

On the Concept of Service Reliability

Erik W. Aslaksen

Principal, Sinclair Knight Merz, and Adjunct Professor,
Graduate School of Engineering, University of Technology, Sydney

Abstract. The concept of reliability is traditionally tied to the failure of physical components. However, in the top-down approach to the design of complex systems, one starts out with the user requirements, i.e. the requirements on the service, without at first having any knowledge of the physical system which is to provide the service. This paper investigates how the concept of reliability can be defined in terms of service characteristics.

1 Introduction

In order to define the scope of this paper, it is convenient to consider engineering as composed of two reasonably distinct groups of activities; one consisting of the activities related to the development of technology, the other consisting of activities concerned with the application of technology to meet the needs of society or particular sections of society. Within the latter group, an important activity is *design*; a process which converts a set of requirements into a *system* which, through its operating lifetime, provides a *service* that satisfies those requirements.

The design process starts out with two given entities; the technology base in the form of a set of construction elements, or elements in the *physical domain*, and the requirements on the service, i.e. a set of requirements in the *functional domain*. The objective of the process is to link the two in the most cost-effective manner; that is, the process involves a transition from the functional domain (the given requirements) to the physical domain (a particular combination of construction elements). A crucial issue is then at what point in the process this transition should take place and, if some of the design process is to take place in the functional domain, how the concept of *design in the functional domain* could be developed into a well-defined process, on a par with the process of design in the physical domain. The outline of such a development has been discussed in [1,2,3].

In the present paper only a small, albeit important, part of the development of the basis for design in the functional domain is addressed, arising from the need to define two basic concepts - reliability and availability - in the functional domain. Both concepts are defined in terms of failure, but how should failure of a service be defined? In the physical domain, failure of components is defined in terms of failure modes and failure mechanisms, and the failure of combinations of elements is defined in terms of reliability block diagrams, in accordance with the bottom-up approach. In the functional domain, the starting point is a set of user requirements - are there generally valid and useful definitions of failure, reliability, and availability in terms of such a set of requirements?

2 The Quality and Availability of Service

The user requirements, which define the service to be provided, will generally be complex and involve a large number of parameters. However, as a first approximation,

one can define two parameters which characterises the service delivery (i.e. the system performance) - *Availability* and *Quality of Service (QOS)*. Roughly speaking, the availability is the probability of the service being available to the users at any point in time, and the QOS is a measure of how closely the service matches the users' expectations when it is available. It will be a function of all or some of the parameters involved in the user requirements.

That it is always, 'at least in principle, possible to find such a parameter as the QOS can be seen from the following procedure: Vary the values of all the parameters used to define the user requirements, and for each combination of values, determine what value the users place on that particular level of service, i.e. what they are willing to pay for it. That value is a single parameter characterising the level of service delivery.

However, it is important to distinguish between the two concepts QOS and value. Firstly, because the value is dependent on the particular user group and is therefore not a measure of the performance of the system alone and, secondly, because the value is often a very non-linear function of system performance, and that non-linearity is easier to handle if it is shown explicitly (see also [2], Ch. 7).

The QOS shall be denoted by S , and it shall be normalised such that $0 \leq S \leq 1$, where the value 1 corresponds to the service delivered by a fully intact system under ideal operating conditions. The availability shall be denoted by A , and as a probability it is always restricted to the range $0 \leq A \leq 1$. The values that the users would like to see for these two parameters, as specified in the user requirements, shall be called the *nominal* values, and be denoted by a subscript zero, i.e. S_0 and A_0 .

3 Stochastic Behaviour

It is in the nature of a system that its performance depends on the performance of a number of diverse, but interacting elements, and the more complex the system is, the more elements are needed to describe the performance. Consequently, in analogy with the statistical treatment of complex mechanical systems, as in statistical mechanics, it is reasonable to approach the high-level description of system behaviour through the use of statistics and probability theory.

