

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/220439243>

End-To-End Arguments in System Design

Article in *ACM Transactions on Computer Systems* · November 1984

DOI: 10.1145/357401.357402 · Source: DBLP

CITATIONS

1,696

READS

428

3 authors, including:



D.P. Reed
TidalScale, Inc.

39 PUBLICATIONS 5,629 CITATIONS

SEE PROFILE

End-To-End Arguments in System Design

J. H. SALTZER, D. P. REED, and D. D. CLARK

Massachusetts Institute of Technology Laboratory for Computer Science

This paper presents a design principle that helps guide placement of functions among the modules of a distributed computer system. The principle, called the end-to-end argument, suggests that functions placed at low levels of a system may be redundant or of little value when compared with the cost of providing them at that low level. Examples discussed in the paper include bit-error recovery, security using encryption, duplicate message suppression, recovery from system crashes, and delivery acknowledgment. Low-level mechanisms to support these functions are justified only as performance enhancements.

CR Categories and Subject Descriptors: C.0 [General] Computer System Organization—*system architectures*; C.2.2 [Computer-Communication Networks]: Network Protocols—*protocol architecture*; C.2.4 [Computer-Communication Networks]: Distributed Systems; D.4.7 [Operating Systems]: Organization and Design—*distributed systems*

General Terms: Design

Additional Key Words and Phrases: Data communication, protocol design, design principles

1. INTRODUCTION

Choosing the proper boundaries between functions is perhaps the primary activity of the computer system designer. Design principles that provide guidance in this choice of function placement are among the most important tools of a system designer. This paper discusses one class of function placement argument that has been used for many years with neither explicit recognition nor much conviction. However, the emergence of the data communication network as a computer system component has sharpened this line of function placement argument by making more apparent the situations in which and the reasons why it applies. This paper articulates the argument explicitly, so as to examine its nature and to see how general it really is. The argument appeals to application requirements and provides a rationale for moving a function upward in a layered system closer to the application that uses the function. We begin by considering the communication network version of the argument.

This is a revised version of a paper adapted from End-to-End Arguments in System Design by J. H. Saltzer, D.P. Reed, and D.D. Clark from the 2nd International Conference on Distributed Systems (Paris, France, April 8–10) 1981, pp. 509–512. © IEEE 1981

This research was supported in part by the Advanced Research Projects Agency of the U.S. Department of Defense and monitored by the Office of Naval Research under contract N00014-75-C-0661.

Authors' address: J. H. Saltzer and D. D. Clark, M.I.T. Laboratory for Computer Science, 545 Technology Square, Cambridge, MA 02139. D. P. Reed, Software Arts, Inc., 27 Mica Lane, Wellesley, MA 02181.

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1984 ACM 0734-2071/84/1100-0277 \$00.75

In a system that includes communications, one usually draws a modular boundary around the communication subsystem and defines a firm interface between it and the rest of the system. When doing so, it becomes apparent that there is a list of functions each of which might be implemented in any of several ways: by the communication subsystem, by its client, as a joint venture, or perhaps redundantly, each doing its own version. In reasoning about this choice, the requirements of the application provide the basis for the following class of arguments:

The function in question can completely and correctly be implemented only with the knowledge and help of the application standing at the endpoints of the communication system. Therefore, providing that questioned function as a feature of the communication system itself is not possible. (Sometimes an incomplete version of the function provided by the communication system may be useful as a performance enhancement.)

We call this line of reasoning against low-level function implementation the *end-to-end argument*. The following sections examine the end-to-end argument in detail, first with a case study of a typical example in which it is used—the function in question is reliable data transmission—and then by exhibiting the range of functions to which the same argument can be applied. For the case of the data communication system, this range includes encryption, duplicate message detection, message sequencing, guaranteed message delivery, detecting host crashes, and delivery receipts. In a broader context, the argument seems to apply to many other functions of a computer operating system, including its file system. Examination of this broader context will be easier, however, if we first consider the more specific data communication context.

