# 1

# PERFORMANCE EVALUATION IN TELECOMMUNICATIONS

## 1.1  INTRODUCTION: THE TELEPHONE NETWORK

The performance of a telecommunications system is a subject of considerable importance. Before going further, we must first describe the system under study.[1] The model of a generic telephone system as depicted in Figure 1.1 shows four basic components: *customer premises equipment*, the *local network*, the *switching plant*, and the *long-haul network*. As we look in detail at each of the categories, we see increasing complexity reflecting the explosive growth of the telecommunications industry.

### 1.1.1  Customer Premises Equipment

The most common piece of equipment on the customer's premises is, of course, the ordinary telephone, the handset. Today, a number of other items can also be found. There is scarcely a telephone without an answering service, which is provided by a box next to the phone. Up until the early 1990s the Internet was the domain of the technical types. Today the personal computer with the attendant *modem* is an appliance found in many homes. Higher-capacity modems convey data from more

---

[1]The Further Reading section at the end of the chapter includes several telecommunication sources that we have found to be useful. The history of telecommunications is itself an interesting study. We have found two Websites on the subject: *http://www-stall.rz.fht-esslingen.de/telehistory/* and *http://www.webbconsult.com/timeline.html*.
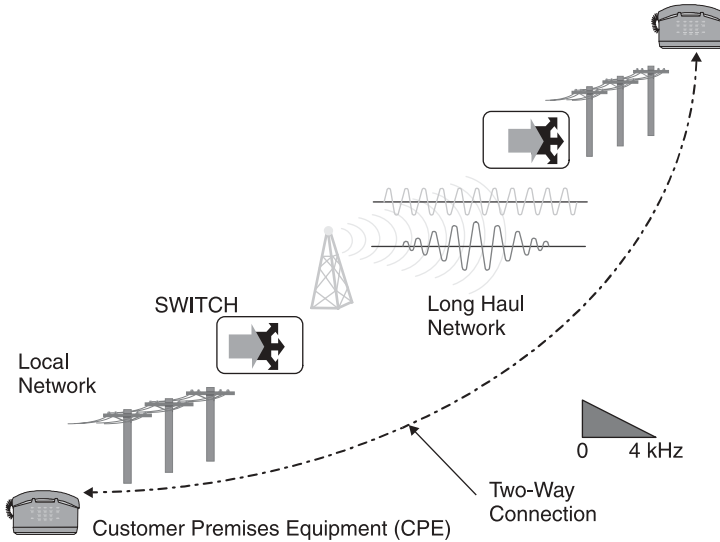
**Figure 1.1**  Telephone system.

complex processing center industrial sites (see discussion below). *Facsimile* (Fax) *machines* have also experienced explosive growth in deployment, and the private branch exchange (PBX), which is simply a local switchboard for connecting the phones in an office building to the outside world, has in many instances been replaced by central national switching hubs.

In industrial, commercial, and scientific sites an ubiquitous installation is the *local-area network* (LAN). A LAN may range from a simple bus with a few computers to a complex of switches and routers linking together hundreds of computers and specialized pieces of equipment. A LAN may or may not be connected to the telephone network, but it is a very important part of modern telecommunication. A great deal of effort has been devoted to a study of its performance.[2]

The complexity and growth of the network precludes easy categorization. Where do we place the ubiquitous cell phone? Considering its function, we can put it in the same category as the ordinary wired phone as a piece of customer service equipment.

### 1.1.2   The Local Network

The *local network* is the means by which the customer premises equipment is connected to the telephone network. All of the ordinary telephones are connected to the local telephone end office through pairs of 22- or 26-gauge wire twisted together

---

[2]Three papers from which modern LANs have evolved are listed in the Further Reading section. For example, the ubiquitous Ethernet may be seen as a lineal descendant of ALOHA.

in order to minimize crosstalk. Hundreds of these twisted pairs are bundled in a cable. In order to improve the quality of the analog voice signal, *loading coils* are attached to the twisted pairs. These same coils are removed for digital transmission of modem signals. Since they are by far the most common means of access to the premises, a great deal of work is going into increasing the information-carrying capability of twisted pairs. This would allow the classic telephone network to play a role in what might be called the Internet revolution with a relatively modest increase in investment. The generic term for these techniques is *digital subscriber links* (DSLs). Asynchronous *digital subscriber links* (ADSLs) recognize that most Internet users receive far more information than they transmit; accordingly, the rates provided are 1.536 megabits per second (Mbps) on the downlink and 400 kilobits per second (kbps) on the uplink. These datastreams coexist with the ordinary voice signal on the same line. The latest development in this area is the very high-data-rate digital subscriber line (VDSL), which provides uplink and downlink rate of 52 Mbps over optical links close to the subscriber premises.

