

Technology and Networks

On Saturday mornings, Henri used to sleep until noon and only after a late breakfast think about amusing evening activities. This time, he takes a brisk walk already before nine o'clock and continues his project with a fresh mind. Although he is really inspired by some of the technical parts of his novel application, he knows that in reality most innovations will never be successful regardless of their technical merits. Thus, Henri decides on Friday evening that he will spend Saturday designing the technical architecture of the application in a way that it supports both the creation of a pleasant user experience and the dissemination of the product.

He starts with an idea that most of the content of the advice given by Flourishator shall be created by users rather than by himself or by any paid personnel. Surely, there must be a collection of pre-defined advice that covers the most reasonable combinations of experienced and desired emotions. Some of the advice might be rather funny than serious. Henri may easily devise some humorous advices to convert love to frustration, if someone happens to want that kind of development. Still, his prime aim is to give genuinely beneficial advice.

But how could he know what advice is good and what is bad? Fortunately, Henri had read a couple of weeks ago "The how of happiness" which provides excellent recommendations based on scientific research. He can use that as a starting point to sketch seven reasonable pieces of advice to be used in Flourishator.

Still Henri is aware of his limited creativity and credibility as a well-being guru. It is obvious for him that he needs a community that is interested in the effort of improving people's lives. That sounds like a bold but yet an achievable goal. It means also that part of the functions of the application have to be implemented in servers located somewhere in the network. The servers provide an interface towards the application running in the user devices, an interface to advice developers, and an interface for the operation and management of the system. As far as Henri is able to judge there should not arise any significant performance problems with the servers or in the connection between users and servers.

As a software engineer, Henri is fascinated by the idea that some pieces of advice could be created artificially. But, no, he is also aware of the insurmountable

ble obstacles of artificial intelligence. Maybe something funny could be created by combining parts of already available pieces of advice. That might serve as an entertaining part of the Flourishator. However, the more serious advices shall be written and selected by human beings.

Ideally, a developer community will take the responsibility of creating new advices. Of course, Henri judges, there must be a strict set of rules that makes it certain that only reasonable advices will be given to those users that want honest advices. What could serve as a reliable indicator that an advice is valuable?

Henri is inclined to believe that a simple opinion scheme is the most feasible approach, just like/dislike. As an extra feature, users might recommend special advices for their friends. That might be an important feature for the dissemination of Flourishator.

Altogether, Henri notices, this means four levels of satisfaction: *dislike*, *no reaction*, *like*, and *recommend*. If there are enough users, that would create a lot of information for a thorough analysis. The most acute problem might occur at the beginning, because then there are only few users, few advices, and not much information to be exploited. In the first phase, only some curious people might be inclined to become excited about novel things.

Henri remembers the concepts of innovators and early adopters. He has a vague idea that a deliberate effort should be done to inspire a big enough group of innovators that can also be used to develop the product. Innovators, so the authors say, are often more than willing to participate to the development of the product. The risk at that stage of dissemination process is, a lecturer has warned, that innovators want to develop the product for themselves and the result will likely be too difficult for the majority of potential users. Henri notices that he needs to be careful with this, particularly as he considers himself an innovator as well.

How could he solve this dilemma? Henri thinks that he needs to rely on objective measures that give a reliable indication both about the acceptability of the advices and the interest in the product itself in several user segments. If possible, he should try to categorize the users to three groups: innovators, people having strong influence on other people, and people that usually follows what others have already done.

Henri's plan is to handle every innovator as a valuable assistant in the development task. If they also happen to invent emotional advices offered to others, that would be nice, but it might be that those advices are not the most suitable ones for majority of people. He would trust Irene more on that area of emotional advices. Then the most important job of the early adopters is to recommend the product to other people.

His intuition is that a set of unique advices has to found: they need to be creative, funny and serious at the same time, even at the risk that someone may consider some of the advices unpleasant, scary, or bizarre. Emotions must be induced, mostly positive but a slight flavor of negativity might be unavoidable

because of the differences in people's emotional reactions. Also, some scary movies are popular. Should the users be profiled to avoid undesirable reactions? All kinds of questions appear in Henri's mind.

Anyway, a clear plan for coping with the challenges at each of the three stages would seem professional in the business plan needed in the assignment report. Whether or not the plan will work in reality is another matter—only time will tell.

Selection of topics

The specific target group for this book contains the students in the International Master Programme of Communications Ecosystem at Aalto University. A majority of those students have already studied networking technology for a couple of years. Still, I assume that many of you only have a superficial knowledge about communications networks and technology in general. The aim of this chapter is to serve both groups. This chapter starts with a general discussion about engineering and technology, particularly from the viewpoint of the seven rules and 4-level structure of metrics consisting of capability, performance, efficiency, and worth.

The remaining part of the chapter concentrates on the quality and control aspects of networking services. In that part, I go deeper than in most of the other topics discussed in this book. But why do we need to be, in the first place, concerned with quality and control? The justification is related to my professional background. I spent almost twenty years to design controlling systems that could be used to improve the quality of communications services. The main result of those studies—in addition to a doctoral thesis, a textbook, and several patents—was a deep conviction that in order to achieve anything relevant in that field we have to consider the complex interplay between technology, human needs, and economics. This book is an outcome of a ten years effort to understand that complex matter called the communications ecosystem.

That kind of specific history does not imply that quality and control should have any specific role in the discussion of technical matters in general or networking in particular. There are, however, other reasons to consider Quality of Service (QoS) an important topic. First of all, service quality is the main way through which the properties of network technologies are experienced by the users of services. In the terminology of this book, QoS is closely linked to the technical methods that are used inside the networks to provide special types of service quality for users. To put it simply, QoS is in the core of communications ecosystems as a mediator between technical, human, and economic aspects. Another similar topic that serves as a mediator between separate fields, usability, is discussed in Chapter U.

Numerous other branches of technology, of course, have a momentous effect on the characteristics of communications services, including optical, semiconductor, radio, and router technologies. In many respects, these technologies have been much more influential for the development of communications ecosystem than any QoS method or mechanism. We may still argue that after the initial phase new versions of all these technologies primarily offer higher

performance with lower cost, without complex interactions with human behavior and business models. When complex human or business interactions emerge, the issues to be solved typically are somehow related to either service quality or control mechanisms, that is, how resources are divided between different actors in the ecosystem.

Terms

There are a huge number of specific technical terms, even in the limited field of communications. Many standardization organizations have published extensive dictionaries. This small selection does not substitute those dictionaries; instead, it defines the basic technical terms used in this book as consistently as possible. A communications ecosystem expert (CEE) shall master the key networking terminology (which does not mean, however, that the definition of every term shall be memorized).

Many distinctions are important to notice and remember. For instance, *communications* (in plural form) is used in the specific meaning while *communication* (in singular form, discussed in Chapter S) has meaning that is more general. Note also the essential difference between *invention* and *innovation* (discussed in Chapter C). Finally, the Internet refers to the one and only global network, while internet refers more generally to a packet network. Furthermore, Internet protocols (IP) can be used in many kinds of networks and devices.

The selected terms to describe the main topics of this chapter are:

capability: measure of the ability of an entity or system to achieve its objectives, especially in relation to its overall mission,

communications: the science and technology of communicating, especially by electronic means,

data: information in numerical form,

engineering: the use of science in the design, planning, construction, and maintenance of buildings, machines, and other manufactured things,

Internet: a global information network that consists of a large number of smaller internets,

invention: the idea of a new product, or a new method of producing an existing product,

link: a physical connection between two network nodes,

network: a collection of nodes and links that provide connections between access points,

node: a device attached to a network with the capability to make connections to other devices,

packet: an information unit that contains enough information to transmit it through a network,

performance: a quantitative attribute of a system that describes how well the system is able to fulfill its predefined purpose,

technology: the entire collection of devices and engineering practices available to a culture,

throughput: the amount of data transmitted or processed over a given period, and

topology: the physical arrangement of network nodes and links within an organization's networking structure.

In addition, the following terms are defined in Glossary:

access point	connection oriented	guaranteed service	RFC
availability	connectionless	implementation	router
bandwidth	coverage	infrastructure	scalability
base station	deep craft	internet	signal
best effort	delay	interoperability	software
billing	device	IT	technology push
bit rate	differentiated services	production	telecommunication
blocking	digital	QoS	traffic
channel	equipment	reliability	transmission
compatibility	Erlang	requirement	utilization

Engineering

“Engineers like problems they can solve,” as Walter Vincenti have said according to W. Brian Arthur (2009, p. 15) referring to the lack scientific foundation in engineering work. What kinds of problems can engineers then solve? Technical problems, of course. It is unfair to require that an ordinary engineer shall be able solve scientific, business, or social problems.

A typical engineering activity means a project in which known methods, principles, and technical building blocks are used to construct a new version of an existing technology. Brian Arthur calls this *standard engineering* (2009, p. 91). Cases in which an invention creates a new technology that is not based on existing technologies are extremely rare (the compass might be an example, although the stages before the invention of the compass are uncertain). For instance, the invention of self-propelled flight made by the Wright brothers was heavily based on a long course of developing necessary technologies. What the brothers achieved was that they were able to develop certain key technologies necessary for successful flight, like an efficient propeller and managing the stability of the flight. They solved many tough engineering problems, but they were not really trying to solve any business or social problem, or satisfy any business or social need.

In contrast, I would argue that a communications ecosystem expert has to be able to, at least to some extent, address both business and social problems related to communications service, by using different skills and knowledge. This kind of multi-purpose capability is hard to acquire. You may also think that when a high enough skill level is acquired, that skill

becomes precious. After a long enough period of analyzing ecosystems, CEE may acquire something that can be called *deep craft* in an area that connects communications networks and the business of service providers. Deep craft is not just knowledge about something but knowing what is likely to work and what is not likely to work (see Arthur 2009, p. 159). In the framework of this book, deep craft refers to the capability to select the appropriate tools and models to evaluate the business potential of new communications products.

Deep craft of a certain area of an ecosystem may unfortunately have a negative influence on the ability to look the ecosystem from other perspectives. This may lead to a situation in which an organization, like a device vendor with a large number of technology experts, tends to overrate the value of technical capabilities. This together with the relentless demand for new business models may result in:

technology push: technology development that is driven by the ideas or capabilities of the developing organization in the absence of any specific need that customers might have.

The responsibility of a CEE is to be aware of this tendency and remember the first rule for CEEs, the Rule of human benefit.

In the current form of this chapter, the networking aspects might be somewhat overrated, because communications ecosystems include many other fields of technology that could be discussed in more detail. For instance, the performance and quality of software technology is in many cases a critical issue that may determine whether an innovation will be successful or not. The emphasis of networking is partly due to my background in networking and particularly in the area of QoS. The justification for this approach, in addition to the familiarity of the area for the author, is that service quality forms a strong link between technical matters—like the efficient sending and receiving of electromagnetic waves—and human aspects. How the limited capacity of radio channels is divided between users or customers is the issue dealt by QoS engineers. Note that *how* means both how the division of resources is technically realized and what rules shall be used to define the actual division of resources.