2. CAREFUL FILE TRANSFER

2.1 End-to-End Caretaking

Consider the problem of *careful file transfer*. A file is stored by a file system in the disk storage of computer A. Computer A is linked by a data communication network with computer B, which also has a file system and a disk store. The object is to move the file from computer A's storage to computer B's storage without damage, keeping in mind that failures can occur at various points along the way. The application program in this case is the file transfer program, part of which runs at host A and part at host B. In order to discuss the possible threats to the file's integrity in this transaction, let us assume that the following specific steps are involved:

- (1) At host A the file transfer program calls upon the file system to read the file from the disk, where it resides on several tracks, and the file system passes it to the file transfer program in fixed-size blocks chosen to be disk format independent.
- (2) Also at host A, the file transfer program asks the data communication system to transmit the file using some communication protocol that involves splitting the data into packets. The packet size is typically different from the file block size and the disk track size.

- (3) The data communication network moves the packets from computer A to computer B.
- (4) At host B, a data communication program removes the packets from the data communication protocol and hands the contained data to a second part of the file transfer application that operates within host B.
- (5) At host B, the file transfer program asks the file system to write the received data on the disk of host B.

With this model of the steps involved, the following are some of the threats to the transaction that a careful designer might be concerned about:

- (1) The file, though originally written correctly onto the disk at host A, if read now may contain incorrect data, perhaps because of hardware faults in the disk storage system.
- (2) The software of the file system, the file transfer program, or the data communication system might make a mistake in buffering and copying the data of the file, either at host A or host B.
- (3) The hardware processor or its local memory might have a transient error while doing the buffering and copying, either at host A or host B.
- (4) The communication system might drop or change the bits in a packet or deliver a packet more than once.
- (5) Either of the hosts may crash part way through the transaction after performing an unknown amount (perhaps all) of the transaction.

How would a careful file transfer application then cope with this list of threats? One approach might be to reinforce each of the steps along the way using duplicate copies, time-out and retry, carefully located redundancy for error detection, crash recovery, etc. The goal would be to reduce the probability of each of the individual threats to an acceptably small value. Unfortunately, systematic countering of threat (2) requires writing correct programs, which is quite difficult. Also, not all the programs that must be correct are written by the file transfer-application programmer. If we assume further that all these threats are relatively low in probability—low enough for the system to allow useful work to be accomplished—brute force countermeasures, such as doing everything three times, appear uneconomical.

The alternate approach might be called *end-to-end check and retry*. Suppose that as an aid to coping with threat (1), stored with each file is a checksum that has sufficient redundancy to reduce the chance of an undetected error in the file to an acceptably negligible value. The application program follows the simple steps above in transferring the file from A to B. Then, as a final additional step, the part of the file transfer application residing in host B reads the transferred file copy back from its disk storage system into its own memory, recalculates the checksum, and sends this value back to host A, where it is compared with the checksum of the original. Only if the two checksums agree does the file transfer application declare the transaction committed. If the comparison fails, something has gone wrong, and a retry from the beginning might be attempted.

If failures are fairly rare, this technique will normally work on the first try; occasionally a second or even third try might be required. One would probably consider two or more failures on the same file transfer attempt as indicating that some part of this system is in need of repair.

Now let us consider the usefulness of a common proposal, namely, that the communication system provide, internally, a guarantee of reliable data transmission. It might accomplish this guarantee by providing selective redundancy in the form of packet checksums, sequence number checking, and internal retry mechanisms, for example. With sufficient care, the probability of undetected bit errors can be reduced to any desirable level. The question is whether or not this attempt to be helpful on the part of the communication system is useful to the careful file transfer application.

The answer is that threat (4) may have been eliminated, but the careful file transfer application must still counter the remaining threats; so it should still provide its own retries based on an end-to-end checksum of the file. If it does, the extra effort expended in the communication system to provide a guarantee of reliable data transmission is only reducing the frequency of retries by the file transfer application; it has no effect on inevitability or correctness of the outcome, since correct file transmission is ensured by the end-to-end checksum and retry whether or not the data transmission system is especially reliable.

Thus, the argument: In order to achieve careful file transfer, the application program that performs the transfer must supply a file-transfer-specific, end-to-end reliability guarantee—in this case, a checksum to detect failures and a retry-commit plan. For the data communication system to go out of its way to be extraordinarily reliable does not reduce the burden on the application program to ensure reliability.