***Blurring of Distinctions***   Until recently, the telephone network and the relatively primitive cable television network (CTV) have been physically separate and distinct in function; however, the dizzying pace of regulatory as well as technical development has led to the CTV and the telephone companies competing for the same business. The immediate point of contention is Internet access. The *cable modem* operation over the coaxial cable entering customer premises has the potential to allow rapid access to Internet material. The same technology would allow the same cable to carry a wide range of services. There is a potential topological problem. The CTV network is a *fanout*, where one point feeds many; accordingly, techniques for controlling uplink traffic, specifically, many to one, are required.

CTV companies are certainly not secure in their base market. Direct-access satellites and ground radio channels leapfrog the local cable network. Although these services now focus on distribution of video images, it is only a matter of time until there is an uplink channel and access to the Internet for a range of services.

***Wireless Transmission***   An area of explosive growth is wireless transmission. It seems that it is not possible to walk down the street without seeing someone, possibly on a bicycle, chatting away with someone on a cell phone. The handset receives and transmits from a base station that serves a limited geographic area, the cell. The base station serves the same function as the end office in the wire network, connecting to the system at large. The connection could be to a mobile user in the same cell or, through switching, to the long haul. In developing countries, wireless services are a great economic benefit since they avoid the installation of local (landline) wire distribution, which is the most expensive part of the network. Up to this point, wireless services have emphasized voice and low speed data of the order of 19.2 kbps; however, the future growth is in high speed multimedia data services of the order of Mbps. The current deployment of 3G wireless network strives to provide 2 Mbps per user and the leading research on 4G wireless networking technologies aims at providing QoS based services at a rate of the order of 20 Mbps.

### 1.1.3   Long-Haul Network

The long-haul network carries traffic from one telephone end office to another end office. In general, the long-haul network is a mesh of interconnected links (see Fig. 1.2). There may be a number of links in this path. The switch, which will be discussed in the next subsection, serves to route the flow of traffic between links.

A number of different kinds of transmission media may be used to implement the links in the long-haul network: twisted pairs with repeaters, coaxial cable, microwave radio, satellite, and optical fiber. Increasingly, optical fiber is replacing the metallic media and microwave radio as the transmission medium. This is the case for transoceanic cable as well. The dominance of optical fiber is not difficult to understand since it provides virtually error-free transmission at gigabit rates. The only challenge to the hegemony of fiber is the satellite system in its area of application. Satellites allow direct access to any point in their footprint; thus, new networks can be set up quickly. Further, earthly impediments and distances constitute no barrier. Satellites are unsurpassed in linking to remote areas, for example.

### 1.1.4   Switching

It is clear that a telephone system of any size consisting of hard connections between all pairs of subscribers would be impossibly large; accordingly, there is a need for trunks that are transmission lines that can be used by different pairs of users at different times (see Fig. 1.2). The switching onto these trunks is carried out at exchanges. The first manual exchange was developed in 1878. The first automatic exchange was put into operation in Illinois in 1892.

Until the late 1960s the traffic was routed using *circuit switching* in which an electrical path with a nominal 4 kHz bandwidth is established between transmitter and receiver (see Fig. 1.3). This technique works well for voice traffic since the time required to set up the path, which is approximately 0.5 s, is small compared to the duration of a typical call whose average is about 3 minutes.
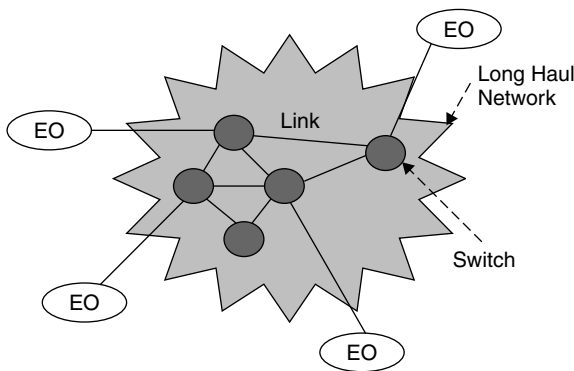
**Figure 1.2**    Long-haul network.