The first question is clearly technical and may require deep understanding about the properties of radio and semiconductor technologies. The other question is essentially non-technical, particularly if it is considered from a user or customer perspective. Both issues, technical and non-technical, are relevant. The task of QoS engineers is to find solutions that are acceptable and efficient both from technical and non-technical perspectives. Thus the studies on Quality of Service, as defined in the context of communications networks, have a similar, but more limited role as Communications Ecosystem studies. Quality of Services has both economic and human aspects as demonstrated in later in this chapter.

Thanks to my long background in QoS domain my mind seems to have an extensive amount of resources dedicated for considering questions related to QoS (remember that a great majority of what our minds do is automated). That is sometimes a very useful ability because I can use intuition to answer questions that are difficult to analyze by means of formal reasoning. There also is another side to this capability: it is hard for me to think about QoS issues from any other viewpoint than what the unconscious processes in my mind automatically generate. The long period of researching quality of service and network performance

means that efficient but stubborn automatized subconscious processes have evolved in my mind.

For instance, just before writing about the traffic control methods in this chapter, I wrote about the pros and cons of marketing discussed in Chapter C. A change from marketing to technology meant a clear shift in my mindset as well. Performance and efficiency are the key attributes of the technical QoS mindset. I can almost immediately imagine network elements, links and nodes, and something quite abstract, called traffic, going through the network. The objective of a network engineer is to transmit as much traffic as possible through the network: the more the better. There are, of course, many constraints including QoS requirements, but even they are usually presented in technical terms: packet loss ratio, delay, bit rate, reliability, and so on. Some results of this engineering mindset are presented later in this chapter.

Thus, if I need to define one key objective for CEE, it is to make technology more human and more social. As Donald A. Norman has expressed it (2011, p. 117):

“Machines need to show consideration for the people with whom they interact, understand their point of view, and above all, communicate so that everyone understands what is happening.”

Engineers, finally, are responsible for the development of more human and sociable technology. One may argue that the problems of modern technology and economy do not have any clear optimal solution from any given viewpoint. Systems are just too complex. This hardly is a special property of modern technology. It is an inevitable state of affairs with any ecosystem: biological, economic, social, technical, or a combination of these. Ecosystem problems are ambiguous by nature, and even when the problem can be defined clearly, the ecosystem itself defies any accurate analysis. The point is not so much to prove that certain technology is more suitable for human beings and for societies, but to develop processes that lead to machines and systems more suitable for human needs. Norman (2011, p. 115) calls this *sociable design*, which includes asking for compliance and tolerance (used as social not technical terms). Note that the effect of terminological selections can be profound, as the results of many psychological studies have shown.

Thus, the recommendation of this book for CEE is to use terms that are solidly human and cannot be easily interpreted as technical terms or translated to technical parameters. Your terms shall work appropriately even if your mindset is otherwise technical. *Quality* does not work well in that respect. Emotions shown in Figure H.3 may serve better that purpose better.

Metrics for networks

The objective of communications networks is to transport information between separate entities by using signals. Thus, the most obvious metric to assess communications network is throughput: the more information is transmitted through a link, node, or network, the better. Then if we consider the communications purely on the technical level, we tend to replace information by data. Note, however, that there is a crucial difference between the usage of terms *information* and *data*. Information is defined as “a difference that makes a difference”

whereas data is information in numerical form. The somewhat cryptic definition of information proposed by Gregory Bateson stresses the active nature of information. Information is something that can be used to make a decision by an agent. We may also ask: is there any information without conscious beings that could use the information for some purpose? I tend to think that there is no information without the possibility of using it by someone. Thus, the criterion for successfulness of communications service is that information can be finally utilized by someone. Information is hard to measure if the true usefulness shall be taken into account. In practice, the quantity of communication is measured by the amount of information in numerical form, that is, by the amount of data. As a result, a typical unit for throughput is bit per second.

The importance of throughput is illustrated by the names of network types, such as broadband network or Gigabit Ethernet (instead of quick network or 10-ms-delay-network). Then there are additional parameters that are important depending on the purpose of the network. We can call these parameters by a general name of *capability*. Examples of network capabilities are: a network is able to transmit data within a certain delay limit (say, 10 ms), another network is able realize differentiation rules between traffic flows, while a third network is able support mobility of customer devices. Networks also have more general properties, including scalability, flexibility, and compatibility with other networks. However, we may consider those properties secondary in the sense they are not typically relevant matters for the users of the network if the throughput is high enough.

We must keep in mind that *capability* and *performance* are different categories: a capability is a feature of the network or a part of the network, whereas performance depends both on the network capabilities and on the load and operational processes. For instance, the core network illustrated in Figure T.2 is required to offer service with negligible loss of data. The network and device vendors may realize various mechanisms to achieve that goal, but it is up to the operator (and to some extent, users) how these mechanisms are used.

Thus, the core network capabilities do not solely define the network performance. Performance was defined in the previous section as a set of quantitative attributes of a system that describes how well the system is able to fulfill its predefined purpose. Is performance then a metric? I would say no, because performance typically consists of several attributes, which does not allow us to put systems in an unequivocal order of performance. Throughput and end-to-end delay are examples of appropriate metrics, although even they are difficult to define exactly in the case of a large network. This ambiguousness of performance might be circumvented by introducing a formula that combines all performance attributes to one parameter that describes how well the system fulfills its predefined purpose in general. The simplest solution is to list all performance requirements and then count the number of fulfilled requirements: the bigger number the better. This, indeed, is often the solution applied in real cases. Note also that the same problem of multidimensionality concerns both capability and performance.

We can continue the terminological system towards business and human aspects by introducing two additional concepts: *efficiency* (a concept in the area of economics) and *worth* (a concept related to social systems).

Efficiency can be defined as the extent to which a resource is used for the intended purpose. Use of resource indicates that there is a cost to be paid to achieve the purpose. In reality,

all essential resources, such as money, network capability, time, and human effort are limited. Correspondingly, the purpose of an actor can be to create wealth, to transmit packets, or increase user satisfaction. In many cases efficiency is more appropriate metric than corresponding performance metric, because efficiency takes in account cost aspects, not only what has been achieved in positive terms.

Finally, if we define efficiency straightforwardly as the positive output to cost ratio, we must also be aware of the difference between *outcome* and *value*. *outcome* is something easily observable and measurable (like money) while *value* is something that describes the benefits of something from the perspective of a human being (see Chapter E). I use the term *worth* when the result is defined primarily on the scale of human benefit (or eudemony) and all relevant cost of factors are taken into account. In the end, something is worth doing if it serves the fundamental purpose of life.

In reality, metric is often directly related to the concepts of performance and efficiency. The metric can be the same as the most obvious performance parameter of the system or activity to be assessed. For instance, the metric to be used to assess universities is the production of students and scientific publications: the more and the higher quality, the better university. That is a natural approach, but still, it is a performance rather than efficiency metric, and certainly, it is not a worth metric. We might ask many relevant questions, like: “Should the assessment of universities consider the share of high quality “student material” as a cost, because excellent students are in short supply in any society, or as an indication of the quality of the university?” or “Does a high salary level of key personnel mean high cost or is it an indication of high performance (because it means that the university can compete efficiently on the academic market)?”

Table T.1: Fundamental terms to describe the characteristics of a system.

<i>Term</i>	<i>Description</i>	<i>Depends on</i>	<i>Measurability</i>	<i>Examples in networking context</i>
<i>Capability</i>	Property of a system in isolation	System itself	Straight-forward	Traffic scheduling algorithm
<i>Performance</i>	Property of a system in a pre-defined environment	System and load	Moderate	Packet loss ratio, end-to-end delay
<i>Efficiency</i>	Desirable outcome relative to a pre-defined input	System, load, and cost	From moderate to difficult	Throughput divided by total cost
<i>Worth</i>	Value or benefit compared to all relevant cost aspects	System, load, price, customer, and society	Difficult, or even impossible	Long-term benefits of a service divided by all sacrifices including time and money

Let us take a further example to illustrate the difference between capability, performance, efficiency, and worth. You may consider it desirable to own a car with the capability to

accelerate from 0 to 100 km/h in less than 8 seconds and to keep a maximum speed of 250 km/h. Still, those are mere capabilities in the context of this framework whereas performance is measured in real situations. The maximum speed does not really have much practical use unless you are driving in Germany.

Let us consider a case in which an engineer wants to drive to his summer cottage located 500 km away from his home. A natural way to assess the performance of the car is to measure the time it takes to drive from home to the cottage. As an engineer, he maximizes the average speed by driving whenever possible 25 km/h above the speed limit. Moreover, he had estimated that if he drives faster, there is too high a probability to lose his driving license. He exploits, naturally, the acceleration capability of the car to the fullest. As a result, he reaches his cottage in 4 hours and 22 minutes. He will be proud of this great achievement, and his self-confidence is improved significantly, for a while.

What could be the corresponding efficiency metric? It might be that during a weekend trip the engineer notices that the car consumes enormous amount of gasoline because of his aggressive driving behavior. Gasoline is expensive, about 1.50 Euros per liter, and he has other needs that require money. Thus when driving towards home he starts to make an efficiency analysis. First he assess that his sport utility vehicle consumes about 14 liters per 100 kilometers with aggressive driving habits. That means 195 Euros per weekend. Then he decides to make an experiment on the next weekend. He will not exceed speeding limits or the speed of 100 km/h at motorways, he will not accelerate unless really necessary, and he will sometimes drive behind trucks and busses to save gasoline. As a result, the driving time is increased up to 5 hours and 46 minutes. To his surprise, gasoline consumption is decreased down to 7 liters per 100 kilometers.

A straightforward calculation reveals that during a weekend fast driving saves 2 hours and 48 minutes, but costs 95 Euros (plus possible speeding tickets and the higher probability of traffic accidents). That means that every saved hour costs about 37 Euros. As an engineer, he can make a more extensive efficiency analysis by including all relevant cost factors and the value of his own experiences during the driving. The result of the evaluation is that the optimum driving time is about 5 hours and 25 minutes with an average gasoline consumption of 7.5 liters per 100 kilometers. However, that is not the final stage of the analysis.

The engineer should also consider the question: is the driving to the cottage worth all the sacrifices he and other people need to pay? This assessment goes beyond any direct economic calculation and depends on many details of the situation. Why does he maintain his own cottage? Do the regular visits to the cottage really serve his most fundamental needs? Does the habit of driving to the cottage affect his social interactions in positive or negative way? These are essentially different questions than what were asked in the beginning of this example. Note also that the mindset of the engineer changes during the process from appreciating pure capabilities to assessing the worthiness of his habits in general.

