective for applications using nearest-neighbor communication; see Fig. 6.6(d). It is prevalent for proprietary interconnects. Example of use: Six-D Mesh/torus of Fujitsu K supercomputer.

Switched networks have multiple layers of switches connecting the nodes: (i) a crossbar switch has N^2 crosspoint switches; see Fig. 6.7(a); (ii) Omega (Butterfly, Benes, Banyan, etc.) have $(N \log N)/2$

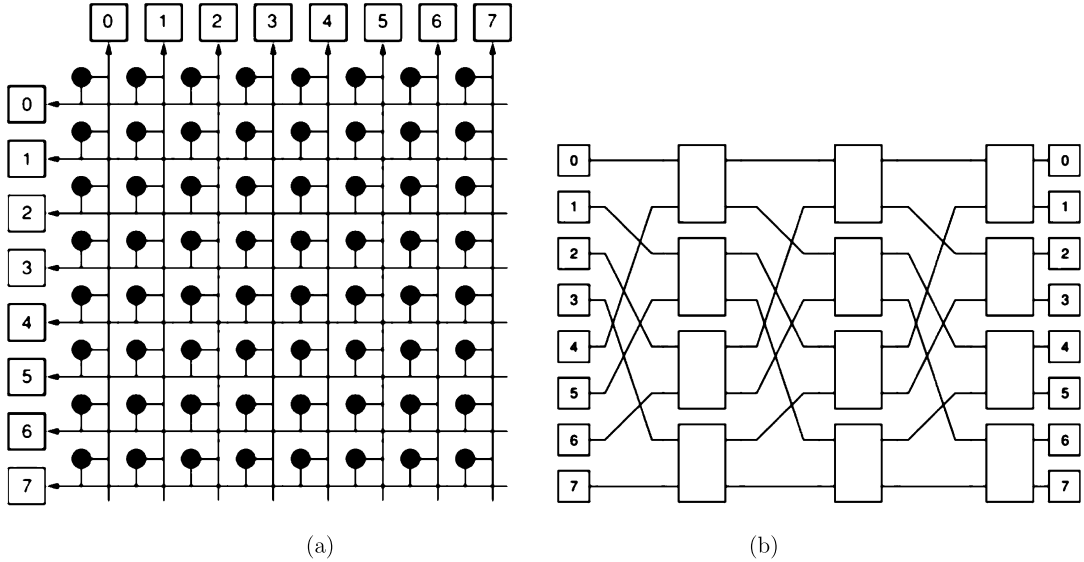


FIGURE 6.7

Switched networks: (a) An 8×8 crossbar switch. Sixteen nodes are interconnected by 49 switches represented by the dark circles; (b) An 8×8 Omega switch. Sixteen nodes are interconnected by 12 switches represented by white rectangles.

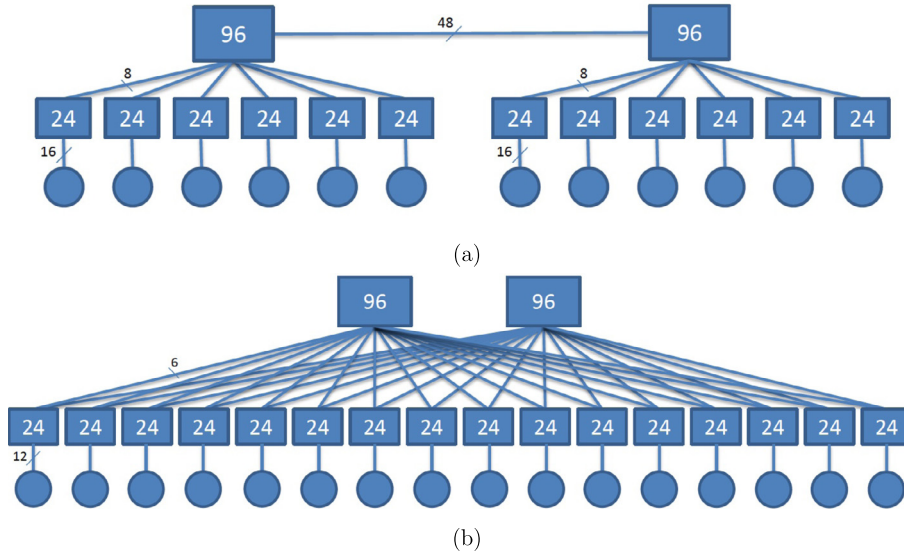
switches; see Fig. 6.7(b). The cost is $\mathcal{O}(N \log N)$ and the latency is $\mathcal{O}(\log N)$. Omega networks are small diameter networks.

Cloud interconnection networks. The cloud infrastructure consists of one or more warehouse-scale computers discussed in Section 4.2. WSCs have an infrastructure with a very large number of servers interconnected by high-speed networks.

The networking infrastructure of a cloud must satisfy several requirements, including scalability, cost, and performance. The network should allow low-latency, high-speed communication, and, at the same time, provide *location transparent communication* between servers. In other words, each server should be able to communicate with every other server with similar speed and latency. This requirement ensures that *applications need not be location aware*, and, at the same time, it reduces the complexity of the system management.

Typically, the networking infrastructure is organized hierarchically. The servers of a WSC are packed into racks and interconnected by a rack router. Then, rack routers are connected to cluster routers, which in turn are interconnected by a local communication fabric. Finally, inter-data center networks connect multiple WSCs [279]. Clearly, in a hierarchical organization, true location transparency is not feasible. Cost considerations ultimately decide the actual organization and performance of the communication fabric.

Oversubscription, defined as the ratio of the worst-case achievable aggregate bandwidth among the servers to the total bisection bandwidth of the interconnect, is a useful measure of the fitness of an

**FIGURE 6.8**

One hundred ninety-two node fat-tree interconnection network with two 96-way and 12 24-way switches in a computer cloud. (a) The two 96-way switches at the root are connected with one another and with six 24-way switches; (b) Each 96-way switch is connected to each one of the 12 24-way switch.

interconnection network for a large cluster. An oversubscription of one-to-one indicates that any host may communicate with an arbitrary hosts at the full bandwidth of the interconnect. An oversubscription value of four to one means that only 25% of the bandwidth available to servers can be attained for some communication patterns. Typical oversubscription figures are in the 2.5 to 1 and 8 to 1 range.

The cost of routers and of cables interconnecting the routers are major components of the overall cost of an interconnection network. Wire density has scaled up at a slower rate than processor speed, and wire delay has remained constant over time, thus better performance and lower costs can only be achieved with innovative router architecture. This motivates us to take a closer look at routers design.

Fat-trees are optimal interconnects for large-scale clusters and, by extension, for WSCs. Servers are placed at the leaves, and switches populate the root and the internal nodes of the tree. Fat-trees have additional links to increase the bandwidth near the root of the tree. A fat-tree interconnect can be built with inexpensive commodity parts since all switching elements of a fat-tree are identical. A fat-tree is a particular instance of Clos topology discussed in Section 6.7.

Two 192 node fat-trees built with two 96-way switches and 12 24-way switches are shown in Fig. 6.8. The two 96-way switches at the root are connected via 48 links and each 24-way switch has 6×8 uplink connections to the root and 6×16 down connections to 16 servers in Fig. 6.8(a), while in Fig. 6.8(b) each 96-way switch is connected via six links to each one of the 12 24-way switches connected via 12 links with servers.

Table 6.1 from [232] summarizes the performance/cost for a 2D-mesh, a 2D-torus, a hypercube of order seven, a fat-tree, and a fully connected network. Two dimensions of interconnection network

Table 6.1 A side-by-side comparison of performance and cost figures of several interconnection network topologies for 64 nodes.

Property	2D mesh	2D torus	Hypercube	Fat-tree	Fully connected
BW in # of links	8	16	32	32	1 024
Max/Avg hop count	14/7	8/4	6/3	11/9	1/1
I/O ports per switch	5	5	7	4	64
Number of switches	64	64	64	192	64
Total number of links	176	192	256	384	2 080

performance, the bisection bandwidth and the average and maximum number of hops, along with three elements affecting the cost, the number of ports per switch, the number of switches, and the total number of links, are shown. Fat-tree interconnects have the largest bisection bandwidth with the smallest number of I/O ports per switch, while a fully connected interconnect has a prohibitively large number of links.

6.7 Multistage interconnection networks

Low-radix routers have a small number of ports. *High-radix* routers have a large number of ports, a considerably smaller number of stages and enjoy lower latency and reduced power consumption. High-radix chips divide the bandwidth into larger number of narrow ports, while low-radix chips divide the bandwidth into a smaller number of wide ports. Every five years during the past two decades, the pin bandwidth of the chips used for switching has increased by an order of magnitude as a result of signaling rate increase and larger number of signals.