Circuit-Switching

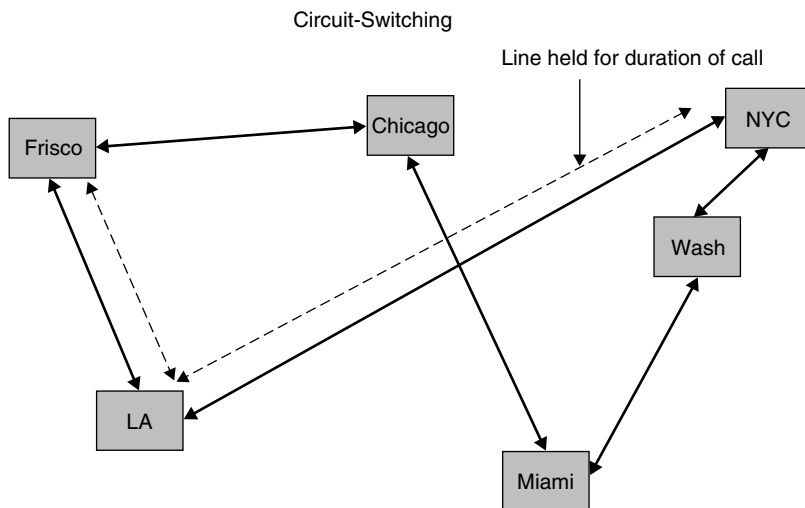Line held for duration of call



**Figure 1.3** Circuit switching.

The call setup time is one of the problems encountered by data traffic connections. A typical data message may only contain 1000 bits; accordingly, even at the low rate of 4.8 kbps, the transmission time is of the same order as the setup time. The remedy *is packet switching*, whereby information is segmented into fixed-size blocks. Source and destination addresses are appended to the data payload in addition to control fields for recovering from errors and other desired functions to form a packet. Transmission facilities are dedicated only for the time required to transmit the packet.

The addition of parity check bits combats the high bit error rate, which was a common impairment of the old voice network. The packets are routed through the network in *store-and-forward* fashion over links connecting packet switches (see Fig. 1.4). At each switch along the way the packet can be checked for errors before being routed to the destination. Implementation of early packet-switched networks, notably the ARPANET, have evolved to the present Internet. Seminal papers in the development of the Internet are listed at the end of the chapter.

There is a form of packet switching, which emulates circuit switching. In *asynchronous transfer mode* (ATM) the packets, called cells in this context contain 424 bits including 40 overhead bits. The small, constant cell length allows processing associated with routing and multiplexing to be done very quickly with current technology. The basic objective is to have a degree of flexibility in handling a wide range of traffic types. The cells making up a call follow the same path for its duration, which is called a *virtual circuit*.

In the early implementation of ARPANET-based IP networks, the focus was on handling data traffic in a best-effort manner. Since the IP network provided a connectionless network service and then the applications such as ftp (File Transfer
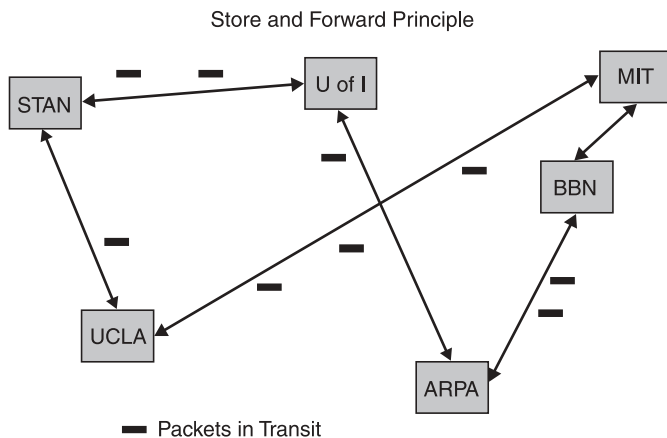
Store and Forward Principle



**Figure 1.4**   Packet-switched networks.

Protocol) and Telnet (a terminal application program for TCP/IP networks) were the only applications widely used, the requirements of service from the network for such applications were not so stringent. Although users would have preferred better response times, for their real-time sessions, such as Telnet, or chat (UNIX-based), the focus of the network implementers was on improving the ad hoc techniques developed from the point of implementation. With World Wide Web (WWW) traffic and the explosive growth of new (in particular e-commerce) applications, an IP network with just best-effort service was not sufficient. This lead to the definition of *differentiated* services and *integrated* services in the context of next-generation IP networks.

### 1.1.5   The Functional Organization of Network Protocols

A useful view of a telecommunications system is by the functions that it must perform in order to convey information from one point to another. The Open Systems Interconnection (OSI)[3] protocol structure provides a delineation of these functions (see Fig. 1.5). In OSI closely related functions are grouped into one of seven layers. This grouping allows simple interface between the layers. Although the seven-layer structure can be fitted to any telecommunications network, it is most suited to packet-switched networks. Fortunately, these are the networks that we are most interested in.