Some of you may question the relevance of this example in the context of this book. You may still appreciate the fact that the example is able to create feelings because it describes a realistic situation. In general, I encourage you to use examples from other fields of life to understand the essence of a phenomenon in a specific field, like networking technology. This important method might be called *structural portability* (the term was proposed to me by Esa

Saarinen). Note particularly that structural portability can be used successfully in all directions and between all areas of life. The architecture of communication networks resembles the structure of body. For instance, it is not just an accident that cellular networks consist of cells. The resemblance extends beyond simple analogies towards complex phenomena, like, on the one hand, the interaction between instances of applications competing for limited resources in the Internet by means of Transmission Control Protocol (TCP) and, on the other hand, the members of a society competing for limited resources of health services through various rules and conventions.

Quality in technical context

The standpoint of this book is that quality should primarily refer to something that is closely related to human experience. Thus, it might be better to use terms capability, performance, or efficiency in a purely technical context. An additional trouble is that if quality is used as the only attribute of something (e.g., in the case of QoS), then it automatically appears to be a sign of positive characteristics (which is anyway a proper usage of *quality* as a common word). For marketing purposes, that kind of association certainly is advantageous. However, from user perspective the quality of a network with QoS capabilities might be low or high depending on how the QoS mechanisms are used and how the network is operated.

Yet, we need to accept the fact that in the field of communications Quality of Service is used to refer to technical parameters of the service. The solution to this controversy in this book is that acronym QoS is used in the way defined by International Telecommunication Union (ITU). In contrast, when quality is written with a small initial letter (the quality of a device) it refers to the experience of a human being.

Furthermore, quality may refer to the intentions and capabilities of the actor responsible for the production of the entity, in addition to the properties of the entity itself. If we consider the quality of a mobile phone, we think not only about the tangible phone in our hand but also the vendor that has been responsible for the design and manufacturing of the product (this can be interpreted as a placebo effect). Consequently, in case of commercial products, quality comes close to the term *brand*. If a vendor has a high brand value, we expect that the products of the brand are of high quality or that they are particularly useful. In that case, a customer expects that all parts of the product are carefully designed and that the product is able to satisfy the needs that it is planned to satisfy. In general, quality stands for something that is felt and experienced rather than something that is measurable in any objective way.

The following sections discuss three topics that are closely related to quality of services and that affect networking business in general: traffic control principles used in the network, network dimensioning, and the reliability of network services. Figure T.1 outlines the position of these activities in the overall business of communications service providers. Note that in addition to the terms brought up in the figure, numerous other issues affect the outcome from a business perspective, that is, the profit of the service provider.

All these activities affect the QoS of network services and, thus, user experience and customer satisfaction. More resources, and improved control and reliability likely have positive effects on the business. However, all activities also create negative effects, both monetary cost

and other types of cost. The main objective of this section is to give an overview of the relationship between these activities and the profit of service provider.

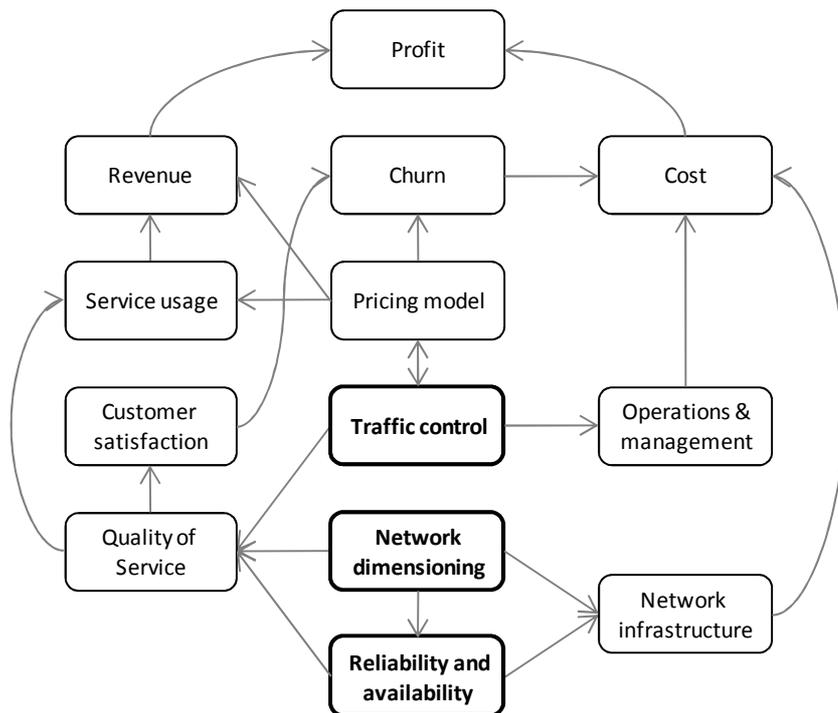


Figure T.1: The position of traffic control, reliability and network dimensioning in the business of a service provider.

Traffic control

Why should we consider traffic control in an introductory book at all? The main reason is the fundamental difference between the principles of the Internet and networks stemming from the field of telecommunication. The differences are not only important technically, but they are also reflected in the business models adopted by Internet providers and mobile operators. The differing principles pose a critical problem both when Internet services are offered through mobile networks and when telephone services are offered through the Internet (voice over IP or VoIP). Furthermore, the appeal for increasing control has generic importance beyond the fields of communications.

Let us first consider the question: what is the most reasonable way (if any) to control the traffic transmitted through a network? Figure T.2 illustrates a communications network that can be divided roughly into two parts: core network and access network. In this discussion, we assume that both the core and access networks are managed and operated by a single network operator. Because customer devices typically are owned by other entities, the network operator

has more limited ability to control the behavior of those devices than the core network. Moreover, there are human beings, users and customers that have their own intentions and goals and make their own decisions. The main aspects to be discussed in this section are the means to control the traffic load in the core network, the reasonability of these means in different situations.

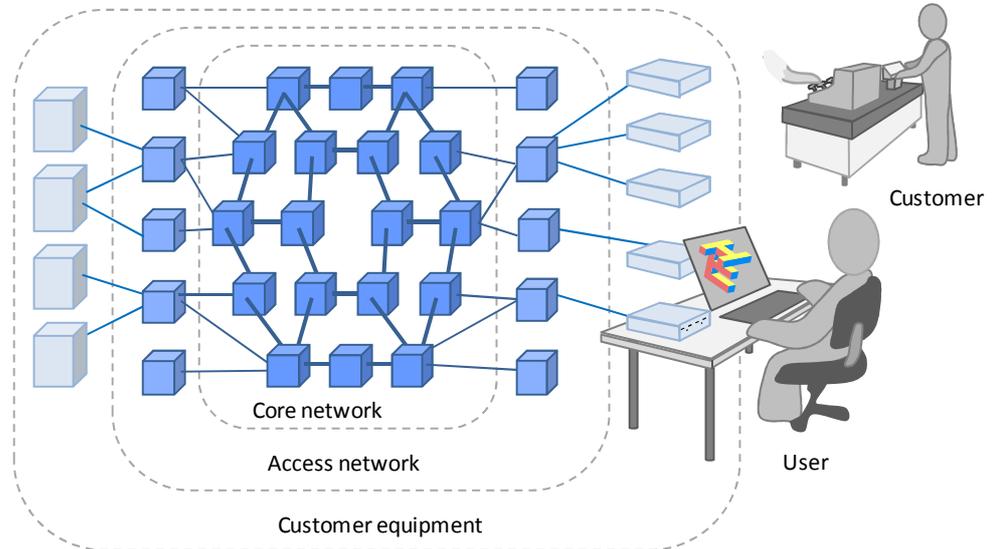


Figure T.2: Network, customer equipment, user, and customer.

The current Internet relies largely on the abundance of network capacity and on the best effort service model. Although the building of excessive network capacity compared to average demand seemingly wastes resources, under some conditions it is the most reasonable approach. The main advantage of this “throw bandwidth” –model is that there is no need for complex control of traffic in any part of the network, not even at the user side. In an extreme case, users are allowed to surf the net, watch streaming videos, and download huge files without thinking about any consequences at all. Similarly, network operators have to care only about the amount of network capacity, without any special mechanisms needed to regulate traffic and users. Even the business part of the ecosystem is the simplest possible: connect customers to the network, send a fixed bill each month, and let the users fulfill their needs in whatsoever manner they want and invent.

Because this Internet model has proved to be successful, my recommendation is to take it as the reference model for any communications network. A more complex model in any part of the ecosystem has to be justified by careful analysis, instead of taking a complex system as the starting point and assuming that the improved capabilities necessarily provide significant advantages. However, there are situations in which additional complexity is justified.

But before going on to the logic of analysis, let us start with sketching the available mechanisms to control the amount of traffic in the network. Here we assume that the network

operator believes (for instance, based on the advice of external consultants) that traffic must be somehow controlled to avoid unnecessary capacity, additional expenses, and possible problems in service quality. The traffic control approaches can be classified along several dimensions, including:

- proactive vs. reactive mechanisms,
- place of realization: core network, access network, or customer equipment,
- granularity of control: packet, flow, user, or customer,
- time scale of control: from milliseconds to months, and
- connection-oriented vs. connectionless mechanisms.

The first dimension is always crucial: the operator may try to implement mechanisms that prevent all problems and conflicts beforehand, or it may try to implement mechanisms that react as efficiently as possible when problems or conflicts occur. Someone may automatically think that in principle proactive is better than reactive. Sometimes that surely is the case, but in general, we shall not jump to any hasty conclusion without defining carefully the criterion for better and worse.

Another key dimension is the location in which the control actions are primarily performed. Many control mechanisms require co-ordination between different entities located in different parts of the ecosystem. For instance, connection admission control is usually done together with a suitable pricing scheme, and Differentiated services calls for cooperation between access and core nodes. Some control mechanisms are presented in Table T.2. Note that similar methods can be used in other parts of the network or even in other parts of the communications ecosystem.

Table T.2: Examples of control mechanisms to limit the load in core network.

	<i>Core network</i>	<i>Access network</i>	<i>Customer equipment</i>	<i>User</i>	<i>Customer</i>
<i>Proactive</i>	Connection admission control	Access rate control	Reservation requests based on Integrated services	Usage-based pricing	Limitation of the number of customers
<i>Reactive</i>	Packet dropping	Marking based on differentiated services	Transmission Control Protocol (TCP)	Congestion pricing	Sporadic unavailability of service based on service class and load

Proactive control methods

Connection admission control is a method that limits the network load by restricting the number of and the capacity reserved for connections going through the network based on the load situation in the network. Traditional telephone networks have always relied on admission control. There are both technical and psychological reasons for this reliance. Some electro-mechanical components in telephone exchanges had to be reserved for the whole duration of a voice call, because their speed did not allow time-division multiplexing or shifting resources continuously between calls. Similarly, a wire (or a channel) had to be reserved for the whole duration of a call. In this respect, current technology is much more flexible and allows numerous ways to share resources between calls while still providing acceptable user experience even when the number of simultaneous calls fluctuates.