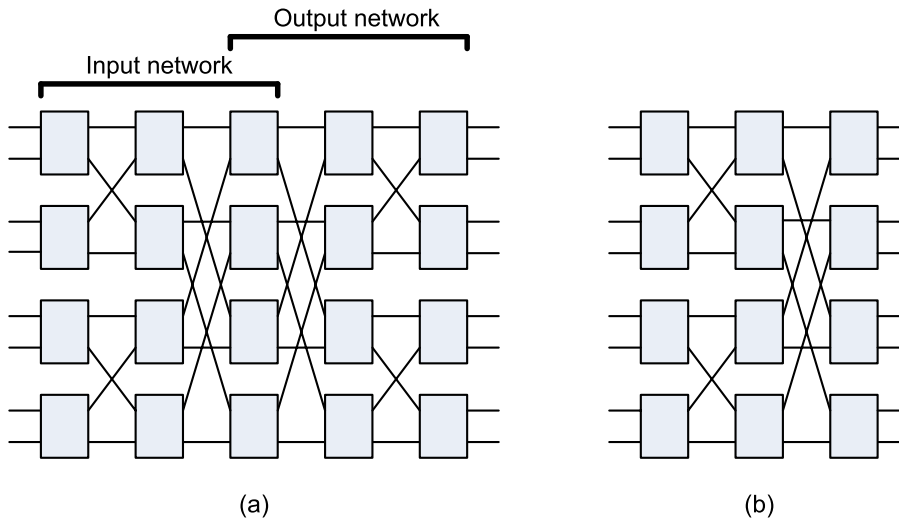
Clos networks were invented in early 1950s by Charles Clos from Bell Labs [110]. A Clos network is a multistage nonblocking network with an odd number of stages; see Fig. 6.9(a). The network consists of two butterfly networks where the last stage of the input is fused with the first stage of the output.

The name *butterfly* comes from the pattern of inverted triangles created by the interconnections, which look like butterfly wings. A butterfly network transfers data packets using the most efficient route, but it is blocking, it cannot handle a conflict between two packets attempting to reach the same port at the same time. In a Clos network, all packets overshoot their destination and then hop back to it. Most of the time, the overshoot is not necessary and increases the latency, so a packet takes twice as many hops as it really needs.

In a *folded Clos* topology the input and the output networks share switch modules. Such networks are sometimes called *fat-tree*; many commercial high-performance interconnects, such as Myrinet, InfiniBand, and Quadrics, implement a fat-tree topology. Some folded Clos networks use low-radix routers, e.g., the Cray XD1 uses radix-24 routers. The latency and the cost of the network can be lowered using high-radix routers.

The *Black Widow* topology extends the folded Clos topology and has a lower cost and latency; it adds side links, and this permits a statical partitioning of the global bandwidth among peer subtrees [443]. The Black Widow topology is used in Cray computers.

Flattened butterfly network topology [272] is similar to the generalized hypercube proposed in early 1980s and it is able to exploit high-radix routers. In a flattened butterfly, we start with a conventional

**FIGURE 6.9**

(a) A five-stage Clos network with radix-2 routers and unidirectional channels equivalent to two back-to-back butterfly networks. (b) The corresponding folded-Clos network with bidirectional channels; input and output networks share switch modules.

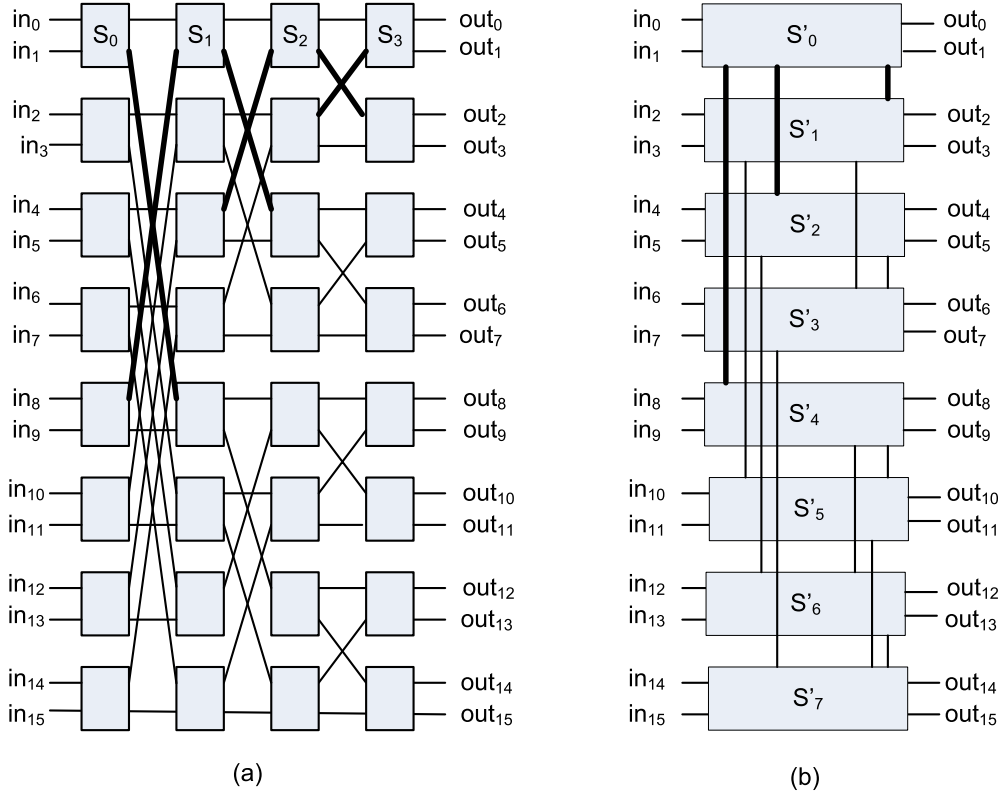
butterfly and combine the switches in each row into a single, higher-radix one. Each router is linked to more processors, and this halves the number of router-to-router connections.

Data sent by one processor can reach its destination with fewer hops and latency is reduced, though the physical path may be longer. For example, in Fig. 6.10(a), we see a two-ary four-fly butterfly; we combine the four switches S_0 , S_1 , S_2 and S_3 in the first row into a single switch, S'_0 . The flattened butterfly adaptively senses congestion and overshoots only when it needs to. On adversarial traffic patterns, the flattened butterfly has a similar performance as the folded Clos but provides over an order of magnitude increase in performance compared to the conventional butterfly.

The network cost for computer clusters can be reduced by a factor of two when high-radix routers (radix-64 or higher), and the flattened butterfly topology are used [273]. Flattened butterfly topology does not reduce the number of local cables, e.g., backplane wires from the processors to routers, but it reduces the number of global cables. The cost of the cables represents as much as 80% of the total network cost, e.g. for a 4K system, the cost savings of the flattened butterfly exceed 50%.

6.8 InfiniBand and Myrinet

InfiniBand is an interconnection network used by supercomputers and computer clouds backed by top companies in the industry, including Dell, HP, IBM, Intel, and Microsoft. Intel manufactures InfiniBand host bus adapters and network switches. InfiniBand uses a switched fabric rather than a shared communication channel and has a high throughput and low latency. Every link of the fabric has exactly one

**FIGURE 6.10**

(a) A two-ary four-fly butterfly with unidirectional links. (b) The corresponding two-ary four-flat flattened butterfly is obtained by combining the four switches S_0, S_1, S_2 and S_3 in the first row of the traditional butterfly into a single switch S'_0 and by adding additional connections between switches [272].

device connected at each end of the link, thus the worst case is the same as the typical case. A switched fabric architecture is fault tolerant and scalable, while in a shared channel architecture, all connected devices share the same bandwidth and the larger the number of devices, the less bandwidth is available to each one of them.

InfiniBand architecture implements the “bandwidth-out-of-the-box” concept and delivers bandwidth traditionally trapped inside a server across the interconnect fabric. InfiniBand supports several signaling rates, and the energy consumption depends on the throughput. Links can be bonded together for additional throughput, as shown in Table 6.2. InfiniBand specifications define multiple operational modes: single data rate (SDR), double data rate (DDR), quad data rate (QDR), 14 data rate (FDR), and enhanced data rate (EDR).

The signaling rates are: 2.5 Gbps in each direction per connection for an SDR connection; 5 Gbps for DDR; 10 Gbps for QDR; 14.0625 Gbps for FDR; 25.78125 Gbps for EDR per lane. The SDR, DDR,

Table 6.2 Evolution of high-speed interconnects. Several networking equipment suppliers offered proprietary solutions for 200 GE and 400 GE in 2016. Ethernet Alliance's 2020 technology roadmap expects 800 Gbps and 1.6 Tbps to become IEEE standard after 2020.