At the lowest level, the *physical layer*, the signal that is to be conveyed from one point to another is transformed into a signal suitable for transmission over the medium at hand. In digital communication systems, our main concern is that zeros

---

[3] The original purpose of OSI was to provide a common framework for all telecommunications. Because of its complexity and the rapid advance of the technology, it never functioned in this way.
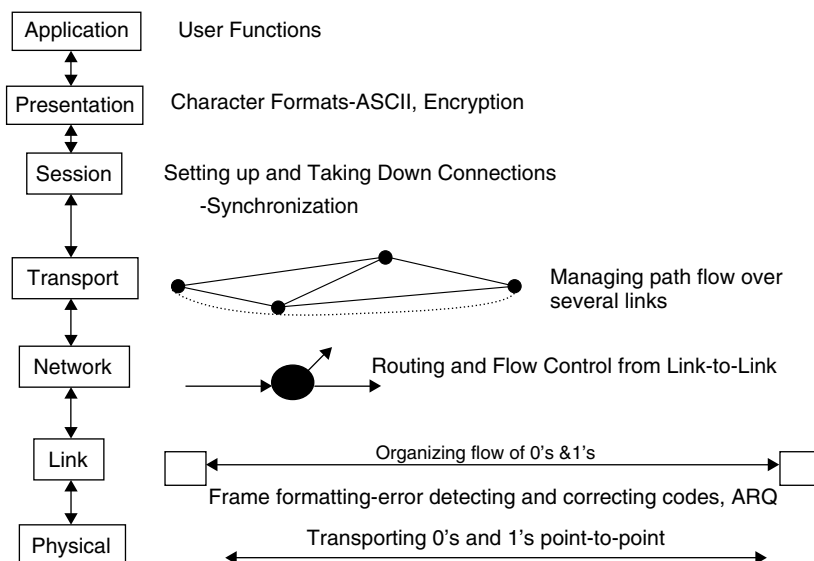
| | |
|---|---|
| Application | User Functions |
| Presentation | Character Formats-ASCII, Encryption |
| Session | Setting up and Taking Down Connections -Synchronization |
| Transport | Managing path flow over several links |
| Network | Routing and Flow Control from Link-to-Link |
| Link | Organizing flow of 0's &1's / Frame formatting-error detecting and correcting codes, ARQ |
| Physical | Transporting 0's and 1's point-to-point |

**Figure 1.5** OSI seven-layer reference model.

and ones of a digital stream modulate electrical or optical pulses that are transmitted over the respective media. The implementation of functions such as signal filtering and phase and timing recovery and tracking are handled at this level. The performance issue is to perform these functions in such a way that the probability of error is minimized within the implementation constraints.

The physical layer delivers raw bits to the *link level*, whose function is to organize the flow of bits over a segment of the path called a *link*. A link could be, for example, an optical fiber between buildings on a university campus or an intercontinental satellite channel. The basic format of bits at the link level is the *frame*, which consists of user information bits together with a number of overhead bits. The overhead bits perform several functions. In general, frames are not of fixed length; accordingly, framing bits indicating the beginning and end of a frame must be provided. Parity bits can be included in a frame in order to detect and/or correct errors that have occurred in transmission. In order to ensure that frames are delivered in order, sequencing bits must also be included in the frame. The primary performance issue at the link level is the efficacy of error detection and correction.

The frames, which are supposed to be error-free, are passed from the link layer to the *network layer*, which is responsible for routing and flow control of the links in the path between the source and the destination. At this level, a packet is formed by adding addresses as well as other overhead to the frame.

In packet-switched networks buffering is required to smooth flow between links. The *probability of buffer overflow* is a function of the rate of traffic flow and the information-carrying capacity of the link. A routing strategy could be based on minimizing this probability.

At the *transport layer* the end-to-end flow over all of the links in a path is managed. Flow control can be used for efficient operation between systems whose speeds are mismatched and as a type of congestion control. In some systems, errors are detected and corrected at this level. Further, there is a requirement to provide enough buffering to bring probability of buffer overflow to an acceptable level.

Calls are established at the *session layer.* In a circuit-switched network, the task is to find a suitable path through the network. The criterion that is relevant here is the *probability of call blocking*, specifically, the probability that a suitable path is not available to an arriving call. In packet-switched networks where resources are assigned on demand, the criterion is whether there will be enough resources available throughout the duration of a call. This is the *admission control* problem. The other functions of the session layer include data transfer control using tokens, in the case of request-response type dialogs (i.e. half duplex), synchronization in case of data transfer failures.

The performance issues that we will be considering do not arise in the two upper layers. At the *presentation layer* the information is formatted for purposes other than communication. Encryption is a form of formatting, for example. The *application layer* deals with the particular function that the user is exercising, such as mail or image transfer, for example. Since many functionalities of different applications have a common structure for the communication, the application layer supports these common functionalities in terms of protocols called application service elements (ASEs) called common ASEs (CASEs) and those functionalities which are specific to the application are modeled as specific ASEs (SASEs). The overall control of these ASEs to implement a particular application is done by using Control Functions (CFs). For modeling applications OSI recommends the concept called abstract service definition conventions (ASDC).