2.2 A Too-Real Example

An interesting example of the pitfalls that one can encounter turned up recently at the Massachusetts Institute of Technology. One network system involving several local networks connected by gateways used a packet checksum on each hop from one gateway to the next, on the assumption that the primary threat to correct communication was corruption of bits during transmission. Application programmers, aware of this checksum, assumed that the network was providing reliable transmission, without realizing that the transmitted data were unprotected while stored in each gateway. One gateway computer developed a transient error: while copying data from an input to an output buffer a byte pair was interchanged, with a frequency of about one such interchange in every million bytes passed. Over a period of time many of the source files of an operating system were repeatedly transferred through the defective gateway. Some of these source files were corrupted by byte exchanges, and their owners were forced to the ultimate end-to-end error check: manual comparison with and correction from old listings.

2.3 Performance Aspects

However, it would be too simplistic to conclude that the lower levels should play no part in obtaining reliability. Consider a network that is somewhat unreliable, dropping one message of each hundred messages sent. The simple strategy outlined above, transmitting the file and then checking to see that the file has arrived correctly, would perform more poorly as the length of the file increased. The probability that all packets of a file arrive correctly decreases exponentially

with the file length, and thus the expected time to transmit the file grows exponentially with file length. Clearly, some effort at the lower levels to improve network reliability can have a significant effect on application performance. But the key idea here is that the lower levels need not provide "perfect" reliability.

Thus the amount of effort to put into reliability measures within the data communication system is seen to be an engineering trade-off based on performance, rather than a requirement for correctness. Note that performance has several aspects here. If the communication system is too unreliable, the file transfer application performance will suffer because of frequent retries following failures of its end-to-end checksum. If the communication system is beefed up with internal reliability measures, those measures also have a performance cost, in the form of bandwidth lost to redundant data and added delay from waiting for internal consistency checks to complete before delivering the data. There is little reason to push in this direction very far, when it is considered that *the end-to-end check of the file transfer application must still be implemented no matter how reliable the communication system becomes*. The proper trade-off requires careful thought. For example, one might start by designing the communication system to provide only the reliability that comes with little cost and engineering effort, and then evaluate the residual error level to ensure that it is consistent with an acceptable retry frequency at the file transfer level. It is probably not important to strive for a negligible error rate at any point below the application level.

Using performance to justify placing functions in a low-level subsystem must be done carefully. Sometimes, by examining the problem thoroughly, the same or better performance enhancement can be achieved at the high level. Performing a function at a low level may be more efficient, if the function can be performed with a minimum perturbation of the machinery already included in the low-level subsystem. But the opposite situation can occur—that is, performing the function at the lower level may cost more—for two reasons. First, since the lower level subsystem is common to many applications, those applications that do not need the function will pay for it anyway. Second, the low-level subsystem may not have as much information as the higher levels, so it cannot do the job as efficiently.

Frequently, the performance trade-off is quite complex. Consider again the careful file transfer on an unreliable network. The usual technique for increasing packet reliability is some sort of per-packet error check with a retry protocol. This mechanism can be implemented either in the communication subsystem or in the careful file transfer application. For example, the receiver in the careful file transfer can periodically compute the checksum of the portion of the file thus far received and transmit this back to the sender. The sender can then restart by retransmitting any portion that has arrived in error.

The end-to-end argument does not tell us where to put the early checks, since either layer can do this performance-enhancement job. Placing the early retry protocol in the file transfer application simplifies the communication system but may increase overall cost, since the communication system is shared by other applications and each application must now provide its own reliability enhancement. Placing the early retry protocol in the communication system may be more

efficient, since it may be performed inside the network on a hop-by-hop basis, reducing the delay involved in correcting a failure. At the same time there may be some application that finds the cost of the enhancement is not worth the result, but it now has no choice in the matter.¹ A great deal of information about system implementation is needed to make this choice intelligently.

3. OTHER EXAMPLES OF THE END-TO-END ARGUMENT

3.1 Delivery Guarantees

The basic argument that a lower level subsystem that supports a distributed application may be wasting its effort in providing a function that must, by nature, be implemented at the application level anyway can be applied to a variety of functions in addition to reliable data transmission. Perhaps the oldest and most widely known form of the argument concerns acknowledgment of delivery. A data communication network can easily return an acknowledgment to the sender for every message delivered to a recipient. The ARPANET, for example, returns a packet known as *Request For Next Message* (RFNM) [1] whenever it delivers a message. Although this acknowledgment may be useful within the network as a form of congestion control (originally the ARPANET refused to accept another message to the same target until the previous RFNM had returned), it was never found to be very helpful for applications using the ARPANET. The reason is that knowing for sure that the message was delivered to the target host is not very important. What the application wants to know is whether or not the target host acted on the message; all manner of disaster might have struck after message delivery but before completion of the action requested by the message. The acknowledgment that is really desired is an end-to-end one, which can be originated only by the target application—"I did it," or "I didn't."