Network	Year	Speed
Gigabit Ethernet (GE)	1995	1 Gbps
10-GE	2001	10 Gbps
40-GE	2010	40 Gbps
100-GE	2010	100 Gbps
Myrinet	1993	1 Gbps
Fiber Channel	1994	1 Gbps
InfiniBand (IB)	2001	2 Gbps (1X SDR)
	2003	8 Gbps (4X SDR)
	2005	16 Gbps (4X DDR) & 24 Gbps (12X SDR)
	2007	32 Gbps (4X QDR)
	2011	56 Gbps (4X FDR)
	2012	100 Gbps (4X EDR)

and QDR link encoding is 8 B/10 B, every 10 bits sent carry 8 bits of data. Thus single, double, and quad data rates carry 2, 4, or 8 Gbit/s useful data, respectively, because the effective data transmission rate is four-fifths the raw rate.

InfiniBand allows links to be configured for a specified speed and width; the reactivation time of the link can vary from several nanoseconds to several microseconds. Exadata and Exalogic systems from Oracle implement the InfiniBand QDR with 40 Gbit/s (32 Gbit/s effective). InfiniBand fabric is used to connect compute nodes, compute nodes with storage servers, and Exadata and Exalogic systems.

In addition to high throughput and low latency, InfiniBand supports quality of service guarantees and failover—the capability to switch to a redundant or standby system. It offers point-to-point bidirectional serial links intended for the connection of processors with high-speed peripherals, such as disks, as well as multicast operations.

Initially, InfiniBand deployment pushed the limitations of personal computers (PCs) buses. Introduced as a standard PC architecture in the early 90s, the PCI bus has evolved from 32 bit/33 MHz to 64 bit/66 MHz, while PCI-X has doubled the clock rate to 133 MHz. PCI Express 2.0/x16 has a 16 Gbps bandwidth. The number of pins of a parallel bus is quite large; the number of pins per connection is large, e.g., a 64 bit PCI bus requires 90 pins.

Myrinet is an interconnect for massively parallel systems developed at Caltech. It is no longer in use, though it has several remarkable features [65]:

- a. Robustness ensured by communication channels with flow control, packet framing, and error control.
- b. Self-initializing, low-latency, cut-through switches.
- c. Host interfaces that can map the network, select routes, and translate from network addresses to routes, as well as handle packet traffic.

- d. Streamlined host software that allows direct communication between user processes and the network.

Myrinet design benefits from the experience gathered by a project to construct a high-speed local-area network at USC/ISI. A Myrinet is composed of point-to-point, full-duplex links connecting hosts and switches, an architecture similar to the one of the ATOMIC project at USC. Multiple-port switches may be connected to other switches and to the single-port host interfaces in any topology. Transmission is synchronous at the sending end, while the receiver circuits are asynchronous. The receiver injects control symbols in the reverse channel of the link for flow control. Myrinet switches use blocking-cut-through routing similar to the one in Intel Paragon and Cray T3D.

Myrinet supports high data rates; a Myrinet link consists of a full-duplex pair of 640 Mbps channels and has regular topology with elementary routing circuits in each node. The network is scalable, and its aggregate capacity grows with the number of nodes. Algorithmic routing enables multiple packets to flow concurrently and avoids deadlocks.

Store and forward and *cut-through or wormhole* networks are very different. In the former, an entire packet is buffered, and its checksum is verified in each node along the path from the source to the destination. In wormhole networks the packet is forwarded to its next hop as soon as the header is received and decoded. This decreases the latency, but a packet can still experience blocking if the outgoing channel expected to carry it to the next node is in use. In this case the packet has to wait until the channel becomes free.

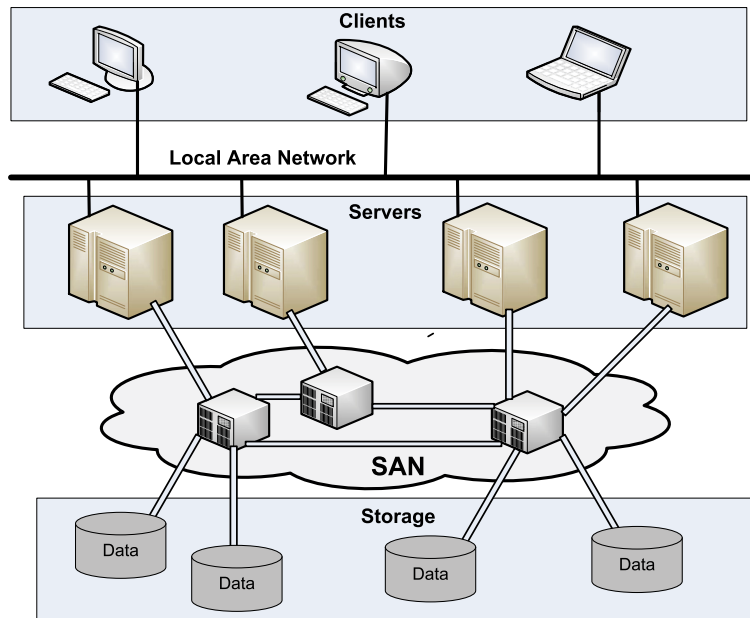
A comparison of Myrinet and Ethernet performance as a communication substrate for MPI libraries in [323] shows that the MG NAS benchmark,⁵ over Myrinet, only achieves 5% higher performance than TCP over Ethernet. MPI library implementations for Ethernet have a higher message latency and lower message bandwidth because they use the OS network protocol stack.

6.9 Storage area networks and the Fibre Channel

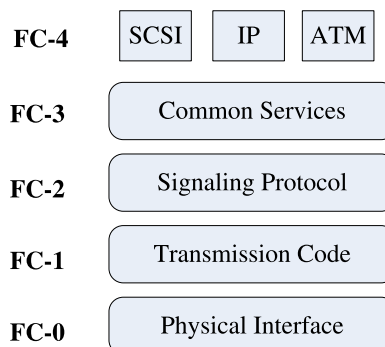
A storage area network (SAN) is a specialized, high-speed network for data block transfers. SANs interconnect servers to servers, servers to storage devices, and storage devices to storage devices; see Fig. 6.11. A SAN consists of a communication infrastructure and a management layer. The Fibre Channel (FC) is the dominant architecture of SANs. FC is a layered protocol with several layers, as depicted in Fig. 6.12, discussed next.

- A. Lower layer protocols: FC-0 physical interface; FC-1 transmission protocol responsible for encoding/decoding; and FC-2 signaling protocol responsible for framing and flow control. FC-0 uses laser diodes as the optical source and manages point-to-point fiber connections; when a fiber connection is broken, the ports send a series of pulses until the physical connection is reestablished and the necessary handshake procedures are followed. FC-1 controls the serial transmission and integrates

⁵ NASA Parallel Benchmarks, NAS, are used to evaluate the performance of parallel supercomputers. The original benchmark included five kernels: IS—Integer Sort, random memory access; EP—Embarrassingly Parallel; CG—Conjugate Gradient; MG—Multi-Grid on a sequence of meshes, long- and short-distance communication, memory intensive; and FT—discrete 3D Fast Fourier Transform, all-to-all communication.

**FIGURE 6.11**

A storage area network interconnects servers to servers, servers to storage devices, and storage devices to storage devices. Typically, it uses fiber optics and the FC protocol.

**FIGURE 6.12**

FC protocol layers. FC-4 supports communication with Small Computer System Interface (SCSI), IP, and Asynchronous Transfer Mode (ATM) network interfaces.

data with clock information. It ensures encoding to the maximum length of the code, maintains DC-balance, and provides word alignment. FC-2 provides the transport methods for data transmitted in 4-byte ordered sets containing data and control characters. It handles the topologies, based on the presence or absence of a fabric, the communication models, the classes of service provided by the fabric and the nodes, sequence and exchange identifiers, and segmentation and reassembly.

B. Upper layer protocols: FC-3 common services layer; and FC-4 protocol mapping layer. FC-3 supports multiple ports on a single-node or fabric using:

- Hunt groups—sets of associated ports assigned an alias identifier that allows any frame containing that alias to be routed to any available port within the set.
- Striping to multiply bandwidth, using multiple ports in parallel to transmit a single information unit across multiple links.
- Multicast and broadcast to deliver a single transmission to multiple destination ports or to all ports.

FC supports several classes of service to accommodate various application needs:

Class 1: rarely used blocking connection-oriented service; acknowledgments ensure that the frames are received in the same order in which they are sent, and reserves full bandwidth for the connection between the two devices.

Class 2: acknowledgments ensure that the frames are received; enables the fabric to multiplex several messages on a frame-by-frame basis; since frames can take different routes, it does not guarantee in-order delivery, rather it relies on upper layer protocols to take care of frame sequence.

Class 3: datagram connection; no acknowledgments.