## 1.2   APPROACHES TO PERFORMANCE EVALUATION

There are three general approaches to evaluating the performance of a network. The one with which we are most familiar through course work is analysis. The calculation of the probability of error for a time-invariant channel disturbed by additive white Gaussian noise should be a familiar example. Calculation of the probability of overflow for Poisson arrival of packet to a buffer is another. At the switch level the probability of blocking, that is, the probability that an output line is not available, is of interest.

Analysis is certainly the best approach for simple models since it is fast and accurate. The problem is that mathematical approaches are tractable for only a limited number of models. There is a real art in finding an analyzable model, which approximates a given real system to some degree of accuracy. Analysis is most useful in the early stages of a project when it is necessary to do a rough assessment of available options.

After an analysis, the next step of refinement in the development of a system would be *computer simulation*, in which a more detailed model of the system under study is emulated in software. In most cases of interest, one wishes to assess the effect of random inputs in telecommunications systems. Random path delays in wireless channels and random arrival patterns to a switch are examples. The technique for

dealing with a system with random stimuli is called *Monte Carlo simulation*.[4] The essence of this approach is repeated trials to obtain a set of responses to random inputs. A statistical analysis is performed on the output set in order to estimate a performance measure. Two examples will serve as illustration:

1. The number of bit errors that occur when a channel is disturbed by non-Gaussian noise is counted. This number divided by the total number of trials serves as an estimate of the probability of bit error.
2. The time of arrival to and departure from a system of a sequence of messages is recorded. The average message delay in the system can be estimated from these data.

In no sense is simulation a substitution for analysis—both techniques have their own place. If the program is even moderately complicated, it is difficult to be certain that there are no errors despite thorough checking. An analytic model is a valuable authentication tool. It can ascertain whether the simulation outputs are in the ballpark. Also, it may be possible to run the simulation for cases that can be analyzed. Mathematics is necessary to analyze the results of simulation. For example, how many data points are required for one to judge the accuracy of an estimate? Finally, in addition to verification and analysis of results, mathematics can play a role in implementation. Analytically tractable models may be used for subsystems of a larger system, thereby simplifying the program. Furthermore, in many cases the event being measured is rare. For example, the probability of error on an optical link may be on the order of $10^{-9}$; to obtain valid estimates, $10^{11}$ samples may be required, implying a prohibitively long runtime for the program. Mathematical techniques can be used to obtain accurate results with reasonable run times. The variance reduction technique treated in Chapter 9 are examples.

With enough time and effort expended the most accurate approach to evaluating a system is building a *prototype*. It is the kind of exercise that is best suited to the final stages of a project; it is really too inflexible for early to middle stages of design. Certainly, one would not want to build a prototype to test the large number of alternatives that arise in the early stages of a project. Modern telecommunication systems have a large software component. In this respect the distinction between prototype and simulation could be blurred. A working simulation of certain system components could be directly converted into silicon, for example.

## 1.3 QUEUEING MODELS

### 1.3.1 Basic Form

The primary analytic tool that we will use to evaluate systems is *queueing theory*, which is the application of stochastic processes to the study of

---

[4]The fundamentals are presented in Chapter 9.

waiting lines.[5] The generic model for queueing systems is illustrated in Figure 1.6. It consists of three basic components: (1) an arrival process, (2) a storage facility, and (3) a server. A bakery is an everyday example of a queueing system. Customers arrive at rates that vary according to the season and time of day. If the sales personnel are occupied, the customer queues, that is, stands in line. The time required to fill an order varies according to the customer's demands. Of course, the number of servers is the number of sales clerks.

The *Poisson arrival process* is the most widely applied model of the arrival process and is the most tractable mathematically. In the models that we shall study, the service time is a random variable, which is independent of the arrival process in our study. We develop a general model for the service time. Storage facilities hold customers until servers are available. There are two important cases. The storage facilities can be so large that they may be considered infinite. On the other extreme, we have facilities that hold only customers who are in the process of being served.

Queueing models are widely applicable. A search of a university library for books applying queueing theory showed 96 books on the following topics: service and manufacturing, storage facilities with special emphasis on dams, inventory and maintenance, construction and mining, insurance risk, and social organization. This wide applicability notwithstanding, the most successful and the most important application of queueing models has been and continues to be telecommunications. This is the application that is the subject of this book. The generic model for telephony is shown in Figure 1.7. The arrival process consists of the generation of calls or messages. (The distinction will be made clear later.) Calls and messages are stored in a buffer prior to transmission over telephone lines. The call- and message-handling capacity of these lines is a key element in determination of the performance of the system.

In terms of performance requirements, telephone systems fall into two basic categories: *loss and delay systems.* In the former, calls leave the system if transmission facilities are not available. An everyday example is trying to place a telephone call at a later time when the line is busy. In telephone applications handling data, it is frequently the case that delay is not as critical as loss and messages can be stored until transmission facilities are available.