Another strategy for obtaining immediate acknowledgments is to make the target host sophisticated enough that when it accepts delivery of a message it also accepts responsibility for guaranteeing that the message is acted upon by the target application. This approach can eliminate the need for an end-to-end acknowledgment in some, but not all, applications. An end-to-end acknowledgment is still required for applications in which the action requested of the target host should be done only if similar actions requested of other hosts are successful. This kind of application requires a two-phase commit protocol [5, 10, 15], which is a sophisticated end-to-end acknowledgment. Also, if the target application either fails or refuses to do the requested action, and thus a negative acknowledgment is a possible outcome, an end-to-end acknowledgment may still be a requirement.

3.2 Secure Transmission of Data

Another area in which an end-to-end argument can be applied is that of data encryption. The argument here is threefold. First, if the data transmission system performs encryption and decryption, it must be trusted to securely manage the required encryption keys. Second, the data will be in the clear and thus vulnerable

¹ For example, real-time transmission of speech has tighter constraints on message delay than on bit-error rate. Most retry schemes significantly increase the variability of delay.

as they pass into the target node and are fanned out to the target application. Third, the *authenticity* of the message must still be checked by the application. If the application performs end-to-end encryption, it obtains its required authentication check and can handle key management to its satisfaction, and the data are never exposed outside the application.

Thus, to satisfy the requirements of the application, there is no need for the communication subsystem to provide for automatic encryption of all traffic. Automatic encryption of all traffic by the communication subsystem may be called for, however, to ensure something else—that a misbehaving user or application program does not deliberately transmit information that should not be exposed. The automatic encryption of all data as they are put into the network is one more firewall the system designer can use to ensure that information does not escape outside the system. Note however, that this is a different requirement from authenticating access rights of a system user to specific parts of the data. This network-level encryption can be quite unsophisticated—the same key can be used by all hosts, with frequent changes of the key. No per-user keys complicate the key management problem. The use of encryption for application-level authentication and protection is complementary. Neither mechanism can satisfy both requirements completely.

3.3 Duplicate Message Suppression

A more sophisticated argument can be applied to duplicate message suppression. A property of some communication network designs is that a message or a part of a message may be delivered twice, typically as a result of time-out-triggered failure detection and retry mechanisms operating within the network. The network can watch for and suppress any such duplicate messages, or it can simply deliver them. One might expect that an application would find it very troublesome to cope with a network that may deliver the same message twice; indeed, it is troublesome. Unfortunately, even if the network suppresses duplicates, the application itself may accidentally originate duplicate requests in its own failure/retry procedures. These application-level duplications look like different messages to the communication system, so it cannot suppress them; suppression must be accomplished by the application itself with knowledge of how to detect its own duplicates.

A common example of duplicate suppression that must be handled at a high level is when a remote system user, puzzled by lack of response, initiates a new login to a time-sharing system. Another example is that most communication applications involve a provision for coping with a system crash at one end of a multisite transaction: reestablish the transaction when the crashed system comes up again. Unfortunately, reliable detection of a system crash is problematical: the problem may just be a lost or long-delayed acknowledgment. If so, the retried request is now a duplicate, which only the application can discover. Thus, the end-to-end argument again: If the application level has to have a duplicate-suppressing mechanism anyway, that mechanism can also suppress any duplicates generated inside the communication network; therefore, the function can be omitted from that lower level. The same basic reasoning applies to completely omitted messages, as well as to duplicated ones.

3.4 Guaranteeing FIFO Message Delivery

Ensuring that messages arrive at the receiver in the same order in which they are sent is another function usually assigned to the communication subsystem. The mechanism usually used to achieve such first-in, first-out (FIFO) behavior guarantees FIFO ordering among messages sent on the same virtual circuit. Messages sent along independent virtual circuits, or through intermediate processes outside the communication subsystem, may arrive in a different order from the order sent. A distributed application in which one node can originate requests that initiate actions at several sites cannot take advantage of the FIFO ordering property to guarantee that the actions requested occur in the correct order. Instead, an independent mechanism at a higher level than the communication subsystem must control the ordering of actions.