Class 4: connection-oriented service. Virtual circuits (VCs) established between ports, guarantees in-order delivery and acknowledgment of delivered frames. The fabric is responsible for multiplexing frames of different VCs. Supports Guaranteed Quality of Service (QoS), including bandwidth and latency. This layer is intended for multimedia applications.

Class 5: isochronous service for applications requiring immediate delivery, without buffering.

Class 6: supports dedicated connections for a reliable multicast.

Class 7: similar to Class 2 but used for the control and management of the fabric; a connectionless service with notification of non-delivery.

While every device connected to a LAN has a unique MAC address,⁶ each FC device has a unique Id called the WWN (World Wide Name), a 64 bit address. Every port in the switched fabric has its own unique 24-bit address consisting of: the domain (bits 23–16), the area (bits 15–08), and the port physical address (bits 07–00).

The switch of a switched fabric environment assigns dynamically and maintains the port addresses. When a device with a WWN address logs into a given port, the switch maintains the correlation between that port address and the WWN address of the device using the Name Server. The Name Server is a component of the fabric operating system, which runs inside the switch.

⁶ Typically, a virtualization platform generates a new, random MAC address for each virtual network interface at creation. Moreover, in Linux, it is possible to spoof the MAC address.

Word 0 4 bytes SOF (Start Of Frame)	Word 1 3 bytes Destination port address	Word 2 3 bytes Source port address	Word 3-6 18 bytes Control information	(0-2112 bytes) Payload	CRC	EOF (End Of Frame)
---	---	--	--	---------------------------	-----	--------------------------

FIGURE 6.13

The format of an FC frame; the payload can be at most 2112 bytes; and larger data units are carried by multiple frames.

The format of an FC frame is shown in Fig. 6.13. Zoning permits finer segmentation of the switched fabric; only the members of the same zone can communicate within that zone. It can be used to separate different environments, e.g., a Microsoft Windows NT from a UNIX environment.

SANs use several other protocols, for example, Fibre Channel over IP (FCIP) allows transmission of Fibre Channel data through IP networks using tunneling. *Tunneling* allows a network protocol to carry a payload over an incompatible delivery-network, or to provide a secure path through an untrusted network. A protocol normally blocked by a firewall is wrapped inside a protocol that the firewall does not block. For example, an HTTP tunnel can be used for communication from network locations with restricted connectivity, e.g., behind NATs, firewalls, or proxy servers. Restricted connectivity is a commonly used method to lock down a network to secure it against internal and external threats.

Internet Fibre Channel Protocol (iFCP) is a gateway-to-gateway protocol supporting communication among FC storage devices in a SAN, or on the Internet, using TCP/IP; iFCP replaces lower-layer Fibre Channel transport with TCP/IP and Gigabit Ethernet. Fibre Channel devices connect to an iFCP gateway or switch, and each Fibre Channel session is terminated at the local gateway and converted to a TCP/IP session via iFCP.

6.10 Scalable data center communication architectures

The question addressed now is: How does one organize the communication infrastructure of large cloud data centers for best performance at the lowest cost? *Data center networks* (DCNs) architectures providing an answer to this question face major challenges: (i) aggregate cluster bandwidth scales poorly with the cluster size; (ii) the bandwidth needed by many cloud applications comes at a high price, and the interconnect cost increases dramatically with cluster size; and (iii) the communication bandwidth of DCNs may become oversubscribed by a significant factor depending on communication patterns.

We only mention two DCN architectural styles, *three-tier* and *fat-tree* DCNs. The former has a multiple-rooted tree topology with three layers: core, aggregate, and access. Servers are connected to switches at tree leaves at the *access layer*. Enterprise switches at the root of the tree form the *core layer* and connect switches together at the *aggregate layer* connecting the data center to the Internet. The uplinks of the *aggregate layer* switches connect them to the core layer, and their download links connect to the access layer. Three-tier DCN architecture is not suitable for computer clouds, it is not scalable, the bisection bandwidth is not optimal, and core layer switches are expensive and power-hungry.

Fat-tree topology is optimal for computer clouds, the bandwidth is not severely affected for messages crossing multiple switches, and the interconnection network can be built with commodity rather than enterprise switches. The implementation of the fat-tree topology proposed in [20] is based on the following principles:

- i. The network should scale to a very large number of nodes.
- ii. The fat-tree should have multiple core switches.
- iii. The network should support multipath routing. Equal-cost multi-path (ECMP) routing algorithm [244] performing static load splitting among flows should be used.
- iv. The network should use switches with optimal cost/performance ratios.

Hierarchical and fat-tree interconnection networks cost per GigE⁷ has decreased by one order of magnitude and the cost/performance indicator for fat-tree built with commodity switches is almost an order of magnitude lower than that of hierarchical networks. The choice of multirooted fat-tree topology and multipath routing is justified because in 2008 the largest cluster that could be supported with a single rooted core 128-port router with 1:1 oversubscription would be limited to 1 280 nodes.

A WSC interconnect can be organized as a fat-tree with k -port switches, $k = 48$, but the organization can be supported for any k . The network consists of k pods⁸ each pod has two layers and $k/2$ switches at each layer. Each switch at the lower layer is connected directly to $k/2$ servers. The other $k/2$ ports are connected to $k/2$ of k ports in the aggregation layer. The total number of switches is $k(k + 1)$ and the total number of servers connected to the system is k^2 ; $(k/2)^2$ paths connect each pair of servers.

A WSC with 16 384 servers can be built with 128-port switches, and one with 262 144 servers requires 512-port switches. Fig. 6.14 shows a fat-tree interconnection network for $k = 4$. The core, the aggregation, and the edge layers are populated with *four*-port switches. Each core switch is connected to one of the switches at the aggregation layer of each pod. The network has four pods, four switches at each pod, two at aggregation layer and two at the edge. Four servers are connected to each pod.

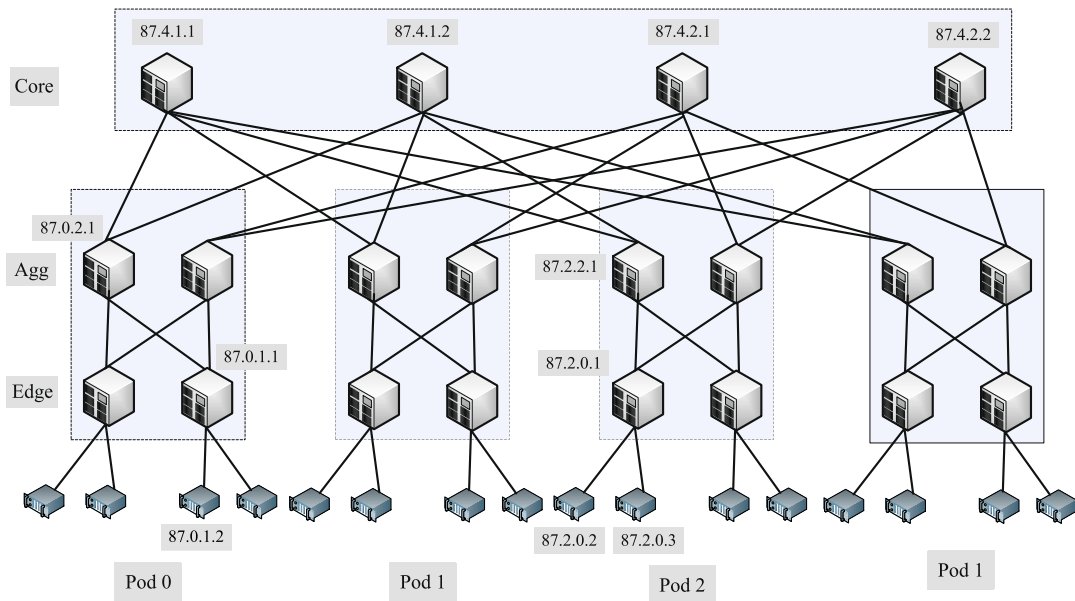
Switch IP addresses are $87.pod.switch.i$; switches are numbered left to right, and bottom to top. Core switches IP addresses are $10.k.j.i$, where j and i denote switch coordinates in the $(k/2)^2$ core switch grid starting from top-left. For example, the four switches of pod 2 have IP addresses 87.2.0.1, 87.2.1.1, 87.2.2.1, and 87.2.3.1. Server IP addresses are $87.pod.switch.serverID$, where *serverID* is server position in the subnet of the edge router starting from left to right. For example, the IP addresses of the two servers connected to the switch with IP address 87.2.0.1 are 87.2.0.2 and 87.2.0.3.

