
INTRODUCTION AND OVERVIEW

The past few decades have seen the merging of computer and communication technologies. Wide-area and local-area computer networks have been deployed to interconnect computers distributed throughout the world. This has led to a proliferation of many useful data communication services, such as electronic mail, remote file transfer, remote login, and web pages. Most of these services do not have very stringent “real-time” requirements in the sense that there is no urgency for the data to reach the receiver within a very short time, say, below 1s. At the other spectrum, the telephone network has been with us for a long time, and the information carried by the network has been primarily real-time telephone conversations. It is important for voice to reach the listener almost immediately for an intelligible and coherent conversation to take place.

With the emergence of multimedia services, real-time traffic will include not just voice, but also video, image, and computer data files. This has given rise to the vision of an integrated broadband network that is capable of carrying all kinds of information, real-time or non-real-time.

Many wide-area computer networks are implemented on top of telephone networks: transmission lines are leased from the telephone companies, and each of these lines interconnects two routers that perform data switching. Home computers are also linked to a gateway via telephone lines using modems. The gateway is in turn connected via telephone lines to other gateways or routers over the wide-area network. Thus, present-day computer networks are mostly networks overlaid on telephone networks. Strictly speaking, the telephone networks that are being used to carry computer data cannot be said to be integrated. The networks are designed with the intention that voice traffic will be carried, and their designs are optimized according to this assumption. A transmission line optimized for voice traffic is not necessarily optimal for other traffic types. The computer data are just “guests” to the telephone networks, and many components of the telephone network may not be optimized for the transport of non-voice services.

The focus of this book is on future broadband integrated networks. Loosely, the terms “broadband” and “integration” imply that services with rates from below one kbps to hundreds of Mbps can be supported. Some of these services, such as video conferencing, are widely known and anticipated, whereas others may be unforeseen and created only when the broadband network becomes available. The broadband network must be flexible enough to accommodate these unanticipated services as well.

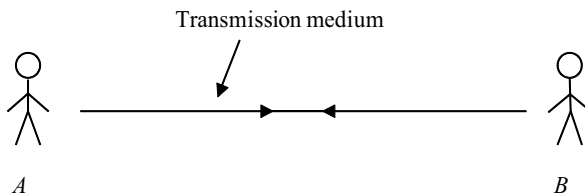
1.1 SWITCHING AND TRANSMISSION

At the fundamental level, a communication network is composed of switching and transmission resources that make it possible to transport information from one user to another. On top of the switching and transmission resources, we have the control functions, which could be implemented by either software or hardware, or both. Among other things, the control functions make it possible to automate the setting up of a connection between two users. At another level, they also ensure efficient usage of the switching and transmission resources. In a real network, the switching, transmission, and control facilities are typically distributed across many locations.

1.1.1 Roles of Switching and Transmission

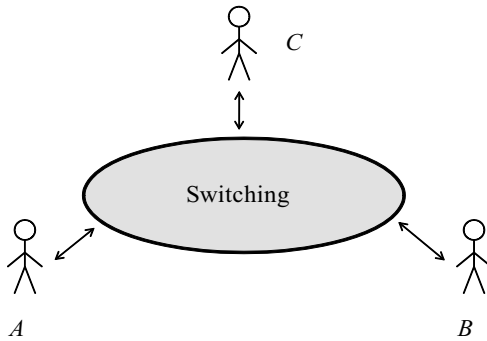
When there are only two users, as shown in Fig. 1.1, information created by one user is always targeted to the other user: switching is not needed and only transmission is required. In essence, the transmission facilities serve to carry information directly from one end of the transmission medium, which could be a coaxial cable, an optical fiber, or the air space, to the other end.

As soon as we have a third user in our network, the question of who wants to communicate with whom, and at what time, arises. With reference to Fig. 1.2, user *A* may want to talk to user *B* at one time but to user *C* later. The switching function makes it possible to change the connectivity among users in a dynamic way. In this way, a user can communicate with different users at different times.



When there are only two users, information from *A* is by default destined for *B*, and vice versa

FIGURE 1.1 A two-user network; switching is not required.