In contrast to technical constraints, the psychological motives for admission control in case of voice calls still remain. An interruption of a voice call due to a technical reason is often a disconcerting experience. It may seem that this obvious fact provides a sufficient justification for strict admission control in any network in which voice calls represent significant part of traffic or business.

Admission control was an extensively studied topic twenty years ago in the context of ATM networks leading both to theoretical formulations and to some practical insight. The fundamental shortcoming of admission control is that it does not work well when the majority of traffic consists of highly variable and dynamic data flows. The assumption before and partly during the development of ATM technology was that majority of traffic in multipurpose networks will have definite performance requirements. Video applications were assumed to be very sensitive of any loss of information or delay variation. The assumption was that there always will be a clearly identified customer and commercial service available between certain points of the network, and ATM would be needed provide an appropriate connection between the customer and the server, or between two customers. Sure enough, there was no concept of web browsing, no experience of YouTube videos, and no problem with peer to peer applications at that time of development or, rather, evolution.

But even if admission control might be technically implementable, and perhaps applicable to all those uses, admission control is difficult to justify without a corresponding pricing scheme. If something is free of charge, every customer may want to just reserve an unlimited amount of resources everywhere. That sort of behavior would certainly destroy all the possible gains obtainable by admission control.

Furthermore, usage based pricing is extremely difficult to design in a way that it would be suitable with highly variable value per used resources (see, e.g., Odlyzko 2004). Note that the willingness to pay for an important text message could easily be several Euros per kilobyte while the willingness to pay for the transmission of a movie through a network likely is about one Euro per gigabyte (GB). The difference is huge. Usage based pricing might be used as an additional way to control the network load, but surely not as a primary way when the same network service is used by all kinds of applications and for all kinds of uses. More examples are presented in Table T.3.

Table T.3: Willingness to pay for transmission through network.

	<i>Unit</i>	<i>Bit rate</i>	<i>Amount of data</i>	<i>Acceptable price</i> € / <i>unit</i>	<i>Acceptable price</i> € / GB
<i>Text message</i>	Message		200 bytes	0.10	500 000
<i>Voice</i>	Call, 3 min	10 kbit/s	0.2 MB	0.20	890
<i>VoIP</i>	Call, 10 min	100 kbit/s	7.5 MB	0.50	67
<i>Video clip</i>	Video, 5 min	500 kbit/s	19 MB	0.20	11
<i>Movie</i>	Movie, 2 h	2 Mbit/s	2 GB	3.00	1.7
<i>Data</i>	Data for a month	10 Mbit/s	50 GB	20	0.4

Note that it is a different question to ask, “How much is a person willing to pay for the possibility to watch a movie?” than to ask, “How much is a person willing to pay for transmitting the movie through a network?” Three Euros might be an acceptable price for cost savings compared to renting a DVD from a store even if movie tickets for two persons may cost 20 Euros.

An operator might be only concerned with the core network load, because it has guaranteed perfect network performance quality through the core network for key customers: no packet losses, delay always below a given limit, and very high availability of the service. The operator can accomplish those technical objectives by limiting the access rate for each customer without any other admission control method. This has been the typical solution in case of leased-line services. The average load level may remain low, below 10 percent. If the core network capacity is not a major cost source, the approach is still reasonable regardless of the low utilization of the network.

As to the Integrated service approach in Table T.2 (Customer equipment column) an alternative solution would be a procedure running in the customer equipment that sends a number of probing packets through the network and then based on the collected information about lost and delayed packets makes a decision whether or not to establish a connection through the network. In a more advanced method, the network would give specific information about the momentary state of the network. However, without cooperation between network nodes and customer equipment and without proper incentives there is no reason to assume that customer equipment will make reasonable decisions from the network operator viewpoint. Any controlled synchronization between the procedures running customer and network equipment might be too complex to make the system useful in reality although in theory the system would provide some advantages.

Finally, the ultimate control method is to limit the number of customers to avoid any excessive load in the network. Whether or not this approach makes any business sense depends on pricing, competition, and alternative options. Because of the substantial fixed (both operational and capital) expenditures, it is not likely that an operator should limit the total number of customers due to the lack of resources in the core network except in exceptional situations. In mobile networks with limited resources at the access, part of the network

the effect of customer base is different and it might be reasonable to limit the total amount of customer to preserve high satisfaction with current customers.

Reactive control methods

If strict admission control for all connections through the network is a kind of ultimate controlling method, then just dropping packets when the next link's buffer full is the other ultimate solution. It is curious that the ATM community tried to combine these two congestion control methods in the same network by means of constant bit rate and unspecified bit rate services. The problem of integration is not primarily a technical one, but related to the business model. It is relatively easy to put a sufficient number of buffers in parallel and design an algorithm that allocates the link capacity among the buffers in any desirable way. Furthermore, the admission control for constant bit rate connections is fairly straightforward to implement.

Surely, the technical and design problems become more difficult when admission control is applied to connections with variable bit rates, but even then both the allocation of buffer and link capacity and admission control can be implemented efficiently under certain conditions. What are those conditions?

First of all, most of the admission control and the resource allocation algorithms require that the properties of the connections are known in advance. Particularly someone must provide a proper estimation of the average bit rate over the whole duration of the session. If this average value is unknown, the network is unable to make any proper reservation. However, it does not make much sense to ask every time when a customer wants to make a connection: what will be the most likely bit rate during the session? Most probably, the customer will not know any credible answer, and even less, he would be willing to accept any punishment after an inaccurate guess.

Secondly, a reservation scheme has to come with suitable pricing that controls the demand for reservations. On the other hand, it is difficult to justify a business model in which web browsing is charged by usage-based pricing. Nevertheless, there were some serious efforts to combine admission-based service with an unspecified service model, also on the business level. An acceptable pricing structure was perhaps the most difficult part, and, as far as I can assess, remained an unsolved problem. The problem has been as difficult for the telecommunication sector (how to add unspecified service to a network based on admission control) as for the Internet sector (how to add guaranteed service to a network based on the best effort model). Neither of those integration efforts have been particularly successful.

Maybe the most logical way of combining the principles of admission control and best effort service has been the Differentiated Services model designed by Internet Engineering Task Force (IETF). Differentiated Services defines technical instruments for network operators to build a service model that lies somewhere in the middle of the guaranteed and best effort models. The essence of the Differentiated Services model can be summarized as follows:

- Packets are marked by 6 bits (differentiated services code points) that can be used to construct service classes with different treatments of IP packets inside the network.
- The original idea was to construct relative preferences rather than absolute quality requirements, although there seems to be still some disagreement about this.
- The main dimensions for packet handling are delay and packet discarding.

It is possible to construct a system in which one bit defines whether the packet has strict delay requirement or not, while 2 or 3 bits are used to define at most 8 drop preference levels. When a packet arrives at the scheduling system serving an outgoing link or other limited resource, the system first defines whether the packet will be accepted or discarded immediately. The decision is made based on the load situation of the link and the buffer and on the drop preference of the packet: the higher the load the more important the packet must be in order to be accepted to the buffering system.

The accepted packets are then delivered to two separate buffers based on the delay marking of the packet. The packets in the strict-delay queue are always sent as fast as possible while the packets in the other queue are sent forward only when the strict-delay queue is empty. This type of construction forms two orthogonal dimensions, one for delay and another for drop preferences. Yet this is just one possibility for the technical construction of differentiated services. The most essential matter both from the user and from the business perspective is the rules that define how the packets are marked. Apparently, if everyone is able to mark his packets freely the whole packet marking system loses its significance. There can be different rules for packet marking resulting in different service models. An extensive account on this topic can be found in Kilkki (1999).

Best effort service model based on TCP

This brief introduction to Transmission Control Protocol (TCP) is based on Kilkki (1999, p. 46 - 50). The reason to include it in this introductory book is TCP's essential role in the functioning of the Internet. Moreover, it is useful to have basic knowledge about the historical background of the Internet. TCP is also an important example of the

end-to-end principle: the principle that, whenever possible, communications protocol operations should be defined to occur at the end-points of a communications system.

It should be stressed that this end-to-end principle is not applied in traditional, connection-oriented telecommunication networks. When the end-points have been able to control the system, for instance by dialing the recipient's number, there always has been a billing system involved in the process. In case of TCP, there is not any specific billing system that encourages the use of proper TCP implementations.

The fairness of the Internet service in the early phase of development relied on the assumption that there is in essence one homogeneous user group consisting of all Internet users. In that case, everyone was allowed to use the network for any sensible purpose, and only when there was not enough capacity for all demand, would there be a need for controlling or limiting the traffic sent to the network. Thus, the engineering philosophy was based on the model of a homogeneous community that had a common interest to design a workable network rather than on a model of service providers and customers.

Even during congestion, it was supposed that all or at least most users behave agreeably. Agreeable behavior could be that users stop transferring large files if they notice any performance problems in the network. The situation was eased by the fact that transferring information through the network was a much more complex process for the users, which limited the usage of the network. It is evident that this sort of approach has serious limitations when the population contains tens of millions of users and the use of the network has become a simple task for anyone even without any knowledge about data networks and protocols.

The next step in the Internet development was to specify protocols that automatically adjust the traffic sent to the network. If everyone is using a similar protocol and does not evade the adjustment control by using other more greedy protocols, the system could partly solve the problem of different user behaviors. Within these limits, any user who has been connected to the network has been allowed to utilize any available network resources independent of the actual purpose of the application or information. The network then provides a service that is called best effort because the network tries to transmit as many packets as possible and as soon as possible but does not give any guarantees. As a result, the realization of best-effort service consists of three main parts:

1. The network transmitting packets.
2. The TCP protocol controlling the bit rate.
3. The application capable of working in changing conditions.

Jon Postel wrote RFC 793 (Request for Comment) defining Transmission Control Protocol in 1981. It is worth noticing what was said about the objective:

“This document focuses its attention primarily on military computer communication requirements, especially robustness in the presence of communication unreliability and availability in the presence of congestion, but many of these problems are found in the civilian and government sector as well.”

Because of this background, TCP provides an effective tool to recover data that is damaged, lost, duplicated, or delivered out of order. This is achieved by assigning a sequence number to all data transmitted in the network, and requiring a positive acknowledgment from the receiving TCP. If the acknowledgment is not received within a timeout interval, the data is retransmitted. As a result, if all TCP implementations function properly, TCP is able to recover from transmission errors.