There are multiple path between any pair of servers. For example, packets sent by server with IP address 87.0.1.2 to server with IP address 87.2.0.3 can follow the following routes:

$$\begin{aligned}
 &87.0.1.1 \mapsto 87.0.2.1 \mapsto 87.4.1.1 \mapsto 87.2.2.1 \mapsto 87.2.0.1 \\
 &87.0.1.1 \mapsto 87.0.2.1 \mapsto 87.4.1.2 \mapsto 87.2.2.1 \mapsto 87.2.0.1 \\
 &87.0.1.1 \mapsto 87.0.1.1 \mapsto 87.4.2.1 \mapsto 87.2.2.1 \mapsto 87.2.0.1 \\
 &87.0.1.1 \mapsto 87.0.1.1 \mapsto 87.4.2.2 \mapsto 87.2.2.2 \mapsto 87.2.0.1
 \end{aligned} \tag{6.5}$$

⁷ The IEEE 802.3-2008 standard defines GigE as a technology for transmitting Ethernet frames. One GigE corresponds to a rate of one gigabit per second, i.e., 10^9 bits per second.

⁸ A pod is a repeatable design pattern to maximize the modularity, scalability, and manageability of data center networks.

**FIGURE 6.14**

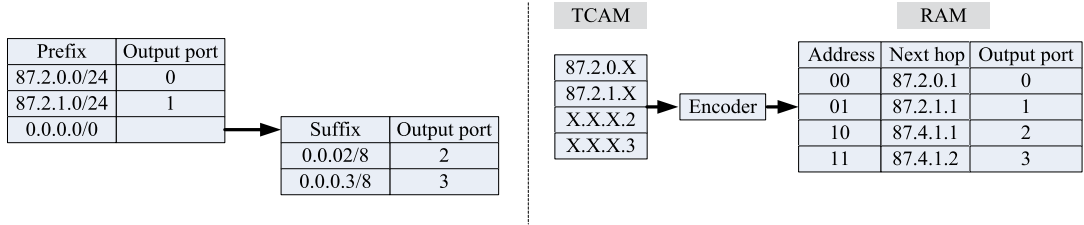
A fat-tree network with k -port switches and $k = 4$. Four core four-way switches are at the root, and there are four pods, two switches at the aggregation layer, and two at the edge layer of each pod. Each switch at the edge of a pod is connected to two servers.

Packet routing uses two-level routing tables and supports two-level prefix lookup. This strategy may increase the lookup latency, but the prefix search can be done in parallel and compensate for the latency increase. The main table entries are of the form $(prefix, outputport)$ and could have an additional pointer to a secondary table or could be terminating if none of its entries point to secondary table. A secondary table consists of $(suffix, outputport)$ entries and may be pointed to by more than one first-level entry.

Fig. 6.15 (Left) shows the two-level routing tables for switch 87.2.2.1 and routing for two incoming packets for servers with the IP addresses 87.2.1.2 and 87.3.0.3; the incoming packets are forwarded on ports 1 and 3, respectively. Lookup engines use a ternary version of content-addressable memory (CAM), called TCAM. Fig. 6.15 (Right) shows that TCAM stores address prefixes and suffixes, which index a RAM that stores the IP address of the next hop and the output port.

The prefix entries are stored with numerically smaller addresses first and the right-handed (suffix) entries in larger addresses. The output of the CAM is encoded so that the entry with the numerically smallest matching address is the output. When the destination IP address of a packet matches both a left-handed and a right-handed entry, then the left-handed entry is chosen.

The k switches in a pod have terminating prefixes to the subnets in that pod. When two servers in the same pod, but on a different subnet communicate, all upper-level switches of the pod will have a terminating prefix pointing to the destination subnet's switch.

**FIGURE 6.15**

(Left) The two-level routing tables for switch 87.2.2.1. Two incoming packets for IP addresses 87.2.1.2 and 87.3.0.3 are forwarded on ports 1 and 3, respectively. (Right) The RAM implementation of a two-level TCAM routing table.

For all outgoing pod traffic, the pod switches have a default /0 prefix with a secondary table matching the least-significant byte of the destination IP address, the server ID. Traffic diffusion occurs only in the first half of a packet's journey. Once a packet reaches a core switch, there is exactly one link to its destination pod, and that switch will include a terminating /16 prefix for the pod of that packet (87.pod.0.0/16, *port*). Once a packet reaches its destination pod, the receiving upper-level pod switch will also include a (10.pod.switch.0/24, *port*) prefix to direct the packet to its destination subnet switch, where it is finally switched to its destination server.

Each pod switch assigns terminating prefixes for subnets in the same pod and adds a /0 prefix with a secondary table matching the *serverIDs* for inter-pod traffic. The routing tables for the upper pod switches are generated with the pseudocode of Algorithm 6.1.

Algorithm 6.1: Generates aggregation switch routing tables

```

1 foreach pod  $x \in [0, k - 1]$  do
2   foreach switch  $z \in [k/2, k - 1]$  do
3     foreach subnet  $i \in [0, k/2 - 1]$  do
4       addPrefix(10.x.z.1, 10.x.i.0/24, i);
5     end
6     addPrefix(10.x.z.1, 0.0.0.0/0, 0);
7     foreach hostID  $i \in [2, (k/2) + 1]$  do
8       addSuffix(10.x.z.1, 0.0.0.i/8, ( $i - 2 + z$ ) mod ( $k/2$ ) + ( $k/2$ ));
9     end
10   end
11 end

```

For the lower pod switches, the /24 subnet prefix in line 3 is omitted since that subnet's own traffic is switched, and intra- and interpod traffic should be evenly split among the upper switches. Core switches contain only terminating /16 prefixes pointing to their destination pods, as shown in Algorithm 6.2. The maximum numbers of first-level prefixes and second-level suffixes are k and $k/2$, respectively.

Flow classification with dynamic port-reassignment in pod switches overcomes cases of local congestion when two flows compete for the same output port, even though another port that has the same

Algorithm 6.2: Generates routing tables for core switches

```

1 foreach  $j \in [0, k - 1]$  do
2     foreach  $i \in [1, k/2]$  do
3         foreach destination pod  $\in [0, k/2 - 1]$  do
4             addPrefix(10.k.j.i, 10.x.0.0/16, x);
5         end
6     end
7 end

```

cost to the destination is underused. Power consumption and heat dissipation are major concerns for the cloud data centers. Switches at the higher tiers of the interconnect in data centers consume several kW, and the entire interconnection infrastructure consumes hundreds to thousands of kW.

6.11 Network resource management algorithms (R)

A critical aspect of resource management in cloud computing is to guarantee the communication bandwidth required by an application as specified by an SLA. The solutions to this problem are based on the strategies used for some time on the Internet to support the data streaming QoS requirements.

Cloud interconnects consist of communication links of limited bandwidth and switches of limited capacity. When the load exceeds its capacity, a switch starts dropping packets because it has limited input buffers for the switching fabric and for the outgoing links, as well as limited CPU cycles. A switch must handle multiple flows, pairs of source-destination endpoints of the traffic, thus a scheduling algorithm has to manage several quantities at the same time: the *bandwidth*, the amount of data each flow is allowed to transport, the *timing* when the packets of individual flows are transmitted, and the *buffer space* allocated to each flow.

Communication and computing require scheduling, therefore it should be no surprise that the first algorithm we discuss can be used for scheduling packets transmission, as well as threads. A first strategy to avoid network congestion is to use a FCFS scheduling algorithm. The advantage of FCFS algorithm is a simple management of the three quantities: bandwidth, timing, and buffer space. Nevertheless, the FCFS algorithm does not guarantee fairness; greedy flow sources can transmit at a higher rate and benefit from a larger share of the bandwidth.

Fair Queuing (FQ). The algorithm ensures that a high-data-rate flow cannot use more than its fair share of the link capacity. Packets are first classified into flows by the system and then assigned to a queue dedicated to the flow. Packet queues are serviced one packet at a time in round-robin (RR) order as depicted in Fig. 6.16. FQ's objective is *max-min* fairness. This means that it maximizes first the minimum data rate and then the second minimum data rate of any data flow. Starvation of expensive flows is avoided, but the throughput is low.

The FQ algorithm guarantees buffer space management fairness but does not guarantee fairness of bandwidth allocation; a flow transporting large packets will benefit from a larger bandwidth [360].