Information from *A* may be destined for *B* or *C*

FIGURE 1.2 A three-user network; switching is required.

It turns out that the locations of the switching facilities in a network have a significant impact on the amount of transmission facilities required in a network. Figure 1.3(a) depicts a telephone network in which the switching facilities are distributed and positioned at the *N* users' locations, and a user is connected to each of the other users via a unique line. Switching is performed when the user decides which of the *N* lines to use. When *N* is large, there will be many transmission lines and the transmission cost will be rather prohibitive.

In contrast, Fig. 1.3(b) shows a network in which each user has only one access line through which it can be connected to the other users. Switching is performed at

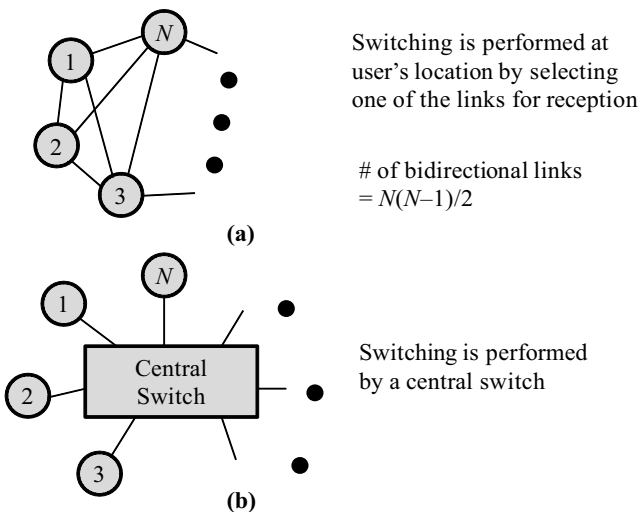


FIGURE 1.3 *N*-user networks with switching performed (a) at user's locations (b) by a central switch.

a central location. To the extent that a user does not need to speak to all the other users at the same time, this is a better solution because of the reduced number of transmission lines. Of course, if a user wants to be able to connect to more than one user simultaneously, multiple lines must still be installed between the user and the central switch.

In practice, a network typically consists of multiple switching centers at various locations that are interconnected via transmission lines. By locating the switching centers judiciously, transmission cost can be reduced.

1.1.2 Telephone Network Switching and Transmission Hierarchy

The switching and transmission facilities in a telephone network are organized in a hierarchical fashion. A simplified picture of a telephone network is given in Fig. 1.4. At the lower end of the hierarchy, we have subscribers' telephones located at business offices and households. Each telephone is connected to a switching facility, called a central office, via a subscriber loop. The switching center at this level is called the local office or the end office. If a subscriber wishes to speak to another subscriber linked to the same local office, the connection is set up by a switch at the local office.

Two local offices may be connected via either direct links or a higher level switching center, called a toll office. In the first case, there must be sufficient voice traffic between the two local central offices to justify the direct connection; otherwise, since the transmission capacities cannot be used by other local offices, they will be wasted. The second solution permits a higher degree of sharing of transmission resources. As illustrated in the figure, local offices *A*, *B*, and *C* are linked together via a toll office *D*. The transmission facilities between *D* and *C* are shared between *A* and *B*.

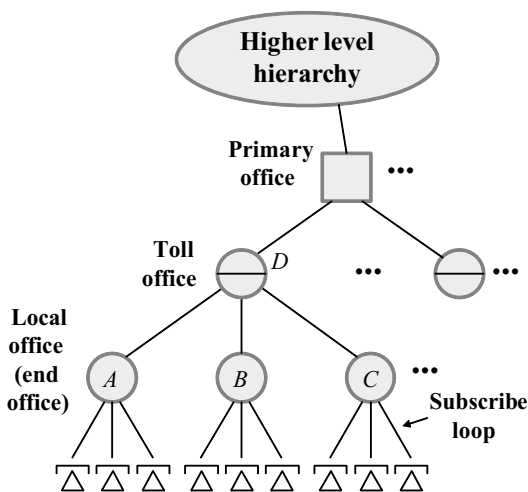


FIGURE 1.4 Telephone network hierarchy.

in the sense that both traffic between A and C and between B and C travel over them.