In addition, TCP provides a means for the receiver to control the amount of data sent by the sender. This property is achieved by returning a “window” indicating a range of acceptable

sequence numbers. The window indicates an allowed number of octets that the sender may transmit before receiving further permission. These characteristics are specified in the original TCP document. The basic TCP scheme does not provide, however, tools for efficiently avoiding or alleviating congestion situations inside the network. In a worst-case scenario, a combination of retransmissions and a rapidly growing load on congested links may lead to a so-called congestion collapse.

The situation may start when a new file transfer begins to fill a buffer assigned to an already loaded link. When this buffer fills up, the round-trip time for all connections rises quickly. In that case, TCP connections suppose that packets are lost, and retransmits them. Finally, several copies of the same packet may exist at the same time in the network. Consequently, the throughput of the network is permanently reduced to a small fraction of the normal throughput. This problem was addressed by RFC 896 in 1984 and various proposals to solve it have been presented and implemented thereafter.

The fundamental problem of the old TCP implementations was that the sender might start a connection by sending lots of data up to the window size advertised by the receiver. Although this simple scheme may work in small networks with large capacity, it may be harmful in large networks with several routers and highly loaded, low capacity links. Slow-start is a solution to this problem. In essence, slow-start means that after connection establishment, the sender is allowed to send only one packet before getting acknowledgement from receiver. When the sender receives a positive acknowledgment from the receiving TCP, the sender is allowed to double the amount of packets to be sent until a packet is discarded and the sender notices that the maximum available capacity in the network is reached.

Congestion can be alleviated by going into a slow-start when the sender notices a congestion situation in the network. Several different schemes are used to increase packet rate after congestion. They all induce a saw-tooth pattern in which the window size and the bit rate go regularly up and down. All connections also encounter intermittent packet losses when the total load exceeds link capacity and the buffers get full. Both the saw-tooth and packet losses are intrinsic characteristics of TCP and usually are insignificant to most end user applications.

An elementary part of the congestion problem is that the network nodes have applied a first in, first out (FIFO) principle in the buffers: packets are discarded only when the buffer is totally full. Although a FIFO principle may seem to be efficient and fair in general, in certain situations it is both inefficient and unfair. A FIFO buffer yields a similar packet-loss ratio to every connection at a certain point in time. When measured during a short period when the buffer is full, the packet-loss ratio could be very high, and consequently, almost all senders get notice of congestion at the same time. If they are also reacting at the same time, the total traffic will drop dramatically. Then for a certain period, the congested link will be underutilized because every connection begins to increase its packet rate from a low value.

One possible solution to this problem is that some randomly selected packets are discarded even before the buffer becomes full. In that way, some senders are informed about the imminent congestion in the network. Because the senders are not synchronized, it is possible to keep the bottleneck link utilized most of the time. The selection of discarded packets can be very random, or some more complicated procedure can be applied. If well-behaved TCP connections must live with other connections with different behavior or logic of control, the

result could be that TCP connections are continuously in the slow-start phase, and aggressive connections without any adjusting mechanism seize most of the capacity.

Lessons for CEE

On a technical level, the TCP/IP protocol stack is still crucial for the operation of the Internet. Contrary to the principle of a traditional telephone network, no central control system is in charge of the allocation of resources in the Internet. A subtle balance between the highly dynamic demand of users and the available, limited resources in the network is largely achieved by means of TCP implementations running simultaneously in millions of computers. There are mechanisms to be used within domains operated by independent network operators to control traffic streams more strictly, most notably, Multiprotocol Label Switching (MPLS). However, MPLS is not used to control end-to-end flows in a similar manner as the connections in telephone network are controlled. On a customer level, the contract between network operator and customer determines the access rate, sometimes also on a level of application type (for instance, critical data, best effort data, streaming, and voice calls).

As a result, TCP produces continuous fluctuations in bit rate, both on a short time scale (mostly invisible to end users) and on a longer time scale (that affects the user experience). In mobile networks, a data connection may often go to a low start mode due to properties of radio channels. Moreover, the intrinsic properties of TCP generate a considerable amount of packet losses inside the network. It would be attractive to think that a guaranteed bit rate without packet losses would bring about a better experience. Here it is especially important to understand that it is not enough to compare individual connections, one with TCP and another one with guaranteed service, but it is necessary to compare the whole systems. The current Internet providing a best effort service based on TCP and a realistic implementation that is able to offer guaranteed services.

Network dimensioning

There are some basic principles in network design that a CEE should be familiar with. The logic behind designing traditional telephone networks differs from that of designing pure data networks, mainly due to the different nature of traffic, and due to business models. Particularly now when mobile networks are more and more often transforming from phone networks to multi-service networks, network designers encounter the grand challenge of integrating different design and operating principles into one network.

Note, however, that although the following discussion is presented in the context of network dimensioning, the models can be used to evaluate and cope with uncertainties in many other similar situations. For instance, the realization of a business model often requires the availability of some resources used by other players. Those resources can also form a structure similar to a communications network. When the total demand consists of individual occurrences of demands, the aggregate demand can often be described by similar models as the traffic in communications networks. The goal of optimization is also similar: enough resources shall be

reserved to meet the customer demand in a great majority of cases, while excessive resources shall be avoided to minimize the total cost of buying and maintaining resources.

The rapid development and construction of telephone networks at the end of 19th and at the beginning of 20th century made it necessary to develop systematic methods for network design and planning. One of the most essential questions was: how many lines are needed for a certain amount of subscribers? A typical situation is depicted in Figure T.3: a number of subscribers (K) are connected to an access exchange that is then connected to other parts of the telephone network via a number of access lines (N). At that time most of the telephone exchanges were manually operated, which meant that the operator also had to decide how many customer service persons are needed to serve customer requests.

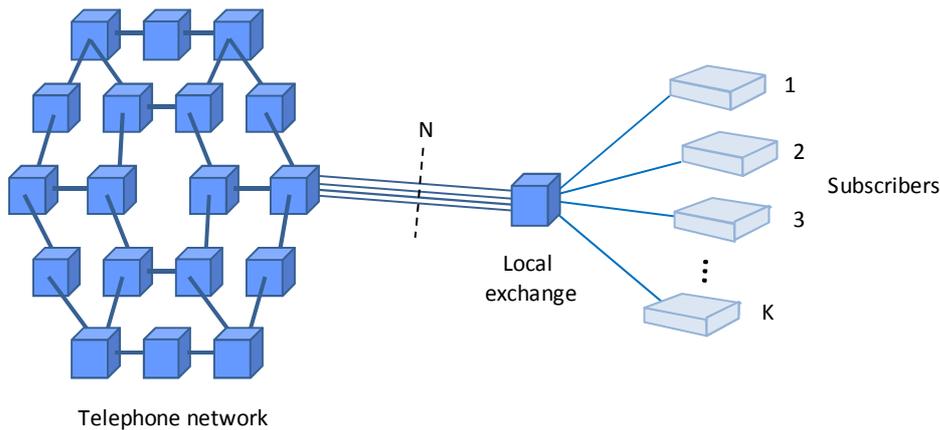


Figure T.3: Subscribers connected to the telephone network via access lines.

If there were too many complaints about the unavailability of phone service, the operator might increase the number of lines or the number of customer service persons. However, because the copper lines were expensive to build, the operators were not willing to add any unnecessary line or person. These kinds of dilemmas created need for optimization methods. Obviously, there were both systematic and random variations in the demand for telephone connections: the demand for telephone calls was higher during daytime than in the middle of the night, but there were considerable, unpredictable variations between days, too.

There are three main approaches to determine the required number of lines or persons. In the simplest approach, an operator may have a simple rule of thumb that says, for instance, that each group of 20 subscribers needs one line. The rule might be based on casual observations of how many lines were reserved in reality and how the subscribers react when a call request could not be satisfied. In the second approach, the operator may carry out systematic recording of traffic load on the links or the share of unmet requests. The operator may have a rule that if the average load exceeds 80 percent, or if more than 10 percent of call requests has to be dropped, the operator would increase the number of lines. However, continuous measurement of a large number of access links and exchanges was not an economically

efficient solution. Thus, the network operator may want to adopt the third approach in which the network dimensioning is based on a theoretical model that can be used to make a credible prediction about the future service needs of network capacity.

A. K. Erlang was able to provide a succinct but realistic formula for the dimensioning of telephone networks in a paper published 1917. An important point here is that the derivation of the formula can be presented without any connection to any specific context; thus, the model can be situated in the realm of mathematics. Still, in order to make the discussion more illuminating I use terms lines and calls although they do not belong to the realm of mathematics.

In order to derive the model we make the following assumptions:

- the probability that a new call attempt is made during a very brief period of time (Δt) is constant ($\mu \Delta t$), where μ describes the rate of new attempts,
- the average duration of calls (h) is known,
- the durations of calls are independent of each other,
- there is a fixed number of lines that serve one call per time, and
- if there is no free line, the call attempt is rejected immediately.

Note also that these assumptions mean that a rejected attempt does not affect the number of new attempts later. In other words, customers have no memory. These assumptions lead to the famous Erlang loss formula that defines the probability that a new call attempt will be rejected:

$$B(A, N) = \frac{A^N}{\sum_{i=0}^N A^i / i!}$$

where N = number of lines, i = number of occupied lines, and $A = \mu h$ = offered traffic measured in Erlangs. In the case of infinite number of lines, the number of occupied lines is Poisson-distributed:

$$P(i, A) = A^i e^{-A}.$$

Without going to the details of traffic and queuing models, some essential limitations of the Erlang model shall be recalled:

1. The formula does not tell directly what the optimal number of phone lines is in any particular case.
2. The Erlang model is able to describe traffic realistically when the traffic remains stable, whereas the model does not include systematic variations that depend on time of day or day of week.
3. The assumption that rejected attempts do not affect the process of new requests later is often unrealistic. Moreover, many other phenomena create dependencies between call attempts.
4. Offered traffic (A) cannot be known accurately in advance.

Regardless of these limitations, the planning and dimensioning of telephone networks was based on the Erlang formula for decades. How did the operators work out the limitations?

The third limitation of the Erlang model, reattempts, has been a source for numerous theoretical studies. As to the planning of telephone networks re-attempts seldom have a significant impact, because normally call blocking is kept on so low a level that reattempts do not considerably affect the traffic process. The fourth limitation of inaccurate knowledge of parameters is common to all models. Systematically made traffic measurements in real networks is the main instrument to alleviate that problem. Thus, the first two limitations are the most relevant for network operators.