### 1.3.2  A Brief Historical Sketch

***The Classical Period—Erlang and Others***   Before getting on with the detailed mathematical models in the rest of the text, we begin with an historical sketch of the field of queueing models. Virtually all the results that we cite will be derived later in the text. We can regard the development queueing theory as falling into three

---

[5]At the end of the chapter several queueing theoretic texts are listed, each shedding light on a different aspect of the subject. While all cover the fundamentals, the first three of these give an interesting historical perspective. We have found that the next four have special tutorial merit. Finally, the last two delve more deeply into the subject.

**Figure 1.6**  Generic queueing model.

distinct phases, which, for purposes of explanation, we call the *classical*, the *romantic*, and the *modern eras*.

As we all know, the genesis of the telecommunications industry was the invention of the telephone by Alexander Graham Bell in 1876. Bell subsequently exploited his invention for practical commercial purposes. The fact that the name Bell and telephony are inextricably linked bears witness to his energy and foresight. It is an often repeated anecdote that the telegraph company declined the opportunity to deploy the telephone because it could see no practical use for speech separated from physical contact.

It is clear that a telephone system of any size consisting of hard connections between all pairs of subscribers would be impossibly large; accordingly, there is a need for trunks, which are transmission lines that can be used by different pairs of users at different times (see Fig. 1.2). The switching onto these trunks is carried out at exchanges. The first manual exchange was developed in 1878. The first automatic exchange was put into operation in Illinois in 1892. These advances gave rise to questions of performance. For example, how many trunks would be required between a pair of exchanges so that subscribers could be connected to one another a high percentage of the time?

Questions of performance were compounded as telephone networks grew in size and complexity. The basic form of the telephone system was the *circuit-switched network* shown in Figure 1.3. For the duration of a call, transmission facilities are dedicated to a call over a fixed path or circuit. The time required to set up this path is an important parameter of service.

The first analysis of telephone traffic was reported by G. T. Blood in an unpublished memorandum in 1898. Although there are other early contributors—Rory, Joannen, and Grinsted—the father of queueing theory is undoubtedly the Danish mathematician, Agner Krarup Erlang (1878–1929). Erlang's models were based on the Poisson arrival process and an exponentially distributed duration of
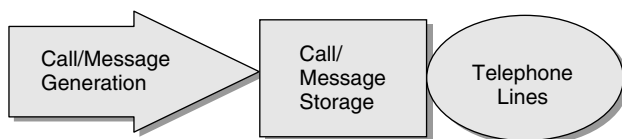
**Figure 1.7**  Generic telephony model.

calls.[6] During the period 1909–1917 he obtained steady-state solutions leading to formulas still used by telephone engineers. The *Erlang B* formula gives the probability that a trunk is not available as a function of demand or load and the number of trunks in a *loss system* in which calls cannot be stored. The same result for a *delay system* in which calls can be held until a line becomes available is given by the *Erlang C* formula.

In 1939, Erlang's analysis and those of subsequent contributors were unified by Feller by the application of the theory of birth and death processes. Among his many insights into the analytical model, Erlang had conjectured that the probability of a lost call (Erlang B) is insensitive to the probability distribution of the call duration. This result was subsequently proved by Kosten in 1948. Extensions and refinements of Erlang's work were carried out by Molina, Engset, and O'Dell, all of whom were on the staff of telephone administrations. This later work was concerned with such questions as retries when a call is initially blocked. The last result of the classical era that we'll cite is the *Pollaczek–Khinchin* formula, which gives the average delay in completing service in a system with Poisson arrivals, a general distribution of service time, and infinite storage. This formula was subsequently derived by Kendall in 1951 using the theory of *imbedded Markov chains.* We cover this material in Chapters 3, 5, and 6 of the text.

***The Romantic Era—Packet Switching***     As discussed above, deficiencies in the operation of the voice-oriented telephone network led to the development of packet switching. The pioneering work in the analysis of this new kind of network was Kleinrock's doctoral thesis. This work has as its basis *Jackson networks*. This concept and its extensions are covered in Chapter 4.

A significant category of packet network developed later has been the local-area network. Much of the queuing theoretic work in this area focused on techniques for sharing a common transmission medium among a number of users. The simplest approach is dedicated capacity, as in *frequency-division multiple access* (FDMA) and *time-division multiple access* (TDMA). An alternative is the *random access technique* introduced in the *ALOHA* network. Other techniques, such as token passing, can be analyzed as variations on *polling*. These media access techniques are analyzed in Chapters 5 and 6.