The toll offices are interconnected via an even higher level office, called the primary office. The primary offices are in turn connected by a yet even higher level office. Each level may move up to a higher level. In this way, a network hierarchy is formed.

The total amount of traffic reduces as we move up the hierarchy because of the so-called community-of-interest phenomenon. For instance, it is generally more likely for a user to make local phone calls than long-distance phone calls. The former may involve only a local switching office while the latter involves a series of switching offices at successive levels.

In short, an important objective achieved with the hierarchical network is the sharing of resources. The resources at the higher level are shared by a larger population of subscribers. The amount of resources can be reduced because it is statistically unlikely that all the subscribers will want to use the higher level resources at the same time.

Another advantage that comes with the hierarchical structure is the simplicity in finding a “route” for a connection between two subscribers. When subscriber i wants to connect to subscriber j , the local office of i first checks if j also belongs to the same office. If yes, switching is completed at the office. Otherwise, a connection is made between the local office and the next level toll office (assuming there are no direct links between the central offices of i and j). This procedure is repeated until an office with branches leading to both i and j is found.

1.2 MULTIPLEXING AND CONCENTRATION

Multiplexing and concentration are important concepts in reducing transmission cost. In both, a number of users share an underlying transmission medium (e.g., an optical fiber, a coaxial cable, or the air space).

As a multiplexing example, frequency-division multiplexing (FDM) is used to broadcast radio and TV programs on the air medium. In FDM, the capacity, or bandwidth, of the transmission medium is divided into different frequency bands, and each band is a logical channel for carrying information from a unique source. FDM can be used to subdivide the capacity of air medium, a coaxial cable, or any other transmission medium. Figure 1.5 depicts the transmission of digital information from a number of sources using FDM. Different carrier frequencies are used to transport different information streams. Receivers at the other end use bandpass filters to select the desired information stream.

Multiplexing can also be performed in the time domain. This is a more widely used multiplexing technique than FDM in telephone networks. Figure 1.6 illustrates a simple time-division multiplexing (TDM) scheme. The N sources take turns in a round-robin fashion to transmit on the transmission medium. Time is divided into frames, each having N time slots. Each source has a time slot dedicated to it in each frame. Thus, time slot 1 is assigned to source 1, time slot 2 to source 2, and so on. The slot positions i in successive frames all belong to the source i .

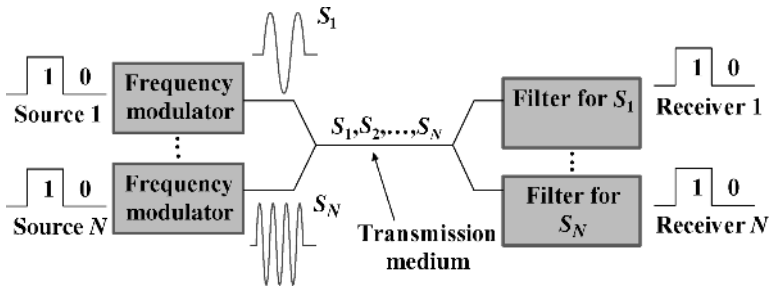


FIGURE 1.5 Frequency-division multiplexing.

In this book, we define switching as changing the connectivity between end users or equipments. The goal of a multiplexing system (consisting of the multiplexer, the transmission medium, and the demultiplexer) is not to perform switching; the goal is to partition a transmission medium into a number of logical channels, each of which can be used to interconnect a transmitter and a receiver. In the two scenarios above, each of the N multiplexed channels is dedicated exclusively to a transmitter–receiver pair, and which transmitter is connected to which receiver does not change over time. As an overall system, an input of the multiplexer is always connected to the same output of the demultiplexer. Thus, functionally, no switching occurs. Such is not the case with a concentrator.