3.5 Transaction Management

We have now applied the end-to-end argument in the construction of the SWALLOW distributed data storage system [15], where it leads to significant reduction in overhead. SWALLOW provides data storage servers called repositories that can be used remotely to store and retrieve data. Accessing data at a repository is done by sending it a message specifying the object to be accessed, the version, and type of access (read/write), plus a value to be written if the access is a write. The underlying message communication system does not suppress duplicate messages, since (a) the object identifier plus the version information suffices to detect duplicate writes, and (b) the effect of a duplicate read-request message is only to generate a duplicate response, which is easily discarded by the originator. Consequently, the low-level message communication protocol is significantly simplified.

The underlying message communication system does not provide delivery acknowledgment either. The acknowledgment that the originator of a write request needs is that the data were stored safely. This acknowledgment can be provided only by high levels of the SWALLOW system. For read requests, a delivery acknowledgment is redundant, since the response containing the value read is sufficient acknowledgment. By eliminating delivery acknowledgments, the number of messages transmitted is halved. This message reduction can have a significant effect on both host load and network load, improving performance. This same line of reasoning has also been used in development of an experimental protocol for remote access to disk records [6]. The resulting reduction in path length in lower level protocols has been important in maintaining good performance on remote disk access.

4. IDENTIFYING THE ENDS

Using the end-to-end argument sometimes requires subtlety of analysis of application requirements. For example, consider a computer communication network that carries some packet voice connections, that is, conversations between digital telephone instruments. For those connections that carry voice packets, an unusually strong version of the end-to-end argument applies: If low levels of the communication system try to accomplish bit-perfect communication, they will probably introduce uncontrolled delays in packet delivery, for example, by re-

requesting retransmission of damaged packets and holding up delivery of later packets until earlier ones have been correctly retransmitted. Such delays are disruptive to the voice application, which needs to feed data at a constant rate to the listener. It is better to accept slightly damaged packets as they are, or even to replace them with silence, a duplicate of the previous packet, or a noise burst. The natural redundancy of voice, together with the high-level error correction procedure in which one participant says “excuse me, someone dropped a glass. Would you please say that again?” will handle such dropouts, if they are relatively infrequent.

However, this strong version of the end-to-end argument is a property of the specific application—two people in real-time conversation—rather than a property, say, of speech in general. If, instead, one considers a speech message system, in which the voice packets are stored in a file for later listening by the recipient, the arguments suddenly change their nature. Short delays in delivery of packets to the storage medium are not particularly disruptive, so there is no longer any objection to low-level reliability measures that might introduce delay in order to achieve reliability. More important, it is actually helpful to this application to get as much accuracy as possible in the recorded message, since the recipient, at the time of listening to the recording, is not going to be able to ask the sender to repeat a sentence. On the other hand, with a storage system acting as the receiving end of the voice communication, an end-to-end argument does apply to packet ordering and duplicate suppression. Thus the end-to-end argument is not an absolute rule, but rather a guideline that helps in application and protocol design analysis; one must use some care to identify the endpoints to which the argument should be applied.

5. HISTORY, AND APPLICATION TO OTHER SYSTEM AREAS

The individual examples of end-to-end arguments cited in this paper are not original; they have accumulated over the years. The first example of questionable intermediate delivery acknowledgments noticed by the authors was the “wait” message of the Massachusetts Institute of Technology Compatible Time-Sharing System, which the system printed on the user’s terminal whenever the user entered a command [3]. (The message had some value in the early days of the system, when crashes and communication failures were so frequent that intermediate acknowledgments provided some needed reassurance that all was well.)

The end-to-end argument relating to encryption was first publicly discussed by Branstad in a 1973 paper [2]; presumably the military security community held classified discussions before that time. Diffie and Hellman [4] and Kent [8] developed the arguments in more depth, and Needham and Schroeder [11] devised improved protocols for the purpose.

The two-phase-commit data update protocols of Gray [5], Lampson and Sturgis [10] and Reed [13] all use a form of end-to-end argument to justify their existence; they are end-to-end protocols that do not depend for correctness on reliability, FIFO sequencing, or duplicate suppression within the communication system, since all of these problems may also be introduced by other system component failures as well. Reed makes this argument explicitly in the second chapter of his Ph.D. dissertation on decentralized atomic actions [14].