***The Modern Era***     Two developments in the technology, optical fiber and very large-scale integration (VLSI), have given rise to the modern era in tele-communications. The first of these provided orders of magnitude increases in the capacity available on transmission facilities, with rates of the order of Gbps and very low error rates. The second allows the implementation of modern digital processing techniques for switching and multiplexing. Further, VLSI is at the heart of modern computer technology, which in itself stimulates applications. The result of these developments is the development of many new services, each with its own traffic

---

[6]In the subsequent chapters of the text, we shall derive the results mentioned in the remainder of this paragraph and in the next subsection, on the romantic era.

characteristic and performance requirements. From the point of view of queueing models the traffic may be divided into two categories: real-time and delay-insensitive, corresponding roughly to the loss and delay systems, respectively, discussed above. In the real-time class we would find teleconferencing and high-definition television (HDTV) as well as conventional voice traffic. Representative of delay-insensitive traffic are forms of voice traffic medical images and TV on demand. Among the new services there is emphasis on visual information since it consumes so much bandwidth.

With huge available bandwidth and multiple types of applications, the emphasis on required quality of service (QoS) and congestion control techniques has grown. Although initially strict QoS-based ATM became an almost sure concept, there was still a great debate on the wastage of bandwidth of ATM cells of almost 10% in terms of cell header. Because of the ease of implementation, Internet Engineering Task Force (IETF) proposals based on IntServ and DiffServ have become widely accepted standards for QoS-based IP networks. While Integrated Services architecture is based on individual flow-level QoS provision using the Resource Reservation Protocol (RSVP), Differentiated Services architecture is based on the QoS provision at the aggregate of flows. DiffServ architecture identifies four different types of aggregates, and defines their per hop behavior at each node, apart from metering, access rate control (e.g., using the "leaky bucket" mechanism, which we discuss in Chapter 7), and/or smoothing functions at the entry of the Differentiated Services network based on the service-level agreement between the user and the Differentiated Services (DS) network, prior to using the DS network. The typical mechanisms used in sharing bandwidth among the aggregates, are based on packetized generalised processor sharing (GPS) schemes and the active queue management to implement different levels of loss probabilities can be based on the random early detection (RED) scheme, which was originally proposed to improve congestion control mechanism associated with TCP (Transmission Control Protocol) dynamics. Most of the mechanisms studied use mathematical models for evaluating the performance for which our book presents the fundamentals.

It has been amply demonstrated that the Poisson model for the new classes of traffic yields misleading results; accordingly, a great deal of effort has been expended on models that capture the salient characteristics of the traffic generated by the new services. In Chapter 7, the *fluid flow model* approximates the discrete flow of digital traffic is by a fluidlike model. In Chapter 8, a more general tool, the matrix analytic technique, is studied.

## 1.4  COMPUTATIONAL TOOLS

The emphasis in the text is on getting numbers in order to evaluate the performance of a system. We will extensively use three different software tools to do this. The simplest is the *Excel* spreadsheet. The virtue of Excel is that it is easy to use and allows interaction. Simple formulas can be quickly evaluated, and curves can be drawn. The effect of changing parameters can be seen immediately.

The second tool we use is *Matlab*. It is more suited than Excel to complex computations, and it has much better graphical capability. We use Matlab extensively to simulate communications systems and techniques by means of Monte Carlo simulation. As its name may indicate, Matlab is particularly well suited to matrix operations. This tool was very valuable in the last three chapters of the text in handling complex matrix operations and in doing simulation.

The last tool that is used in the course is *Maple*. The particular strength of this tool is symbolic manipulation. In Chapter 6, for example, differentiation of complex functions was carried out. Maple was used in Chapter 8 to simplify very complex equations.

## FURTHER READING

### Material on Telecommunications Systems

Freeman, R. L., *Fundamentals of Telecommunications*, Wiley, New York, 1999.

"100 years of communications progress," *IEEE Commun. Soc. Mag.*, **22** (5), (May 1984).

Leon-Garcia, A., and I. Widjaja, *Communication Networks: Fundamental Concepts and Key Architectures*, McGraw-Hill, New York, 2000.

### Genesis of LANs

Abramson, N., "The ALOHA system—another alternative for computer communications," *1970 Fall Joint Computer Conf., AFIPS Conference Proceedings*, Vol. 37, 1970, pp. 281–285.

Farmer, W. D., and E. E. Newhall, "An experimental distributed switching system to handle bursty computer traffic," *Proc. ACM Symp., Problems in Optimization Data Communication Systems*, DP. 1–34, Pine Mountain, GA, 1969.

Metcalf, R. M., and D. R. Boggs, "Ethernet: Distributed packet switching for local computer networks," *Commun. ACM* **19**: 395–404 (July 1976).

Pierce, J., "How far can data loops go?" *IEEE Trans. Commun.* **COM-20**, 527–530 (June 1972).