Concentration achieves cost saving by making use of the fact that it is unlikely for all users to be active simultaneously at any given time. Therefore, transmission facilities need only be allocated to the active users. In the telephone network, for instance, it is unlikely that all the subscribers of the same local office want to use their phones at the same time. An $N \times M$ concentrator, as shown in Fig. 1.7, concentrates traffic from N sources onto M ($M < N$) outputs. A number of concentrators are usually placed at the “front end” of the local switching center to reduce the number

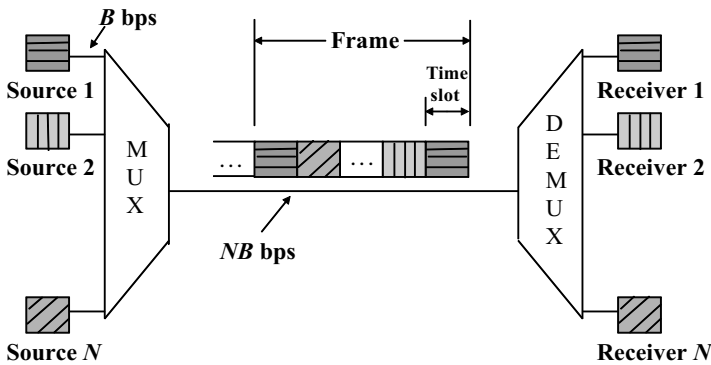
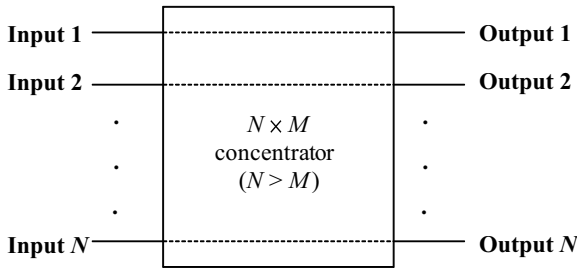


FIGURE 1.6 Time-division multiplexing.



An active input is assigned to one of the outputs. It does not matter which output is assigned.

FIGURE 1.7 An $N \times M$ concentrator.

of ports of the switch. An output of a concentrator, hence an input to the switch, is allocated to the subscriber only when he picks up the phone. The connectivity (i.e., which input is connected to which output) of the concentrator changes in a dynamic manner over time. If more than M sources are active at the same time, then some of the sources may be “blocked.” For telephone networks, M can usually be made considerably smaller than N to save cost without incurring a high likelihood for blocking.

Both multiplexers and concentrators achieve resource sharing, but in different ways. Let us refer to the sums of the capacities (bit rates) of the transmission lines connected to the inputs and outputs as the *total input and output capacities*, respectively. For a multiplexer, the total output capacity is equal to the total input capacity, whereas for a concentrator, the total output capacity is smaller than the total input capacity. The output capacity or bandwidth of the concentrator is said to be shared among the inputs, and that of the multiplexer is not. The concentrator outputs are allocated dynamically to the inputs based on need, and the allocation cannot be foretold in advance. In contrast, although a multiplexer allows the same transmission medium to be shared among several transmitter–receiver pairs, this is achieved by subdividing the capacity of the transmission medium and dedicating the resulting subchannels in an exclusive manner to individual pairs.

Statistical multiplexing is a packet-switching technique that can be considered as combining TDM with concentration. Consider a TDM scheme in which there are M ($M < N$) time slots in a frame. The output capacity is then smaller than the maximum possible total capacities of the inputs. A time slot is not always dedicated to a particular source. Slot 1 of frame 1 may be used by user 1, but user 1 may be idle in the next frame and slot 1 of frame 2 may be used by another user. The same slot positions of different frames can be used by different users, and therefore they can be targeted for different destinations in the overall communication network. To route the information in a slot to its final destination, we therefore need to encode the routing information in a “header” and attach the header to the information bits before putting both into

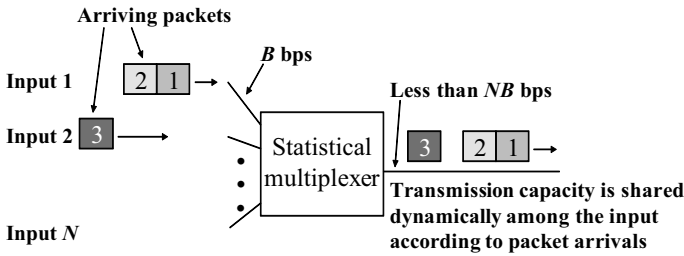


FIGURE 1.8 A statistical multiplexer.

the slot. Routing using header information is a basic operating principle in a packet-switched network. Each time slot is then said to contain a packet. The header in a packet generally contains either the explicit destination address or information from which the route to the final destination can be derived.