End-to-end arguments are often applied to error control and correctness in application systems. For example, a banking system usually provides high-level auditing procedures as a matter of policy and legal requirement. Those high-level auditing procedures will uncover not only high-level mistakes, such as performing a withdrawal against the wrong account, but they will also detect low-level mistakes such as coordination errors in the underlying data management system. Therefore, a costly algorithm that absolutely eliminates such coordination errors may be arguably less appropriate than a less costly algorithm that just makes such errors very rare. In airline reservation systems, an agent can be relied upon to keep trying through system crashes and delays until a reservation is either confirmed or refused. Lower level recovery procedures to guarantee that an unconfirmed request for a reservation will survive a system crash are thus not vital. In telephone exchanges, a failure that could cause a single call to be lost is considered not worth providing explicit recovery for, since the caller will probably replace the call if it matters [7]. All of these design approaches are examples of the end-to-end argument being applied to automatic recovery.

Much of the debate in the network protocol community over datagrams, virtual circuits, and connectionless protocols is a debate about end-to-end arguments. A modularity argument prizes a reliable, FIFO sequenced, duplicate-suppressed stream of data as a system component that is easy to build on, and that argument favors virtual circuits. The end-to-end argument claims that centrally provided versions of each of those functions will be incomplete for some applications, and those applications will find it easier to build their own version of the functions starting with datagrams.

A version of the end-to-end argument in a noncommunication application was developed in the 1950s by system analysts whose responsibility included reading and writing files on large numbers of magnetic tape reels. Repeated attempts to define and implement a *reliable tape subsystem* repeatedly foundered, as flaky tape drives, undependable system operators, and system crashes conspired against all narrowly focused reliability measures. Eventually, it became standard practice for every application to provide its own application-dependent checks and recovery strategy, and to assume that lower level error detection mechanisms, at best, reduced the frequency with which the higher level checks failed. As an example, the Multics file backup system [17], even though it is built on a foundation of magnetic tape subsystem format that provides very powerful error detection and correction features, provides its own error control in the form of record labels and multiple copies of every file.

The arguments that are used in support of reduced instruction set computer (RISC) architecture are similar to end-to-end arguments. The RISC argument is that the client of the architecture will get better performance by implementing exactly the instructions needed from primitive tools; any attempt by the computer designer to anticipate the client's requirements for an esoteric feature will probably miss the target slightly and the client will end up reimplementing that feature anyway. (We are indebted to M. Satyanarayanan for pointing out this example.)

Lampson, in his arguments supporting the *open operating system*, [9] uses an argument similar to the end-to-end argument as a justification. Lampson argues

against making any function a permanent fixture of lower level modules; the function may be provided by a lower level module, but it should always be replaceable by an application's special version of the function. The reasoning is that for any function that can be thought of, at least some applications will find that, of necessity, they must implement the function themselves in order to meet correctly their own requirements. This line of reasoning leads Lampson to propose an "open" system in which the entire operating system consists of replaceable routines from a library. Such an approach has only recently become feasible in the context of computers dedicated to a single application. It may be the case that the large quantity of fixed supervisor functions typical of large-scale operating systems is only an artifact of economic pressures that have demanded multiplexing of expensive hardware and therefore a protected supervisor. Most recent system "kernelization" projects have, in fact, focused at least in part on getting function out of low system levels [12, 16]. Though this function movement is inspired by a different kind of correctness argument, it has the side effect of producing an operating system that is more flexible for applications, which is exactly the main thrust of the end-to-end argument.

6. CONCLUSIONS

End-to-end arguments are a kind of "Occam's razor" when it comes to choosing the functions to be provided in a communication subsystem. Because the communication subsystem is frequently specified before applications that use the subsystem are known, the designer may be tempted to "help" the users by taking on more function than necessary. Awareness of end-to-end arguments can help to reduce such temptations.

It is fashionable these days to talk about *layered* communication protocols, but without clearly defined criteria for assigning functions to layers. Such layerings are desirable to enhance modularity. End-to-end arguments may be viewed as part of a set of rational principles for organizing such layered systems. We hope that our discussion will help to add substance to arguments about the "proper" layering.