### Queuing Texts of Historical Interest

Bear, D., *Principles of Telecommunication Traffic Engineering*, Peter Peregrinus Ltd, 1976

Saaty, T. L., *Elements of Queueing Theory*, McGraw-Hill, New York, 1961.

Syski, R., *Congestion Theory in Telephone Systems*, Oliver & Boyd, 1960.

### Queuing Texts with Tutorial Value

Allen, A. O., *Probability, Statistics and Queuing Theory*, Academic Press, New York, 1978.

Cox, D. R., and W. L. Smith, *Queues*, Methuen, London, 1961.

Gross, D., and C. M. Harris, *Fundamentals of Queueing Theory*, Wiley, New York, 1998.

Kleinrock, L., *Queueing Systems*, Vol. 1: *Theory*, Wiley, New York, 1975.

## Comprehensive Theoretical Treatments of Queuing Theory

Cohen, J. W., *The Single Server Queue*, North-Holland, 1968.

Takagi, H., *Queueing Analysis*: *A Foundation of Performance Evaluation*, Vols. 1–3, Elsevier Science Publishers B.V, 1993.

## Milestones in the Development of Queueing Theory

Brockmeyer, E., H. L. Halstrom, and A. Jensen, "The life and works of A. K. Erlang," *Trans. Danish Acad. Tech. Sci. ATS*, (2) (1948).

Erlang, A. K., "The theory of probabilities and telephone conversations," *Nyt Tidsskrift Matematik B*, **20**: 33–39 (1909).

Erlang, A. K., "Solution of some problems in the theory of probabilities of significance in automatic telephone exchanges," *Electroteknikeren* **13**: 5–13 (1917) [in English: *PO Electric. Eng. J.*, **10**: 189–197 (1917–1918)].

Feller, W., *An Introduction to Probability Theory and Its Applications*, Vol. 1, 3rd ed., Wiley, New York, 1968.

Kendall, D. G., "Some problems in the theory of queues," *J. Roy. Stat. Ser. B* **13**: 151–185, (1951).

Kendall, D. G., "Stochastic process occurring in the theory of queues and their analysis by the method of the imbedded Markov chain," *Ann. Math. Stat.*, **24**: 338–354 (1953).

Khinchin, A. Y., "Mathematical theory of stationary queues," *Mat. Sbornik* **39**: 73–84 (1932).

Kosten, L. *On Loss and Queueing Problems* (Dutch), thesis, Technicological Univ., Delft, The Netherlands, 1942.

Kosten, L., "On the validity of the Erlang and Engset loss formulae," *Het P. T. T. Bedjijf* (*Netherlands Post Office Journal*).

Pollaczek, F., "Uber eine Aufgab der Wahrscheinlichkeitstheorie," *I–II Math. Zeitschrift*. **32**: 64–100, 729–750 (1903).

## Packet Switching and the Internet

Baran, P., "On Distributed Communications Networks," RAND paper P-2626, Sept. 1962; also, *IEEE Trans. Commun. Sys.* **CS-12**(1): 1–9 (March 1964).

Cerf, V. G. and R. E. Kahn, "A Protocol for Packet Network Interconnection," *IEEE Trans. Comm. Techn.*, Vol. Comm-22, No. 5, 627–641, May 1974.

IEEE Communications Magazine, Issues of March and May 2002.

Kahn, R. E., "Resource-sharing computer networks," *in Computer Networks, a Tutorial*, M. Abrams, R. P. Blanc, and I. W. Cotton, eds., IEEE Press, 1978, pp. 5-8–5-18.

Kleinrock, L., *Communication Nets: Stochastic Message Flow and Delay*, McGraw-Hill, New York, 1964.

Leiner, B. M. et al., "A brief history of the Internet," http://www.isoc.org/internet-history/brief.html.

Roberts, L. "Multiple Computer Networks and Intercomputer Communications," ACM, Gatlinburg Conf., October 1967.

Roberts, L. G., "Data by the packet," *IEEE Spectrum* **11**(2): 46–51 (Feb. 1974).

Proceedings of the IEEE Special Issue on Packet Communications, Vol. 66, No. 11, Nov. 1978.

**The Modern Network**

Leland, W. E. et. al., "On the self-similar nature of Ethernet traffic (extended version)," *IEEE/ACM Trans. Network.*, 1–15 (Feb. 1994).

**On Modern Internet Services**

Blake, S., D. Black, M. Carlson, E. Davies, Z. Wang, and W. Weiss, *An Architecture for Differentiated Services*, RFC 2475, Dec. 1998.

Braden, R., D. Clark, and S. Shenker, *Integrated Services in the Internet Architecture: An Overview*, RFC 1633, June 1994.

Braden, R., L. Zhang, S. Berson, S. Herzog, and S. Jamin, *Resource Reservation Protocol (RSVP)—Version 1 Functional Specification*, RFC 2205, Sept. 1997.