In circuit switching, the TDM scheme makes use of the fact that slots are assigned in a fixed way to derive the source and destination information implicitly. Since the same slot positions of successive frames are used by the same source and targeted for the same destination, they should be routed in a similar way. The route of these slots is determined during the call setup process (e.g., when setting up a voice connection); it will be used repeatedly in successive frames. In particular, no header is needed for routing purposes. In a packet network, the routing information must be incorporated into the header since the positions of a slot cannot be used to derive its route. To extend things even further, since the routing information is now contained in the header, we can even allow the time slot, or packet, to vary in length according to the amount of information transmitted by the source at each shot. The frame structure in TDM can be abandoned altogether in a packet network.

Note that unlike a time-division multiplexer, a statistical multiplexer performs switching by dynamically changing the user of its output (Fig. 1.8). We cannot foretell beforehand which source will use the output in advance. The users may also send data to the multiplexers in a random manner. It is therefore possible that there are more packets arriving at the statistical multiplexer than can be sent out on the output immediately. Thus, buffers or memories are required to store these outstanding packets until they can be cleared.

1.3 TIMESCALES OF INFORMATION TRANSFER

The above discussion of the concentrator and statistical multiplexer alluded to resource sharing at different timescales. In a telephone network, an output of a concentrator is assigned to a user only when the user picks up the phone and wants to make a call. The output is assigned for the duration of the call that typically lasts several minutes. In a packet-switched network, the output of a statistical multiplexer is assigned only for the duration of a packet, which typically is much less than 1s. A user may send out

packets sporadically during a communication session. It is important to have a clear concept of the timescales of information transfer to appreciate the fact that resource sharing and network control can be achieved at different timescales.

1.3.1 Sessions and Circuits

Before two end users can send information to each other, a communication session needs to be established. A telephone call is a communication session. As another example, when we log onto a computer remotely, we establish a session between the local terminal and the remote computer.

Network resources are assigned to set up a circuit or connection for this session. Some of these resources, such as an output of a concentrator in a circuit-switched network, may be dedicated exclusively to this connection while it remains active. Some of these resources, such as the output of a statistical multiplexer in a packet network, may be used by other sessions concurrently. In the latter, the transmission bandwidth is not dedicated exclusively to the session and is shared among active sessions, and the associated circuit is sometimes called a virtual circuit.

Some packet networks are not connection-oriented. It is not necessary to preestablish a connection (hence a route from the source to the destination) before data are sent by a session. In fact, successive packets of the session may traverse different routes to reach the destination. Although the concept of a connection is absent within the network, the end users still need to set up a session before they start to communicate. The setup time, however, can be much shorter than in a connection-oriented network because the control functions inside the network need not be involved for connection setup.

1.3.2 Messages

Once a session is set up, the users can then send information in the form of messages in an on-off manner. For a remote login session, the typing of the carriage-return key by the end user may result in the sending of a line of text as a message. Files may also be sent as a message. So, messages tend to vary in length.

For a two-party telephone session, for example, it is known that a user speaks only 40% of the time. The activity of the user is said to alternate between idle period and busy period. The busy period is called a talkspurt, which can be viewed as a message in a voice session. A scheme called time assigned speech interpolation (TASI) is often used to statistically multiplex several voice sources onto the same satellite link on a talkspurt basis.

1.3.3 Packets and Cells

Messages are data units that have meaning to the end users and have a logical relationship with the associated service. It could be a talkspurt, a line of text, a file, and so on. Packets, on the other hand, are transport data units within the network.