The FQ algorithm in [138] proposes a solution to this problem. First, it introduces a *Bit-by-bit Round-robin (BR)* strategy; as the name implies, in this rather impractical scheme, a single bit from

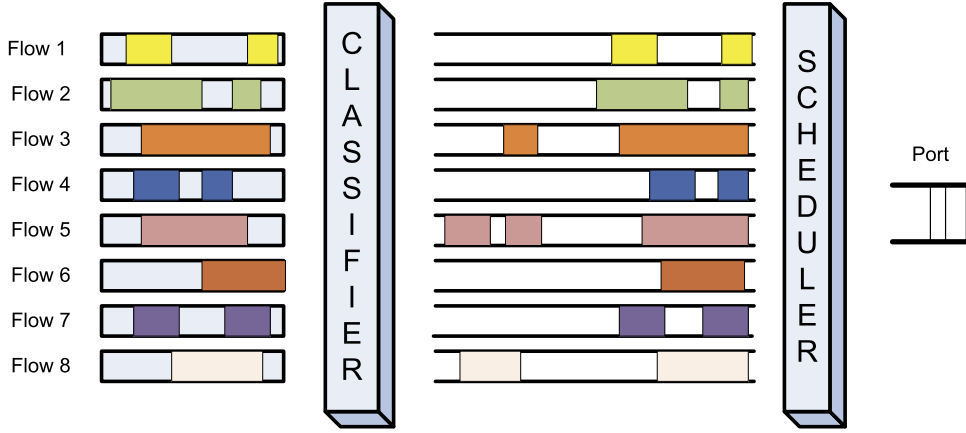


FIGURE 6.16

Fair Queuing—packets are first classified into flows by the system and then assigned to a queue dedicated to the flow; queues are serviced one packet at a time in round-robin order and empty queues are skipped.

each queue is transmitted, and the queues are visited in a round-robin fashion. Let $R(t)$ be the number of rounds of the BR algorithm up to time t and $N_{active}(t)$ be the number of active flows through the switch. Call t_i^a the time when the packet i of flow a , of size P_i^a bits arrives, and call S_i^a and F_i^a the values of $R(t)$ when the first and the last bit, respectively, of the packet i of flow a are transmitted. Then,

$$F_i^a = S_i^a + P_i^a \quad \text{and} \quad S_i^a = \max[F_{i-1}^a, R(t_i^a)]. \quad (6.6)$$

$R(t)$, $N_{active}(t)$, S_i^a , and F_i^a in Fig. 6.17 depend only on the arrival time of the packets, t_i^a , and not on their transmission time, provided that a flow a is active as long as

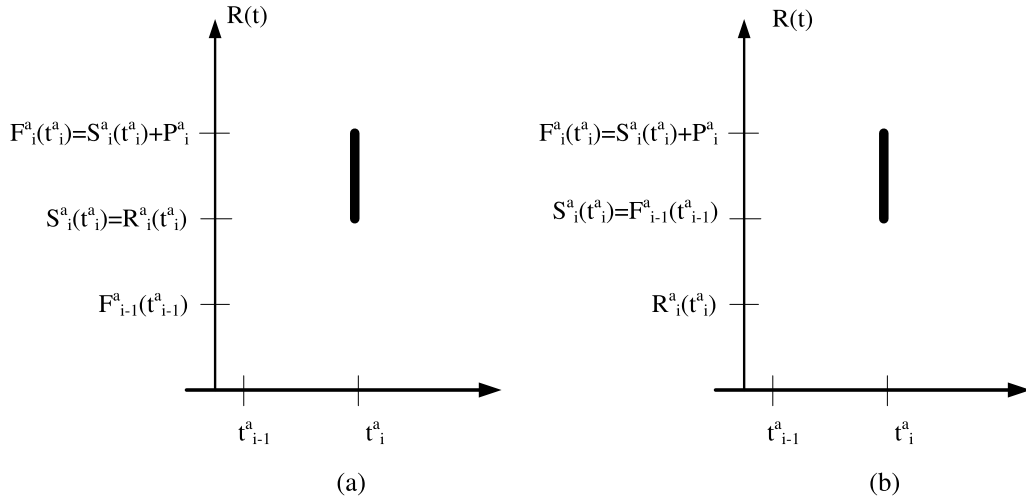
$$R(t) \leq F_i^a \quad \text{when} \quad i = \max(j | t_j^a \leq t). \quad (6.7)$$

The authors of [138] use for the packet-by-packet transmission time the following nonpreemptive scheduling rule which emulates the BR strategy: *the next packet to be transmitted is the one with the smallest F_i^a* . A preemptive version of the algorithm requires that the transmission of the current packet be interrupted as soon as one with a shorter finishing time, F_i^a , arrives.

A fair allocation of the bandwidth does not have an effect on the timing of the transmission. A possible strategy is to allow less delay for the flows using less than their fair share of the bandwidth. The same paper [138] proposes the introduction of a quantity called the *bid*, B_i^a , and scheduling the packet transmission based on its value. The bid is defined as

$$B_i^a = P_i^a + \max[F_{i-1}^a, (R(t_i^a) - \delta)], \quad (6.8)$$

with δ being a nonnegative parameter. The properties of the FQ algorithm, as well as the implementation of a nonpreemptive version of the algorithm, are analyzed in [138].

**FIGURE 6.17**

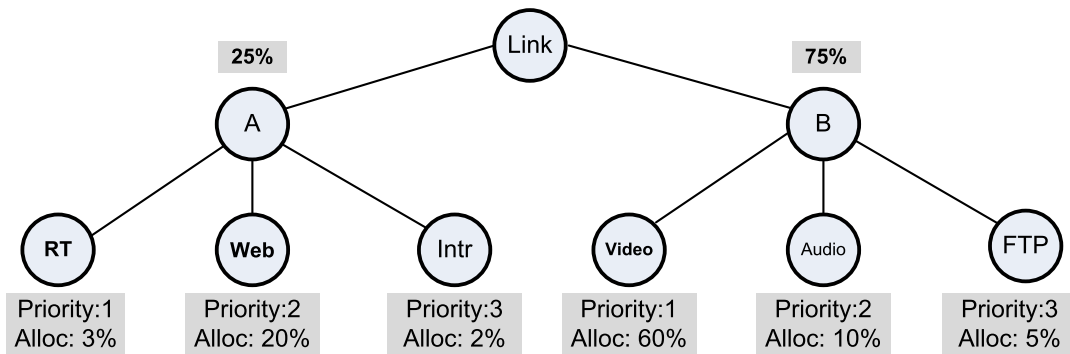
Transmission of packet i of flow a arriving at time t_i^a of size P_i^a bits. Transmission starts at time $S_i^a = \max[F_{i-1}^a, R(t_i^a)]$ and ends at time $F_i^a = S_i^a + P_i^a$ with $R(t)$ the number of rounds of the algorithm. (a) The case $F_{i-1}^a < R(t_i^a)$. (b) The case $F_{i-1}^a \geq R(t_i^a)$.

Stochastic Fairness Queuing (SFQ) requires less calculations; it is a simpler and less accurate implementation of FQ algorithms. SFQ ensures that each flow has the opportunity to transmit an equal amount of data and takes into account data packet sizes [343].

Class-Based Queuing (CBQ). This algorithm is a widely used strategy for link sharing proposed by Sally Floyd and Van Jacobson in 1995 [181]. The objective of CBQ is to support flexible link sharing for applications that require bandwidth guarantees such as VoIP, video streaming, and audio streaming. At the same time, CBQ supports some balance between short-lived network flows, such as web searches, and long-lived ones, such as video streaming or file transfers.

CBQ aggregates connections and constructs a hierarchy of classes with different priorities and throughput allocations. CBQ, uses several functional units for link sharing: (i) a *classifier* uses the information in the packet header to assign arriving packets to classes; (ii) an *estimator* of the short-term bandwidth for the class; (iii) a *selector*, or scheduler, identifies the highest priority class to send next and, if multiple classes have the same priority, to schedule them on a round-robin base; and (iv) a *delayer* to compute the next time when a class that has exceeded its link allocation is allowed to send.

Classes are organized in a tree-like hierarchy; for example, in Fig. 6.18, group A corresponds to short-lived traffic and group B to long-lived traffic. The leaves of the tree are considered Level 1, and this example includes six classes of traffic: real-time, web, interactive, video streaming, audio streaming, and file transfer. At Level 2, there are the two classes of traffic, A and B . The root at Level 3 is the link itself. The link-sharing policy aims to ensure that, if sufficient demand exists, then, after some time intervals, each interior or leaf class receives its allocated bandwidth. The distribution of the “excess” bandwidth follows a set of guidelines but does not support mechanisms for congestion avoidance.

**FIGURE 6.18**

CBQ link sharing for two groups A, of short-lived traffic, and B, of long-lived traffic, allocated 25% and 75% of the link capacity, respectively. There are six classes of traffic with priorities 1, 2, and 3. The RT (real-time) and the video streaming have priority 1 and are allocated 3% and 60%, respectively, of the link capacity. Web transactions and audio streaming have priority 2 and are allocated 20% and 10%, respectively, of the link capacity. Intr (interactive applications) and FTP (file transfer protocols) have priority 3 and are allocated 2% and 5%, respectively, of the link capacity.