ACKNOWLEDGMENTS

Many people have read and commented on an earlier draft of this paper, including David Cheriton, F. B. Schneider, and Liba Svobodova. The subject was also discussed at the ACM Workshop in Fundamentals of Distributed Computing, in Fallbrook, Calif., December 1980. Those comments and discussions were quite helpful in clarifying the arguments.

REFERENCES

1. BOLT BERANEK AND NEWMAN INC. Specifications for the interconnection of a host and an IMP. Tech. Rep. 1822. Bolt Beranek and Newman Inc. Cambridge, Mass. Dec. 1981.
2. BRANSTAD, D.K. Security aspects of computer networks. AIAA Paper 73-427, AIAA Computer Network Systems Conference, Huntsville, Ala. Apr. 1973.
3. CORBATO, F.J., DAGGETT, M.M., DALEY, R.C., CREASY, R.J., HELLIWIG, J.D., ORENSTEIN, R.H., AND KORN, L.K. *The Compatible Time-Sharing System, A Programmer's Guide*. Massachusetts

- Institute of Technology Press, Cambridge, Mass. 1963, p. 10.
4. DIFFIE, W., AND HELLMAN, M.E. New directions in cryptography. *IEEE Trans. Inf. Theory* IT-22, 6 (Nov. 1976), 644-654.
 5. GRAY, J.N. *Notes on database operating systems*. Operating Systems: An Advanced Course. Lecture Notes on Computer Science, vol. 60. Springer-Verlag, New York. 1978. 393-481.
 6. GREENWALD, M. Remote virtual disk protocol specifications. Tech. Memo. Massachusetts Institute of Technology Laboratory for Computer Science, Cambridge, Mass. In preparation.
 7. KEISTER, W., KETCHLEDGE, R.W., AND VAUGHAN, H.E. No. 1 ESS: System organization and objectives. *Bell Syst. Tech. J.* 53, 5, Pt 1, (Sept. 1964), 1841.
 8. KENT, S.T. Encryption-based protection protocols for interactive user-computer communication. S.M. thesis, Dept. of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, Mass., May 1976. Also available as Tech. Rep. TR-162. Massachusetts Institute of Technology Laboratory for Computer Science, May 1976.
 9. LAMPSON, B.W., AND SPROULL, R.F. An open operating system for a single-user machine. In *Proceedings of the 7th Symposium on Operating Systems Principles*, (Pacific Grove, Calif. Dec. 10-12). ACM, New York, 1979, pp. 98-105.
 10. LAMPSON, B., AND STURGIS, H. Crash recovery in a distributed data storage system. Working paper, Xerox PARC, Palo Alto, Calif. Nov. 1976 and Apr. 1979. Submitted for publication.
 11. NEEDHAM, R.M., AND SCHROEDER, M.D. Using encryption for authentication in large networks of computers. *Commun. ACM* 21, 12 (Dec. 1978), 993-999.
 12. POPEK, G.J., et al. UCLA data secure unix. In *Proceedings of the 1979 National Computer Conference*, vol. AFIPS Press, Reston, Va., pp. 355-364.
 13. REED, D.P. Implementing atomic actions on decentralized data. *ACM Trans. Comput. Syst.* 1, 1 (Feb. 1983), 3-23.
 14. REED, D.P. Naming and synchronization in a decentralized computer system. Ph.D. dissertation, Massachusetts Institute of Technology, Dept. of Electrical Engineering and Computer Science, Cambridge, Mass. September 1978. Also available as Massachusetts Institute of Technology Laboratory for Computer Science Tech. Rep. TR-205, Sept., 1978.
 15. REED, D.P., AND SVOBODOVA, L. SWALLOW. A distributed data storage system for a local network. A. West, and P. Janson, Eds. In *Local Networks for Computer Communications, Proceedings of the IFIP Working Group 6.4 International Workshop on Local Networks* (Zurich, Aug 27-29 1980), North-Holland, Amsterdam, 1981, pp. 355-373.
 16. SCHROEDER, M.D., CLARK, D.D., AND SALTZER, J.H. The multics kernel design project. In *Proceedings 6th Symposium on Operating Systems Principles*. *Oper. Syst. Rev.* 11, 5 (Nov. 1977), 43-56.
 17. STERN, J.A. Backup and recovery of on-line information in a computer utility. S.M. thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, Mass. Aug. 1973. Available as Project MAC Tech. Rep. TR-116, Massachusetts Institute of Technology, Jan. 1974.

Received February 1983; accepted June 1983