In a packet network, a long message is often fragmented into smaller packets before it is transported across the network. One possible reason for fragmentation could be that the communication network does not support packets beyond a certain length. For instance, the Internet has a maximum packet size of 64 Kbytes.

Another reason for message fragmentation is that most computer networks are store-and-forward networks in which a switching node must receive the entire packet before it is forwarded to the next node. By fragmenting a long message into many smaller packets, the end-to-end delay of the message can be made smaller, especially when the network is lightly loaded (see Problem 1.4).

Yet another motivation for fragmentation is to prevent a long packet from hogging a communication channel. Consider the output of a statistical multiplexer. While a long packet is being transmitted, newly arriving packets from other sources must wait until the completion of its transmission, which may take an excessively long time if there were no limit on its length. On the other hand, if the long packet has been cut into many smaller packets, the newly arriving packets from other sources have a chance to jump ahead and access the output channel after a short wait for the transmission of the current packet to complete.

Packet length can be variable or fixed. One advantage of the fixed packet-length scheme is that more efficient packet switches can be implemented. For instance, by time aligning the boundaries of the packets across the inputs of a packet switch, higher throughput can be achieved with the fixed packet-length scheme than with the variable packet-length scheme.

The fixed packet-length scheme has a disadvantage when the messages to be sent are much shorter than the packet length. In this case, only a small part of each packet contains the useful message, and the rest is stuffed with dummy bits to make up the whole packet. The observation suggests that small packets are preferable. However, the length of the packet header is largely independent of the overall packet length (e.g., the destination address length in the header is independent of the packet length). Hence, the header overhead (ratio of header length to packet length) tends to be larger for smaller packets. Thus, too small a packet can lead to high inefficiency as well.

In general, the determination of packet size is a complicated issue involving considerations from many different angles, not the least the characteristics of the underlying network traffic. The ITU (International Telecommunication Union), an international standard body, has chosen the asynchronous transfer mode (ATM) to be the information transport mechanism for the future broadband integrated network. An essence of the ATM scheme is that the basic information data unit is a 53-byte fixed-length packet called *cells*. The details of ATM and the motivations for the small-size cell will be covered in the later chapters.

1.4 BROADBAND INTEGRATED SERVICES NETWORK

The discussion up to this point forms the backdrop of the focus of this book—broadband integrated services networks. As the name suggests, an integrated network must be capable of supporting many different kinds of services.

Traditionally, services are segregated and each communication network supports only one type of service. Some examples of these networks are the telephone, television, airline reservation, and computer data networks. There are many advantages to having a single integrated network for the support of all these services. For instance, efficient resource sharing may be achievable. Take the telephone service. More phone calls are made during business hours than during the evening. The television service, on the other hand, is in high demand during the evening. By integrating these services on the same network, the same resources can be assigned to different services at different times.

Traditionally, whenever a new service is introduced, a new communication network may need to be designed and set up. Carrying the information traffic of this service on a network designed for another service may not be very efficient because of the dissimilar characteristics of the traffic. If an integrated network is designed with the forethought that some unknown services may need to be introduced in the future, then these services can be accommodated more easily.

The design of an integrated network taking into account the above concern is by no means easy. Figure 1.9 shows some services with widely varying traffic characteristics; the holding times (durations of sessions) and bit rate requirements of different services may differ by several orders of magnitude. Furthermore, some services, such as computer data transfer, tend to generate traffic in a bursty manner during a session (Fig. 1.10). Other services such as telephony and video conferencing generate traffic in a more continuous fashion.

The delay requirements may also be different. Real-time services are highly sensitive to network delay. For example, if real-time video data do not arrive at the display monitor at the receiver within certain time, they might as well be considered as lost.