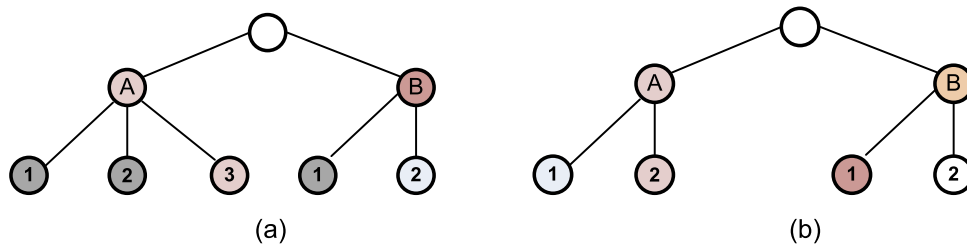
A class is *overlimit* if, over a certain recent period, it has used more than its bandwidth allocation (in bytes per second), *underlimit* if it has used less, and *atlimit* if it has used exactly its allocation. A leaf class is *satisfied* if it is underlimit, has a persistent backlog, and is *unsatisfied* otherwise; a non-leaf class is unsatisfied if it is underlimit and has some descendent class with a persistent backlog. A precise definition of the term “persistent backlog” is part of a local policy. A class does not need to be *regulated* if it is underlimit or if there are no unsatisfied classes; the class should be regulated if it is overlimit and if some other class is unsatisfied, and this regulation should continue until the class is no longer overlimit or until there are no unsatisfied classes; see Fig. 6.19 for two examples.

Linux kernel implements a link sharing algorithm called *Hierarchical Token Buckets* (HTB) inspired by CBQ. In CBQ every class has an *assured rate* (AR); in addition to AR every class in HTB has also a *ceil rate* (CR)⁹; see Fig. 6.20. The main advantage of HTB over CBQ is that it allows *borrowing*. If a class *C* needs a rate above its AR, it tries to borrow from its parent; then the parent examines its children and, if there are classes running at a rate lower than their AR, the parent can borrow from them and reallocate it to class *C*.

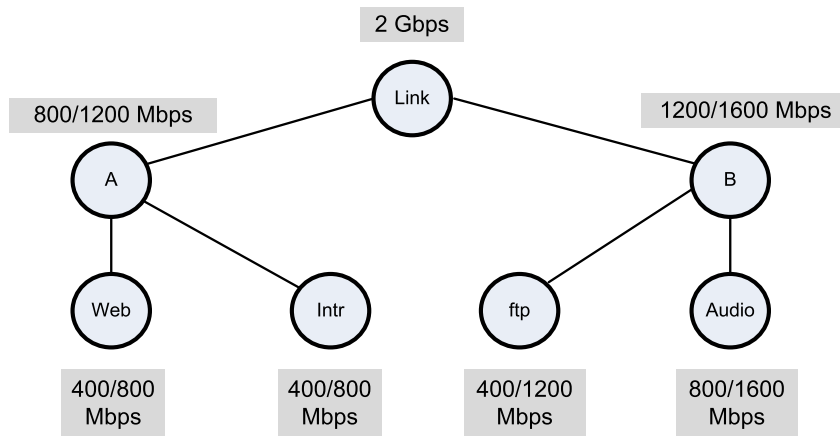
6.12 Content delivery networks

Computer clouds support network-centric computing and network-centric content. The vast amount of data stored on the cloud has to be delivered efficiently to a large user population and Content Delivery

⁹ The ceil rate is the rate a class is allowed to borrow from other classes.

**FIGURE 6.19**

There are two groups *A* and *B* and three types of traffic, e.g., web, real-time, and interactive, denoted as 1, 2, and 3. (a) Group *A* and class *A.3* traffic are underlimit and unsatisfied; classes *A.1*, *A.2*, and *B.1* are overlimit, unsatisfied and with persistent backlog, and have to be regulated; type *A.3* is underlimit and unsatisfied; group *B* is overlimit. (b) Group *A* is underlimit and unsatisfied; Group *B* is overlimit and needs to be regulated; class *A.1* traffic is underlimit; class *A.2* is overlimit and with persistent backlog; class *B.1* traffic is overlimit and with persistent backlog and needs to be regulated.

**FIGURE 6.20**

HTB packet scheduling adds a ceil rate to the allowed rate of every node.

Networks (CDNs) offer fast and reliable content delivery and reduce communication bandwidth by caching and replication. A CDN receives the content from an *Origin* server, then replicates it to its *Edge* cache servers; the content is delivered to an end-user from the “closest” Edge server.

CDNs support scalability, increase reliability and performance, and provide better security. The volume of transactions and data transported by the Internet increases dramatically every year; additional resources are necessary to accommodate the additional load placed on the communication and storage systems and to improve the end-user experience. CDNs place additional resources pro-

visioned to absorb the traffic caused by *flash crowds*¹⁰ and, in general, to provide capacity on demand.

The additional resources are placed strategically throughout the Internet to ensure scalability. The resources provided by a CDN are replicated, and when one replica fails, the content is available from another one; replicas are “close” to the consumers of the content, and this placement reduces the start-up time and the communication bandwidth. A CDN uses two types of servers: the *origin* server updated by the content provider and *replica* servers, which cache the content and serve as authoritative reference for client requests. Security is a critical aspect of the services provided by a CDN; the replicated content should be protected from the increased risk of cyber fraud and unauthorized access.

CDNs deliver static content and/or live or on-demand streaming media. *Static content* refers to media stored using traditional caching technologies because changes are infrequent. Examples of static content are: HTML pages, images, documents, software patches, and audio and/or video files. *Live media* refers to live events when the content is delivered in real time from the encoder to the media server. On-demand delivery of audio and/or video streams, movie files, and music clips provided to the end-users is content-encoded and then stored on media servers. Virtually all CDN providers support static content delivery, while live or on-demand streaming media is considerably more challenging.

CDN providers and protocols. The first CDN was setup by *Akamai*, a company that evolved from an MIT project to optimize network traffic. Since its inception Akamai has placed some 20 000 servers in 1 000 networks in 71 countries; in 2009, it controlled some 85% of the market [391]. Akamai mirrors the contents of clients on multiple systems placed strategically through the Internet. Though the domain name (but not subdomain) is the same, the IP address of the resource requested by a user points to an Akamai server rather than the customer’s server. Then the Akamai server is automatically picked depending on the type of content and the network location of the end user.

There are several other active commercial CDNs including EdgeStream providing video streaming and Limelight Networks providing on-demand and live delivery of video, music, and games. There are several academic CDNs: Coral is a freely available network designed to mirror web content, hosted on PlanetLab; Globule is an open-source collaborative CDN developed at Vrije Universiteit in Amsterdam.

The communication infrastructure among different CDN components uses a fair number of protocols including: Network Element Control Protocol (NECP), Web Cache Coordination Protocol (WCCP), SOCKS, Cache Array Routing Protocol (CARP), Internet Cache Protocol (ICP), Hypertext Caching Protocol (HTCP), and Cache Digest described succinctly in [391]. For example, caches exchange ICP queries and reply to locate the best sites to retrieve an object; HTCP is used for discovering HTTP caches, caching data, managing sets of HTTP caches, and monitoring cache activity.

CDN organization, design decisions, and performance. There are two strategies for CDN organization; in the so-called *overlay*, the network core does not play an active role in the content delivery. On the other hand, the *network* approach requires the routers and the switches to use dedicated software to identify specific application types and to forward user’s requests based on predefined policies.

The first strategy is based exclusively on content replication on multiple caches and redirection based on proximity to the end-user. In the second approach the network core elements redirect con-

¹⁰ The term flash crowd refers to an event disrupting the life of a significant segment of the population, such as an earthquake in a populated area; it causes the Internet traffic to increase dramatically.

tent requests to local caches or redirect data center's incoming traffic to servers optimized for specific content type access. Some CDNs including Akamai use both strategies.

The most important design and policy decisions for a CDN are: (i) the placement of the edge servers; (ii) the content selection and delivery; (iii) content management; and (iv) request routing policies. The placement problem is often solved with suboptimal heuristics using as input the workload patterns and the network topology. The simplest, but costly, approach for content selection and delivery is the *full-site* replication suitable for static content; the edge servers replicate the entire content of the origin server. On the other hand, the *partial-site* selection and delivery retrieves the base HTML page from the origin server and the objects referenced by this page from the edge caches. The objects can be replicated based on their popularity, or on some heuristics.

Content management depends on the caching techniques, the cache maintenance, and the cache update policies. CDNs use several strategies to manage the consistency of content at replicas: periodic updates, updates triggered by the content change, and on-demand updates. Request-routing in a CDN directs users to the closest edge server that can best serve the request; metrics, such as network proximity, client perceived latency, distance, and replica server load, are taken into account when routing a request. Round-robin is a nonadaptive request-routing that aims to balance the load; it assumes that all edge servers have similar characteristics and can deliver the content.