As to the first limitation, most operators adopted an engineering method in which an authority defines an unambiguous requirement for the service quality measured as the blocking probability given by the Erlang formula. For instance, in Finland the dimensioning standard for links between local exchanges and the next level of network hierarchy was 1 percent for links with at least 10 lines, otherwise it was 2 percent. Thus when the average demand was estimated, based either on traffic measurements or on the number of customers, it was an easy task to look from a pre-calculated table how many lines or channels were needed. Some examples of allowed load levels are shown in Table T.4.

It is also recommendable for a CEE to create a capability to quickly assess the acceptable average traffic for a given capacity N (measured in simultaneous connections). My rule of thumb is to accept average traffic of $N - 2\sqrt{N}$, or if high availability is needed to accept traffic of $N - 3\sqrt{N}$. Note also that this model can be used in any context in which identical service requests come independently of each other and the requests are rejected if there are not enough resources to serve the new request.

Table T.4: Allowed traffic according Erlang formula.

<i>Lines (N)</i>	<i>Standard for call blocking (B)</i>				
	<i>0.1%</i>	<i>0.5%</i>	<i>1%</i>	<i>2%</i>	<i>5%</i>
<i>1</i>	0.001	0.005	0.01	0.02	0.05
<i>3</i>	0.19	0.35	0.46	0.60	0.90
<i>10</i>	3.09	3.96	4.46	5.08	6.22
<i>30</i>	16.6	19.0	20.3	21.9	24.8
<i>100</i>	75.2	80.9	84.0	87.9	95.2
<i>300</i>	258	270	277	285	302
<i>1000</i>	930	955	971	991	1 036

Furthermore, we may ask, “What is the metric that leads to this kind of standard?” or “What was the objective of the authorities when they defined the standard?” It seems that the reason for the metric was neither the profit for network operator nor the net benefit for customers. A

direct business metric would have led to a more complex algorithm that should have somehow taken into account the price paid by the customers and the cost of building the network infrastructure. The difference between the standards for larger and smaller links is a kind of attempt to include business aspects in the dimensioning standard (remember that statistical multiplexing is inefficient when the number of lines is small).

We may argue that the standard is most reasonable from the viewpoint of managing the network planning. The given standard is unambiguous and simple enough to be applied without any deeper understanding about business optimization, human behavior, or mathematics. However, it is also complex enough to create credibility and respect for the persons responsible for the design and planning of the network. Besides, the formula served its purpose well until new fancy needs and applications emerged at the end of 20th century.

We can make some general observations about the applicability of the Erlang formula. The allowed load for a link with only one line (serving several subscribers) is so low that there had to be at least two lines to satisfy any reasonable quality requirement. However, in cases with only a few subscribers the Erlang formula is somewhat unrealistic because it assumes an infinite number of customers.

Note that with call blocking of 5 percent, the allowed load exceeds the number of lines if the number of lines is above 300. In practice, a blocking probability above 5 percent likely generates disturbances that make the usage of the Erlang formula inappropriate even for lower traffic loads. The effects might either increase the traffic load (due to reattempts) or decrease the traffic load (because the service is unreliable).

Furthermore, the Erlang formula hardly is applicable if the number of lines is significantly greater than 100, because then other variations (not included in the Erlang model) start to dominate the traffic variations. Finally, if the call blocking standard is significantly below 0.1 percent, any deviation from the independence assumption of incoming calls may affect the real call blocking.

Consequently, as a rule of thumb the Erlang formula is a good starting point for planning a telephone (or other similar connection-oriented) network when the number of lines is between 2 and 200 and the call blocking requirement is between 0.1 percent and 5 percent. Note also that the Erlang formula is directly applicable only when the capacity requirement is the same for every call. Moreover, any possibility to wait for a free line changes the system behavior significantly. In numerous realistic cases, some of these criteria are not fulfilled. Here I discuss briefly about two aspects that can be included in a realistic but still quite simple model, namely, variations of average traffic demand and the variations in the capacity requirement of calls.

As to the long-term variations of traffic demand, there are two types of cases. First, in almost any communications network we can observe systematic variations that repeat themselves from day to day or from week to week. Typically, traffic demand is at the highest levels either in the afternoon or during some evening hours while the lowest levels occur during early morning hours. Similarly, we can often observe a clear pattern between days of the week. The conventional solution to cope with these variations is to define the quality requirement for the

busy hour: the busiest hour of day measured by the total traffic volume.

Figure T.4 shows an artificial measurement of telephone traffic over 5 days. The figure is created in a way that (theoretical) traffic demand is defined for each hour of the day based on typical traffic behavior in real networks. Random variations are added for each period of one hour according to a theoretical model using a normal distribution. The amount of variation depends on the average traffic and on the average duration of calls (here 3 minutes). Thus, there are no systematic variations between days, but all observed differences between days are created by random variations of the traffic process.

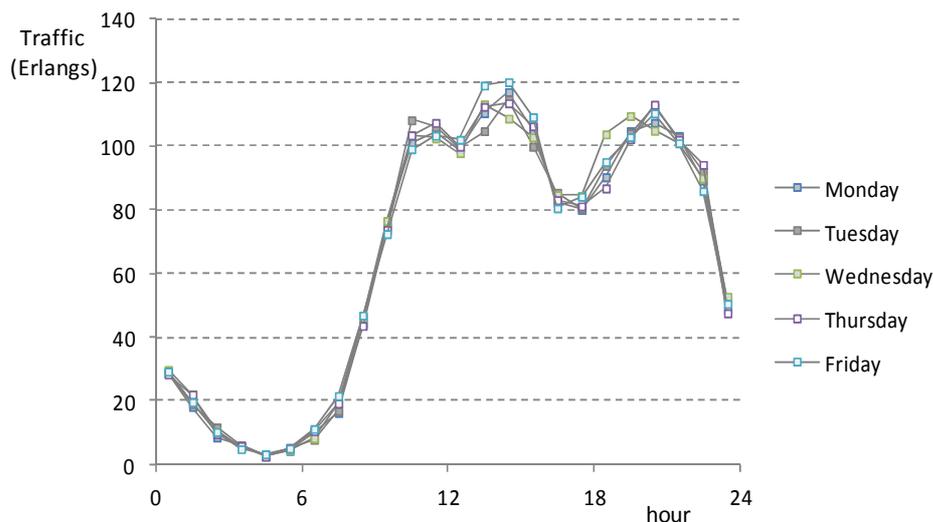


Figure T.4: Result of an artificial traffic measurement lasting five days.

How can we apply this insight in probabilistic processes to interpret the results of traffic measurements? If we conduct an average traffic demand measurement in a telephone network for one hour, we can ask: how much variation shall we expect in the measurement results? Obviously, this question can be answered by an appropriate mathematical model. However, we need some additional knowledge about the traffic process to obtain a credible answer, because it is not at all evident what the number of independent events is in this case.

Let us consider a measurement of a traffic process in which new calls are generated according to Poisson model and the length of the calls is exponentially distributed with a mean length of b . In this case, the variance of the measurement results is approximately $2Ah/T$, where A is the offered traffic in Erlangs, and T is the measurement period. Thus if the average call duration is 3 minutes, the length of the measurement is 60 minutes, and the result of the measurement is 100 Erlangs, the variance of the measurement results is 5 and standard deviation is 2.24. If we had assumed that all calls have fixed length of 3 minutes, the number of

independent measurements would be double (that is, T/h). In general, the details of the traffic process may have a significant impact on the accuracy of the measurement results. In particular, if the assumption of a Poisson process is not valid, the real accuracy might be much worse than what simple formulas indicate.

Network dimensioning in multi-service networks

This section presents a model that extends the area of applicability of the simple Erlang models. Once again, the following models can be used to evaluate and cope with uncertainties in many other similar situations although the models are presented in the context of network dimensioning.

As to the traffic in current communications networks, the most fundamental change compared to telephone networks is that there are huge variations in the capacity requirements of different applications. Even voice calls may use a different amount of network capacity depending on the coding scheme. More notably any video application requires a much higher bit rate than what is usually required by a voice call. Finally, the momentary bit rate consumed by a typical data application varies between zero and the maximum available bit rate. Thus, there are several vital questions that a network designer has to consider, particularly: “How can a network operator estimate the traffic variations, and how do those variations affect the amount of resources needed to satisfy customer needs?”

In order to make a simple analysis, let us assume that there are three categories of applications:

1. Voice sessions with fixed bit rate (R_1).
2. Video sessions with higher bit rate ($R_2 \gg R_1$).
3. Data sessions with on/off nature (bit rate is R_3 when the session is in active state and 0 when the session is inactive, and the probability of being active is q).

We use the term session here because there is not necessarily any connection admission procedure, but the sessions may start and end without any interaction between the network and user application. Typically the data bit rate (R_3) is determined by the capacity limitations of the access network, user device or server. In some cases, the limiting factor can be inside the network, for instance, in the gateways connecting networks of different Internet Service Providers (ISPs).

What is the variation of total bit rate when the average number of users in each category is known? Let us denote the average number of users by K_1 , K_2 , and K_3 , for voice, video, and data applications, respectively. If we assume that there are pure random variations (in a similar way as in the Erlang model) then the number of sessions in each category obeys a Poisson distribution. However, the number of active data sessions (sending or receiving data) is $q \cdot K_3$ not K_3 . Still, the number of active sessions is Poisson-distributed.

Now it is straightforward to calculate the mean and variance for the total bit rate if we assume that sessions start and end independently of each other:

$$E(B) = R_1K_1 + R_2K_2 + qR_3K_3,$$

$$V(B) = R_1^2K_1 + R_2^2K_2 + qR_3^2K_3.$$

If the number of sessions is large enough it can be assumed that the total traffic can be approximated by a normal distribution.

We need to remember that in reality data traffic tends to fill whatsoever capacity is available. The operator may try, however, to choose where the bottlenecks for data connections are located in the network. Customer access is typically the most reasonable location as the main bottleneck. Because the customer contract anyway defines the capacity the operator is obligated to provide, the operator has to have some control mechanisms that can be used to restrict the traffic volume. It might also be possible to limit traffic separately for several traffic classes, particularly independently for voice and data services. Why should anyone limit traffic at the consumer access point, if there is enough capacity in the core network?

There are two main reasons. The first one is business optimization. The network operator (or the network provider in the business role, or network operator as a part of Internet Service Provider) has to implement some incentives for the customers to limit the usage of network resources. Note that “limit the usage of resources” is still a technical viewpoint. From a business perspective, a better expression would be “utilize the differences in customers’ needs and willingness to pay.”

A technical expression (like “limit usage”) easily leads to an assumption that price shall be proportional to the resources used by customers. On the contrary, a business expression (like “utilize willingness to pay”) leads likely to a highly nonlinear pricing scheme at least in the consumer segment, because willingness to pay is a convex phenomenon, that is, the marginal benefit of additional bandwidth is diminished when the bandwidth grows larger. Thus, terminology defines our perspective of a problem and our perspective directs our thoughts and conclusions toward that problem.

