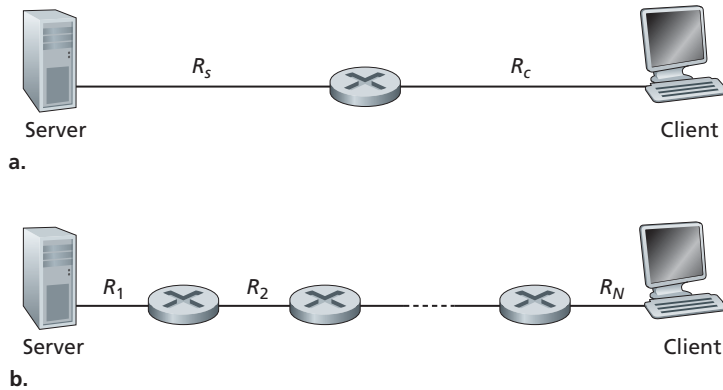


transmit a packet into a shared medium (e.g., as in a WiFi or cable modem scenario) may *purposefully* delay its transmission as part of its protocol for sharing the medium with other end systems; we'll consider such protocols in detail in Chapter 5. Another important delay is media packetization delay, which is present in Voice-over-IP (VoIP) applications. In VoIP, the sending side must first fill a packet with encoded digitized speech before passing the packet to the Internet. This time to fill a packet—called the packetization delay—can be significant and can impact the user-perceived quality of a VoIP call. This issue will be further explored in a homework problem at the end of this chapter.

### 1.4.4 Throughput in Computer Networks

In addition to delay and packet loss, another critical performance measure in computer networks is end-to-end throughput. To define throughput, consider transferring a large file from Host A to Host B across a computer network. This transfer might be, for example, a large video clip from one peer to another in a P2P file sharing system. The **instantaneous throughput** at any instant of time is the rate (in bits/sec) at which Host B is receiving the file. (Many applications, including many P2P file sharing systems, display the instantaneous throughput during downloads in the user interface—perhaps you have observed this before!) If the file consists of  $F$  bits and the transfer takes  $T$  seconds for Host B to receive all  $F$  bits, then the **average throughput** of the file transfer is  $F/T$  bits/sec. For some applications, such as Internet telephony, it is desirable to have a low delay and an instantaneous throughput consistently above some threshold (for example, over 24 kbps for some Internet telephony applications and over 256 kbps for some real-time video applications). For other applications, including those involving file transfers, delay is not critical, but it is desirable to have the highest possible throughput.

To gain further insight into the important concept of throughput, let's consider a few examples. Figure 1.19(a) shows two end systems, a server and a client, connected by two communication links and a router. Consider the throughput for a file transfer from the server to the client. Let  $R_s$  denote the rate of the link between the server and the router; and  $R_c$  denote the rate of the link between the router and the client. Suppose that the only bits being sent in the entire network are those from the server to the client. We now ask, in this ideal scenario, what is the server-to-client throughput? To answer this question, we may think of bits as *fluid* and communication links as *pipes*. Clearly, the server cannot pump bits through its link at a rate faster than  $R_s$  bps; and the router cannot forward bits at a rate faster than  $R_c$  bps. If  $R_s < R_c$ , then the bits pumped by the server will “flow” right through the router and arrive at the client at a rate of  $R_s$  bps, giving a throughput of  $R_s$  bps. If, on the other hand,  $R_c < R_s$ , then the router will not be able to forward bits as quickly as it receives them. In this case, bits will only leave the router at rate  $R_c$ , giving an