Adaptive algorithms perform considerably better but are more complex and require some knowledge of state of the system. The algorithm used by Akamai takes into consideration metrics such as: the load of the edge server, the bandwidth currently available to a replica server, and the reliability of client connection to edge servers. CDN routing can exploit an organization where several edge servers are connected to a service node aware of the load and the information about each edge server connected to it and attempts to implement a *global load balancing policy*.

An alternative is *DNS-based routing*, when a domain name has multiple IP addresses associated to it and the service provider's DNS server returns the IP addresses of the edge servers holding the replica of the requested object; then, the client's DNS server chooses one of them. Another alternative is the *HTTP redirection*; in this case a web server includes in the HTTP header of a response to a client the address of another edge server. Finally, *IP anycasting* requires that the same IP address is assigned to several hosts and the routing table of a router contains the address of the host closest to it.

Critical CDN performance metrics are: (i) cache hit ratio—the ratio of the number of cached objects versus total number of objects requested; (ii) reserved bandwidth for the origin server; (iii) latency—based on the perceived response time by the end users; (iv) edge server utilization; and, last but not least, (v), reliability—based on packet loss measurements.

CDNs face considerable challenges due to increased appeal of data streaming and to the proliferation of mobile devices such as smartphones and tablets. On-demand video streaming requires enormous bandwidth and storage space, as well as powerful servers; CDNs for mobile networks must be able to dynamically reconfigure the system in response to spatial and temporal demand variations.

Content-Centric Networks (CCNs) are related to information-centric networking architectures, such as Named Data Networks (NDNs) discussed in Section 6.4, where content is named and transferred throughout the network. The request for a named content is routed to the producer or to any entity that can deliver the expected content object.

CCN content may be served from any router's cache. Content is signed by its producer and the consumer is able to verify the signature before actually using the content. CCNs supports opportunistic caching. CCNs offer a number of advantages according to <http://chris-wood.github.io/2015/06/>

Table 6.3 VANETs applications, contents, local interest, local validity, and lifetime.

Application	Contents	Local interest	Local validity	Lifetime
Safety warnings	Dangerous Road	All	100 m	10 s
Safety warnings	Accident	All	500 m	30 s
Safety warnings	Work zone	All	1 Km	Construction
Public service	Emergency vehicle	All	500 m	10 min
Public service	Highway information	All	5 Km	All day
Driving	Road congestion	All	5 Km	30 min
Driving	Navigation Map	Subscribers	5 Km	30 min

[16/CCN-vs-CDN.html](#) i.e., content popularity need not be predicted beforehand. CCNs offer benefits not available from by IP-based CDNs: (i) active and intelligent forwarding strategies for routers; (ii) publisher mobility is easily supported via CCN routing protocols; (iii) congestion control can be enforced within the network; (iv) existing/problematic IP stack can be completely replaced with a new set of layers; (v) existing APIs can be completely reworked to focus on content, not on addresses; and (vi) content security is not tied to the channel but to the content itself.

Content service networks (CSN) introduced in [320] are overlay networks built around CDNs to provide an infrastructure service for processing and transcoding.

6.13 Vehicular ad hoc networks

A vehicular ad hoc network (VANET) consists of groups of moving or stationary vehicles connected by a wireless network. Until recently, the main use of VANETs was to provide safety and comfort to drivers in vehicular environments. This view is changing: Vehicular ad hoc networks are seen now as an infrastructure for an intelligent transportation system with an increasing number of autonomous vehicles and for any activity requiring Internet connectivity in a smart city. Also, VANETs allow on-board computers of mostly stationary vehicles, e.g., vehicles at an airport parking, to serve as resources of a mobile computer cloud with minimum help from the Internet infrastructure.

The contents produced and consumed by vehicles has *local relevance* in terms of time, space, and agents involved, the producer and the consumer. Vehicle-generated information has *local validity*, a limited spatial scope, an *explicit lifetime*, a limited temporal scope, and *local interest*, it is relevant to agents in a limited area around the vehicle. For example, the information that a car is approaching a congested area of a highway is relevant only for that particular segment of the road, at a particular time, and for vehicles nearby. The attributes of vehicular contents are summarized in Table 6.3 from [300].

One of the distinguishing characteristics of VANETs is the *content-centric distribution*, meaning that the content is important but the source is not. This is in marked contrast to the Internet, where an agent demands information from a specific source. For example, traffic information floods a specific area, and a vehicle retrieves it without concern for its source, while an Internet request for highway traffic information is directed to a specific site. Vehicle applications collect sensor data, and vehicles collaborate sharing sensory data. Sensory data is collected by vehicle-installed cameras, i.e., by on-board instruments. For example, *CarSpeak* allows a vehicle to access sensors on neighboring vehicles in

the same manner in which it can access its own [286]. *Waze* is a community-based traffic and navigation application allowing drivers real-time traffic and road information.

VANET communication protocols are similar to the ones used by wired networks; each host has an IP address. Assigning IP addresses to moving vehicles is far from trivial and often requires a Dynamic Host Configuration Protocol (DHCP) server, a heresy for ad hoc networks that operate without any infrastructure, using self-organization protocols. Vehicles frequently join and leave the network, and content of interest cannot be consistently bound to a unique IP address. A router typically relays and then deletes content.

6.14 Further readings

“Brief History of the Internet” [289] was written by the Internet pioneers Barry Leiner, Vinton Cerf, David Clark, Robert Kahn, Leonard Kleinrock, Daniel Lynch, Jon Postel, Larry Roberts, and Stephen Wolff. A widely used text of Kurose and Ross [285] is an excellent introduction to basic networking concepts. The book by Bertsekas and Gallager [61] gives insights into performance evaluation of computer networks. The classic texts on queuing theory of Kleinrock [276] are required reading for those interested in network analysis.

A survey of information-centric networking research is given in [522], while [539] is a succinct presentation of NDNs. Several publications related to software-defined networks are available on the site of the Open Network Foundation, <https://www.opennetworking.org/sdn-resources>. Moore’s law for traffic is discussed in [359]. Class-based queuing algorithm was introduced by Floyd and Van Jacobson in [181]. The *Black Widow* topology for system interconnects is analyzed in [443]. Storage Area Networks are analyzed in [472].

Scale-free networks and their applications are described by Barabási and Albert in [17–19,46]. The small-worlds networks were introduced by Watts and Strogatz in [507]. Epidemic algorithms for the dissemination of topological information are presented in [215,258,259]. Erdős-Rényi random graphs are analyzed in [67,165]. Energy-efficient protocols for cooperative networks are discussed in [162].

The future of fiber networks is covered in [290]. Vehicle ad hoc networks and their applications to vehicular cloud computing are discussed in [194,300,511]. An analysis of peer-to-peer networks is reported in [195]. Network management for private clouds, P2P networks, and virtual networks are presented in [356], [357], and [364], respectively.

6.15 Exercises and problems

- Problem 1.** Four ground rules for an open-architecture principles are cited in the “Brief History of the Internet” [289]. Analyze the implication of each one of these rules.
- Problem 2.** Key issues for network design in [289] are: (1) algorithms to prevent lost packets from permanently disabling communications; (2) host-to-host “pipelining” so that multiple packets could be en-route from source to destination, if intermediate networks allow it; (3) end-end checksums, reassembly of packets from fragments, and detection of duplicates, if any; (4) the need for global addressing; and (5) techniques for host-to-host flow control. Discuss how these issues were addressed by the TCP/IP network architecture.

- Problem 3.** Analyze the challenges of transition to IPv6. What will be in your view the effect of this transition on cloud computing?
- Problem 4.** Discuss the algorithms used to compute the TCP window size.
- Problem 5.** Creating a virtual machine (VM) reduces ultimately to copying a file, therefore the explosion of the number of VMs cannot be prevented; see Section 8.10. Virtualization could drastically lead to an exhaustion of the IPv4 address space. Analyze the solution to this potential problem adopted by the IaaS cloud service delivery model.
- Problem 6.** Read the paper describing the stochastic fair queuing algorithm [181]. Analyze the similarities and dissimilarities of this algorithm with the start-time fair queuing algorithm discussed in Section 9.6.
- Problem 7.** The small-worlds networks were introduced by D. Watts and S. H. Strogatz. They have two desirable features, high clustering and small path length. Read [507] and design an algorithm to construct a small-worlds network.
- Problem 8.** Scale-free networks are analyzed in [17–19,46]. Discuss the important features of systems interconnected by scale-free networks discussed in these papers.
- Problem* 9.** Consider the two 192 node fat-trees interconnect with two 96-way and twelve 24-way switches in Fig. 6.8. Compute the bisection bandwidth of the two interconnects.