The other reason for limiting traffic at the consumer access point is that the core network is used to blend all kinds of traffic flows coming from different sources and posing different requirements. In order to manage the core network often with huge total capacity the most practical solution is to dimension the network in a way that it fulfills all quality and quantity requirements with minimal internal differentiation. The operator may still use, primarily as a backup tool, relatively simple differentiation mechanisms that do not have any effect in normal operating conditions, but that gives desirable service differentiation in abnormal conditions.

For this purpose, the operator may estimate the traffic volume and variations for different traffic categories. The estimation shall be based both on theoretical models and real-time measurements. Models are needed to enable predictions and measurements are needed to make the model realistic. Both the measurements and the models should also include several traffic categories because the growth rates of different service demands are typically different.

Now when the network engineer has obtained the best estimations for $E(R)$ and $V(R)$ at the end of the period before the next network upgrade, he may use a simple rule for the required capacity (for more information see Kilki 1995):

$$C(R, V) = (E(R) + \gamma\sqrt{V(R)})/\rho_{max}.$$

Parameters γ and ρ_{max} define the level of safety of the dimensioning. Parameter γ takes into account the random variations in the traffic process while ρ_{max} is needed to cope with the general uncertainties of the prediction task and the systematic variations of traffic demand. Most notably, the prediction of growth rates is notoriously difficult, because external reasons may stimulate the traffic in an unpredictable way.

In case of small links with a limited amount of traffic and customers, factor γ may dominate the result whereas in case of links serving a large number of customers ρ_{max} dominates the result. Although it is somewhat risky to give any recommendations for the values of the parameters, I would use as a starting point values $\gamma = 4$ and $\rho_{max} = 0.8$. In theory (that is, assuming normal distribution), $\gamma = 4$ provides a certainty level of over 99.996 percent that there is enough capacity for all incoming traffic at any given point in time (if $\rho_{max} = 1$). But be cautious, since the optimal values for these dimensioning parameters depend on the context, environment, customer requirements, business model, and so on.

The fundamental problem in this model is the assumption that the bit rate for each data session is constant. In the Internet, that is obviously an incorrect assumption. That fact can be taken into account by introducing several data categories with different peak bit rates. Typically, the highest bit rates dominate the traffic variations. However, there is an additional challenge posed by the nature of the TCP/IP protocol because it creates dependences between data sessions (or flows). In the worst case, the result is that the bit rates of flows going through the same bottleneck are synchronized (see chapter 6 in Kilkki 1999). That means that the assumption about independent sessions is severely violated. Real traffic measurements in the Internet indicate that the assumption of a normal distribution is questionable. In reality, data traffic does not behave as smoothly as this model predicts.

Network reliability and availability

The most important performance parameter for a service is whether it is available or not when a customer wants to use the service. Let us consider an imaginary state of Lakeland illustrated in Figure T.5. The state of Lakeland is famous of its great equality. Even the 18 main cities are equal: they have the same number of inhabitants, they are surrounded by similar areas with urban, sub-urban and rural areas, and the distance to the neighboring area is the same. There is a lake in the middle of the land.

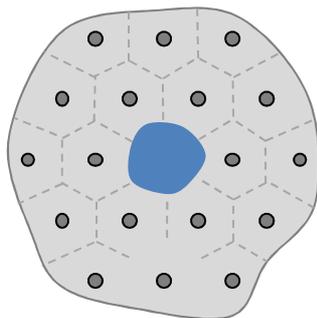


Figure T.5: The imaginary state of Lakeland with 18 cities.

Let us consider the task of designing the core network for an operator that aims to cover the area of Lakeland as well as possible. The first question to be made is: what does *well* mean in this case? In other words, what is the metric the operator should apply? The answer may vary depending on whom we are asking for:

- The technical staff first wants to define technical requirements for the network and then to build a network with a minimum amount of equipment that still fulfills the given requirements. “Well” means a network that is able to meet all the requirements posed for it.
- The staff at the business department first defines the service level offered to customers. Then they want to optimize the network in a way that the total cost including operating expenses (OPEX) and capital expenditure (CAPEX) is minimized while all the service level requirements are satisfied. “Well” means then business efficiency.
- End users want a service that is as cheap as possible but still is able to satisfy their needs well enough. “Well” means inexpensive but good enough service.
- Society wants a network that maximizes the total eudemony of all citizens.

Let us start from the technical department, because that is the main viewpoint of this chapter. The technical design process can be divided into two main parts: first, the operator needs to define what nodes are directly connected with each other, and secondly, the operator needs to determine how much capacity is required on each link.

Because the construction of long distance cables is expensive (note the business term) it seems unreasonable to directly connection cities that are not neighbors. Sea cables over the lake are not included in this analysis, although in practice they might be feasible. If all neighbors are connected with each other, there will be 36 (bi-directional) links. As another extreme, the nodes are connected in a way that they form a chain of 17 links. However, a network in which even one broken link divides the network into two separate parts is hardly ever a feasible solution. Thus, we can safely assume that at least 18 links are required. Figure T.6

shows these two extremes and four other feasible solutions with 21, 24, 27, and 30 links, respectively.

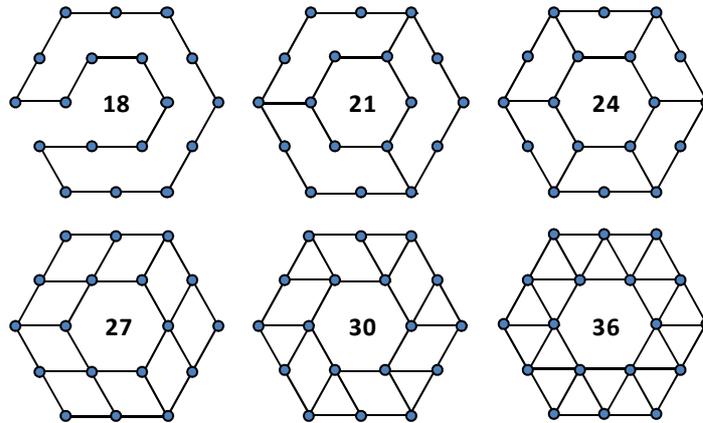


Figure T.6: Six network topologies with 18 nodes, the number in the middle indicates the number of links.

The aim of this discussion is to give some preliminary insight into network design under certain reliability constraints. In order to keep the analysis tractable, let us make the following simplifying assumptions:

1. Traffic demand is evenly distributed between every pair of nodes.
2. Traffic is routed primarily via the shortest path between the nodes.
3. The network operator defines two (or one if two is not possible) secondary paths reserved for situations in which the primary path is broken because of link failure. Each of these paths uses different links and routers than the other primary or secondary paths.
4. All failures are independent of each other.
5. Link failures dominate the network unavailability mainly due to longer repairing times. Note also that routers can be duplicated at each node to achieve better reliability.

We can make the following *technical* analysis to assess which one of the topologies is the most feasible one. First, we shall define an upper limit for the probability that no connection is available between two randomly selected pair of nodes. Then we shall estimate the required availability of an individual link to achieve this goal.

From technical viewpoint, the optimal solution is the topology with the smallest number of links among those topologies that are able to satisfy the availability requirement for a random connection. Thus, the engineer responsible for network optimization needs to know the average unavailability of a link. Let us denote it by p . For instance, we may assume that

according to the knowledge gathered by the network operator during a long a period, the mean time between failures on a link between two cities is 12 years and mean repairing time is 32 hours. The engineer can easily estimate that $p = 0.0003$.

Now the engineer may try to determine three independent paths between each pair of nodes. There are $18 \cdot 17 / 2 = 153$ pairs, but thanks to the symmetries of the networks, there are many similar paths. Let us denote the number links on the three paths by L_1 , L_2 , and L_3 respectively (L_1 is the shortest path, L_2 is the second shortest path, and L_3 is third path). For instance, the primary path from node 1 to node 15 in Figure T.7 consists of three links (1 - 13, 13 - 14, 14 - 15). The secondary path with 5 links goes through nodes 2, 3, 4, and 5. The third option is a path with six links going through nodes 12, 11, 18, 17, and 16. Note that because these paths do not have any shared links, it is always possible to find a connection between the nodes if there are at most two simultaneous link failures. However, in cases where either of the two end-nodes is 2, 4, 6, 8, 10, or 12, there are only two independent paths available, because those nodes have only two links to the remaining network.

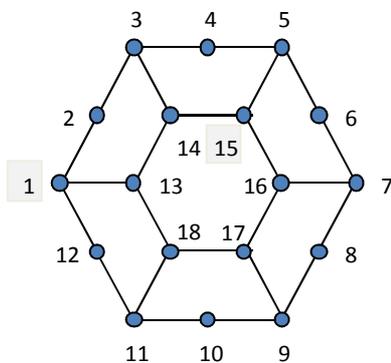


Figure T.7: A network with 18 nodes and 24 links.

Because the engineer is allowed to assume that link failures are independent of each other, the probability that a fixed path with L links is available at a random point of time is:

$$A(L) = (1 - p)^L.$$

Then if the primary path is broken, traffic is routed through the secondary path with the same or a bigger number of links. The probability that at least one of three independent paths is available is:

$$A(L_1, L_2, L_3) = 1 - (1 - A(L_1))(1 - A(L_2))(1 - A(L_3)).$$

If there are only two independent paths available, then $A(L_3) = 0$. With these assumptions the engineer is able to estimate the availability of connectivity between every pair of nodes.

Moreover, by assuming that each pair of nodes is equally important, the engineer can assess whether a certain topology satisfies the availability requirement set by the business department. As a result, he obtains the results shown in Table T.5.

Table T.5: Performance of different topologies for $p = 0.0003$.

<i>Number of links</i>	<i>Average number of hops in the primary path</i>	<i>Average number of connections going through a link</i>	<i>Probability of at least 3 independent paths (percent)</i>	<i>Average unavailability of a connection</i>
18	4.76	40.5	0	$5.12 \cdot 10^{-6}$
21	3.16	23.0	10	$1.65 \cdot 10^{-6}$
24	2.86	18.3	43	$8.61 \cdot 10^{-7}$
27	2.84	16.1	69	$4.44 \cdot 10^{-7}$
30	2.73	13.9	100	$2.03 \cdot 10^{-9}$
36	2.47	10.5	100	$1.32 \cdot 10^{-9}$

A higher availability could be achieved if alternative paths are not defined in advance, but a new path is selected optimally for each combination of broken links. However, pre-defined paths may increase the system performance owing to a faster reaction to failures. Pre-defined secondary paths might be also operationally easier.

Based on the information in Table T.5 it is straightforward to answer the question: What is the optimal topology for a given availability requirement? For instance, if the requirement is that the unavailability due to link failures shall be smaller than 10^{-6} , the engineer will answer: A network with 24 links as shown in Figure T.7.