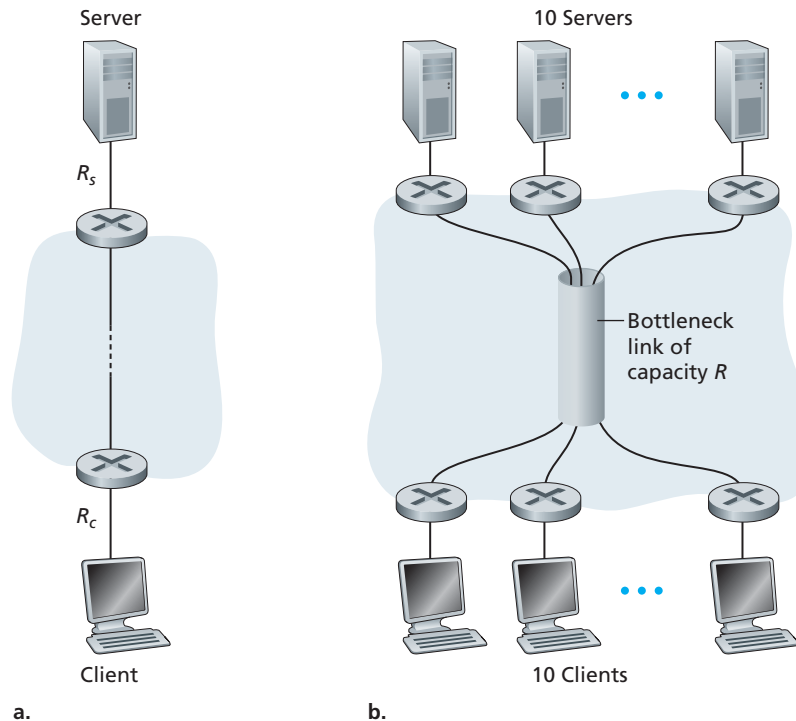


**Figure 1.19** ♦ Throughput for a file transfer from server to client

end-to-end throughput of  $R_c$ . (Note also that if bits continue to arrive at the router at rate  $R_s$ , and continue to leave the router at  $R_c$ , the backlog of bits at the router waiting for transmission to the client will grow and grow—a most undesirable situation!) Thus, for this simple two-link network, the throughput is  $\min\{R_c, R_s\}$ , that is, it is the transmission rate of the **bottleneck link**. Having determined the throughput, we can now approximate the time it takes to transfer a large file of  $F$  bits from server to client as  $F/\min\{R_s, R_c\}$ . For a specific example, suppose you are downloading an MP3 file of  $F = 32$  million bits, the server has a transmission rate of  $R_s = 2$  Mbps, and you have an access link of  $R_c = 1$  Mbps. The time needed to transfer the file is then 32 seconds. Of course, these expressions for throughput and transfer time are only approximations, as they do not account for store-and-forward and processing delays as well as protocol issues.

Figure 1.19(b) now shows a network with  $N$  links between the server and the client, with the transmission rates of the  $N$  links being  $R_1, R_2, \dots, R_N$ . Applying the same analysis as for the two-link network, we find that the throughput for a file transfer from server to client is  $\min\{R_1, R_2, \dots, R_N\}$ , which is once again the transmission rate of the bottleneck link along the path between server and client.

Now consider another example motivated by today's Internet. Figure 1.20(a) shows two end systems, a server and a client, connected to a computer network. Consider the throughput for a file transfer from the server to the client. The server is connected to the network with an access link of rate  $R_s$  and the client is connected to the network with an access link of rate  $R_c$ . Now suppose that all the links in the core of the communication network have very high transmission rates, much higher than  $R_s$  and  $R_c$ . Indeed, today, the core of the Internet is over-provisioned with high speed links that experience little congestion. Also suppose that the only bits being sent in the entire network are those from the server to the client. Because the core of the computer network is like a wide pipe in this example, the rate at which bits can flow



**Figure 1.20** ♦ End-to-end throughput: (a) Client downloads a file from server; (b) 10 clients downloading with 10 servers

from source to destination is again the minimum of  $R_s$  and  $R_c$ , that is, throughput =  $\min\{R_s, R_c\}$ . Therefore, the constraining factor for throughput in today's Internet is typically the access network.

For a final example, consider Figure 1.20(b) in which there are 10 servers and 10 clients connected to the core of the computer network. In this example, there are 10 simultaneous downloads taking place, involving 10 client-server pairs. Suppose that these 10 downloads are the only traffic in the network at the current time. As shown in the figure, there is a link in the core that is traversed by all 10 downloads. Denote  $R$  for the transmission rate of this link  $R$ . Let's suppose that all server access links have the same rate  $R_s$ , all client access links have the same rate  $R_c$ , and the transmission rates of all the links in the core—except the one common link of rate  $R$ —are much larger than  $R_s$ ,  $R_c$ , and  $R$ . Now we ask, what are the throughputs of the downloads? Clearly, if the rate of the common link,  $R$ , is large—say a hundred times larger than both  $R_s$  and  $R_c$ —then the throughput for each download will once again be  $\min\{R_s, R_c\}$ . But what if the rate of the common link is of the same order as  $R_s$  and  $R_c$ ? What will the throughput be in this case? Let's take a look at a specific