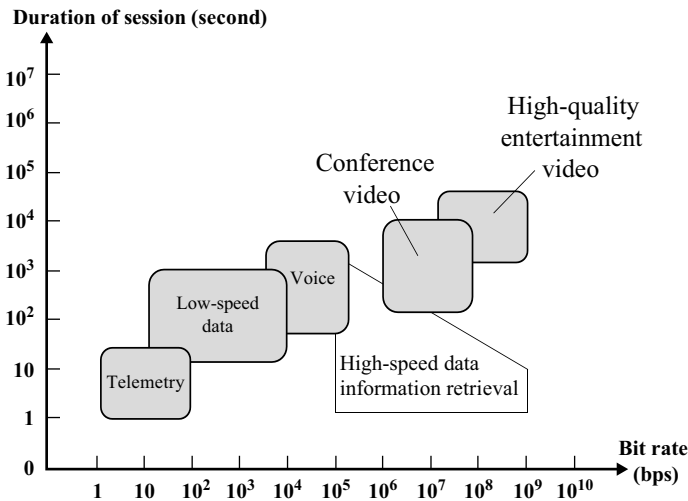


FIGURE 1.9 Holding times and bit rates of various services.

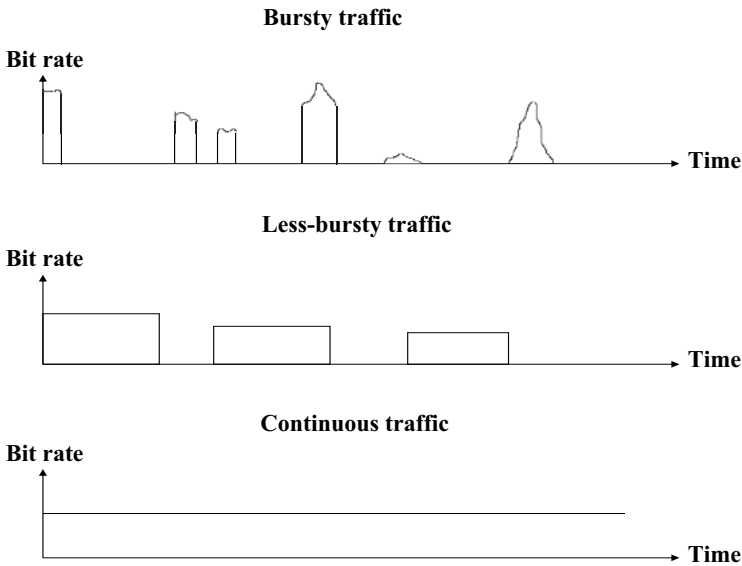


FIGURE 1.10 Traffic characteristics of sources of different burstiness.

How to control the traffic of these services to satisfy their diverse requirements is a nontrivial challenge that is still being worked on actively by researchers in this field.

As mentioned above, part of the challenge is to design the network to be flexible enough for the support of new services and services whose requirements have changed over time. As an example of the latter, advances in video and speech coding algorithms may well influence the characteristics of traffic generated by some services and thus change the service requirements. Finally, the integrated network must also be cost-effective and efficient for it to be successful.

PROBLEMS

- 1.1 This is a simplified problem related to resource sharing. You are among the 10 people sharing four telephones. At any time, a person is using or attempting to use a telephone with probability 0.2. What is the probability that all the telephones are being used when you want to make a call to your girl/boy friend? What if there are 100 people sharing 40 phones? Is it better to have a larger number of people sharing a pool of resources?
- 1.2 Does TDM perform the switching or transmission function?
- 1.3 Each telephone conversation requires 64 kbps of transmission capacity. Consider a geographical region with 6 million people, each with a telephone. We want

to multiplex all the telephone traffic onto a number of optical fibers, each with capacity of 2.4 Gbps. How many fibers are needed?

- 1.4 Consider the problem of fragmenting a message of 1000 bytes into packets of x bytes each. Each packet has a header (overhead) of 10 bytes. The message passes through five links of the same transmission rate in a store-and-forward network on its way to destination. There is no other traffic in the network. What is the optimal packet size x in order to minimize end-to-end delay?
- 1.5 Consider 100 active phone conversations being multiplexed onto 60 lines by TASI. Assume a person speaks only 40% of the time when using a telephone. Talkspurts are clipped when all 60 lines have been used. What is the probability of a talkspurt being clipped?