This discussion provided an example of a technical optimization task with the following phases:

1. Someone gives specific requirements for the system to be designed, such as a required availability of connectivity between all main cities.
2. The environment for the system poses some restrictions, such as the geography of Lakeland as shown in Figure T.5 dictates that there shall be 18 core nodes.
3. Some component level data shall be gathered, such as the expected availability of individual links. The data can be based on the experience of the system owner or other sources.
4. In order to develop a tractable model, the engineer makes some simplifying assumptions, such as that all link failures occur independently of each other.
5. The engineer selects or develops a mathematical model based on the available data, and on the assumptions describing the system behavior.
6. The engineer applies the available data and the model to assess which kind of system satisfies all the requirements defined in advance.

7. The engineer applies some metric to select the best choice among the acceptable systems. Typically, it is a system with the fewest critical components.
8. Finally, the engineer gives the recommendation to select the best choice.

Even though this process seems to be straightforward, in practice there are many obstacles and caveats:

Diverse and conflicting requirements

The requirements given to the engineer are, in practice, just an interpretation of some softer objectives. For instance, there is not any noticeable difference between unavailability of $0.9 \cdot 10^{-6}$ and $1.1 \cdot 10^{-6}$, although the first one meets the requirement whereas the second one does not. In many cases, there are so many diverse requirements that it might be impossible to satisfy all of them at the same time (at least with a reasonable cost).

Changing environment

The environment is continuously changing. For instance, large data centers might be constructed near the big lake due to the opportunity of efficient cooling. That sort of development would significantly change the performance and reliability analysis in future.

Relevance of data

Although the data about the reliability and performance of network components is necessary, the data is only valid for established systems and devices. New components typically provide better performance, but they might be less reliable. The effect of human involvement is often underestimated. Particularly with novel, more capable and more complex systems, the reliability of the system may depend more on human errors than technical faults. Data about human unreliability is difficult to collect but still necessary for guaranteeing the feasibility of the models.

Dangers of simplifications

Simplifications are always necessary, because reality is always too complex to be analyzed in detail. Yet, simplifications are also dangerous. For instance, in the reliability analysis of the example the assumption of independent failures is critical. If link failures were highly correlated due to an external reason, the analysis shown above would become irrelevant. An example of an external reason that may disturb many links and routers at the same time is a large-scale power outage. Note particularly that if the communications network is disturbed that may cause further

problems in the power network. These types of interrelations make a convincing analysis very challenging.

Challenging models

If the engineer has made enough simplifications, he might even be able to make a complete analysis of the simplified model. In practice, phases 4, 5 and 6 form often an iterative structure in which simplifications are made, extensions are added, the model is modified, and preliminary analysis is conducted. This process can be learned only by doing. Sometimes several parallel models are needed to cover different aspects. Some of the models might also include simulations.

Complexity of analysis

In the above case, the analysis part was straightforward. In reality, there are several requirements and many types of resources to be optimized at the same time. Unless the engineer knows the “importance” of each requirement and the “cost” of each resource, he cannot find any unambiguous solution to the optimization problem. Note that although importance and cost can be expressed on different scales, the scales are hardly ever technical. Importance typically refers to customers’ needs while cost refers to business goals.

Selection of metrics

This issue is usually undervalued. It may seem reasonable to assume that the smaller a technical system is the more preferable it is, if the system fulfills all pre-defined requirements. In the case of network design, this typically means as few links and routers with as small capacity as possible. This is obviously a sort of hidden economic metric. However, it is also paradoxical, because from a user perspective exactly the reverse metric is reasonable: the bigger the capacity of a network the better service it can offer and the more preferable it is.

Context of recommendations

Although the engineer could be able to give a brief and clear answer to the original question (in this case: the smallest network that will be able to satisfy the pre-defined requirements consists of 24 links), in reality the decision makers often want to get a better and more elaborate insight into the issue as a whole. The engineer must be prepared for all kinds of questions, even for seemingly irrelevant questions. One of the main objectives of this book is that any prospective expert in communications ecosystem will be able to assess the case under study from many perspectives by using different tools and different terminology.

We may discern some variations of this engineering approach. First, the scope of the modeling effort might be more extensive and/or more detailed thus less simplifications have to be made. Then the analysis and the interpretation of results are much harder to make and explain to those people that finally make the business decisions (I have encountered this problem numerous times during my career). Secondly, the engineer may try to build a real life prototype that may allow the testing of the system without simplifications. However, even if the tested system is as realistic as possible, it is very seldom possible to simulate the real environment, for instance, the effect of human errors or satisfaction of real customers.

The economic optimality of different network architectures can be assessed based on the engineering analysis if the key economic parameters and their realistic values are determined. The main cost is the physical construction of the links including digging the cables and the maintenance expenses. From an economic perspective those cost factors can be summarized as the total cost of a link per year. The price per year for a leased line is a proper indication of that cost. Thus, let us assume that a leased line service is available between any pair of cities at a price of M Euros per year.

The other critical parameter is the cost of unavailability. Now if we consider the issue from the perspective of the network operator, the cost is defined in service level agreements (SLAs) with the most critical customers. The form of the agreement may be complicated and take into account the frequency and length of breaks in the service, and which particular connections are broken. Here we may assume for simplicity that there is a fixed fee of Z Euros per hour for any period when there is no connection between two core nodes. That means that if the whole network is down, the fee will be Z times the number of core node pairs ($153 Z$).

Moreover, although the network provider pays penalties to the customers, there can be additional long-term effects through customer satisfaction and brand value. Thus, the network provider may decide to invest more in the network infrastructure and build a network with 27 links as illustrated in Figure T.6. The network provider may adopt a rule that whenever a critical customer wants highly reliable service, there must be at least three links to other core nodes. A skillful CEE should be able to do this kind analysis and convince the customers that certain topology is the right choice for them.

Note that an engineer tends to estimate the effect of various parameters on the performance of the network. Even if the task given to the research engineer was to assess the feasibility of different topologies, the performance analysis provides other useful information as well. An example is shown in Figure T.8. Now the engineer may be inclined to formulate the following rules of thumb:

- If the traffic load is equal between node pairs, the average availability is dominated by those node pairs that have the fewest number of alternative paths.
- High availability can be achieved even if the availability for individual links is moderate, if there are at least three independent paths between each pair of nodes.

Figures of type T.8 and the above rules of thumb are useful when discussing with other experts in other domains. Further, it would be even more useful to convert this insight to business language.

Then as to the availability requirement for the core network, we must be aware of the fact that many (internal) services necessary for the proper operation of the core network depend on the connectivity between core nodes. Therefore, the requirement for the availability of connectivity in the core network is very high, and essentially higher than the requirement posed by end user connections.

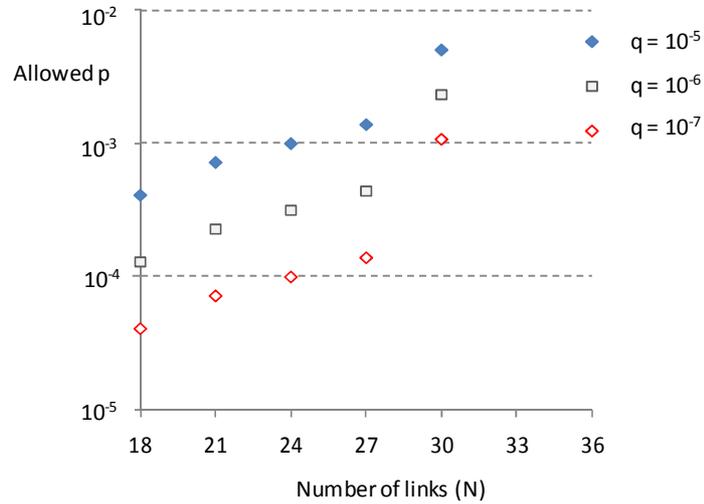


Figure T.8: Allowed availability of a core link (p) as function of service availability (q) and network topology.

As an additional mental exercise, you may try to imagine in your mind the network topology illustrated in Figure T.9. The main point of this exercise is that it gives an intuitive way to deduce that in the given topology the network looks the same from the perspective of every node, if we omit the difference in the length of links 3 - 4, 9 - 10, and 15 - 16 compared to the other links. You may also train your spatial intelligence by imagining the 3-dimensional shape of a network in which the length of each link is equal.

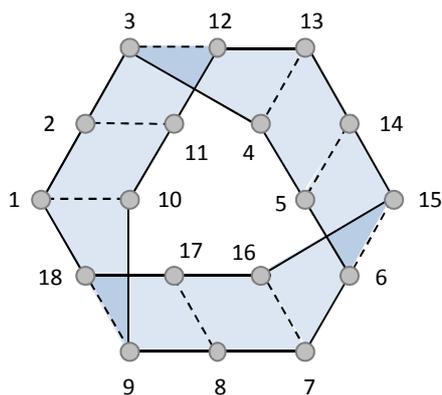


Figure T.9: A network topology that can be interpreted as a Möbius strip.

Book recommendations

W. B. Arthur, 2009, *The Nature of Technology*, New York: Free Press.

What is the true nature of something as artificial as technology? W. Brian Arthur provides an illuminating perspective on this tough question. Arthur outlines the possible processes that lead to new domains of technology. Sometimes the evolution is smooth but quite often it goes through crashes. A historical perspective allows an ecosystem expert to observe and understand the development of new technologies in a balanced way.

A.-L. Barabási, 2003, *Linked*, New York: Penguin Books.

Communications ecosystems include all kinds of networks from physical networks consisting of cables and nodes to networks of social relationships. Albert-László Barabási explains how all these networks share similar properties. These similarities, when valid, provide a special strength for those experts that work in multiple areas including technology, economics and social systems.

S. Talbott, 2004, *In the Belly of the Beast, Technology, Nature, and the Human Prospect*, New York: The Nature Institute.

This is one of those curious books that you probably have not been aware of beforehand. Somehow, Steve Talbott has had a strong impact on my thoughts about the pros and cons of technology. The first essay in Talbott's book "The Deceiving Virtues of Technology" is a splendid illustration of the relationship between human beings and technical devices. Numerous other essays are available at <http://www.netfuture.org/>.

References

- Kilki, K., 1995, Traffic characterization and connection admission control in ATM networks, Doctoral thesis, Helsinki University of Technology, available at <http://kilki.net/3>.
- Kilki, K., 1999, *Differentiated Services for the Internet*, MacMillan Technical Publishing, Indianapolis, available at <http://kilki.net/3>.
- Odlyzko, A., 2004, Pricing and Architecture of the Internet: Historical Perspectives from Telecommunications and Transportation, in Economics of Information Security, available at <http://www.dtc.umn.edu/~odlyzko/doc/networks.html>