The point of departure for any such description is the view that the performance of a system, i.e. the service produced by it, is produced by a set of *ordered* functional elements. The order is expressed as links or relations between the elements; the same set of elements differently ordered will produce a different service. Consider a system consisting of a given set of functional elements, then there is a large number of possible orderings of these elements (the actual number will depend on the extent to which the elements are distinguishable and the capability of each element to interact with other elements); each distinct ordering (i.e. a particular pattern of interactions between the elements) shall, for the present purposes only, be termed a *pure state* of the system. Let the number of such pure states be denoted by n .

The stochastic nature of the system performance is reflected by saying that at any point in time it is not possible to predict what pure state the system will be in, it is only possible to predict the probability of finding the system in state i , $1 \leq i \leq n$, denoted by $P(i)$. Then it is very tempting to use the analogy with statistical mechanics and define the concept of entropy as

$$Entropy = -\sum_{i=1}^n P(i) \ln P(i) .$$

However, it is easy to demonstrate that this simple analogy is meaningless; *any* pure state has entropy zero, and if the system is left to decay without maintenance, the entropy will first increase as order breaks down, but then it will decrease to zero as the system ends up in the fully decayed state with no interactions between the elements (this is, of course, also a pure state). The issue here is one of *intentionality*; the pure states are not equal with regard to achieving the intention of the designer, and in order to describe the stochastic behaviour correctly, the pure states must be assigned a weight reflecting their individual contribution to the intended system performance. The obvious choice in the present context is the QOS introduced in the previous section.

The equilibrium dynamics of any system, as reflected in the level of service provided to the users, arises from a balance between the (most often random) forces that break down the order of the system providing the service and the maintenance effort which restores the order. Without maintenance, the QOS would be a monotonically decreasing function of time, $S(t)$, and at any point on this function, the value of $-dS/dt$ expresses the magnitude of the destructive "force" attempting to drive the system performance downwards for this particular value of S . If this "force" is balanced by an equal restoring "force" provided by maintenance, the QOS will settle down to an equilibrium value, S .

This equilibrium value is the expectation value of a stochastic variable, s , which is the QOS of pure states, i.e. the weighting factor sought above. It is now most convenient to change the terminology slightly, and to define a pure state simply as a value of s . If s can take on values in an interval of the real line, there is a continuum of states, if s can only take on a countable number of values, the pure states are discrete.

The stochastic nature of system performance can be reflected by introducing the *service density distribution*, $\varphi(s)$,

$$P(s' \leq s \leq s'+ds) = \varphi(s')ds ,$$

with the condition

$$\int_0^1 \varphi(s) ds = 1 .$$

The relationship between $\varphi(s)$ and S is given by

$$S = \int_0^1 s \varphi(s) ds .$$

4 Service Failure, Availability, and Reliability

The picture of system behaviour outlined above is similar to that of a (dilute) gas. The state of the gas is described by the three variables volume, V , pressure, P , and

temperature, T , and the thermodynamic behaviour is described by the equation $PV/T = \text{constant}$. These variables are the expectation values of stochastic variables, and at any given point in time any of these, e.g. p , may take on a value that differs from its thermodynamic variable, P . But the volumes V (or more precisely, the number of molecules) for which these fluctuations are of significance compared to the expectation value are so small compared to the spatial resolution of our senses that the fluctuations are not normally noticeable.

However, in the case of most systems the situation is different, in so far as the magnitude and duration of the fluctuations are usually quite within the users' range of perception. A user wanting to use the service at any particular point in time sees the momentary value, s , of the QOS, and even if the expectation value, S , of the QOS equals the specified value, say S_0 , the value of s may still be less than the minimal acceptable value, say s_1 . This picture leads, then, quite naturally to a definition of service failure:

*The service is **faulty** whenever the QOS is less than a minimal acceptable value, s_1 .*

In terms of this definition, the availability of the service is immediately given by the expression

$$A = \int_{s_1}^1 \varphi(s) ds . \quad (1)$$

System reliability is also tied directly to the concept of system failure. If one considers what causes the fluctuations in the service - the random failure of individual components in the physical system providing the service - and that after an initial settling-in period, the system will consist of a mixture of components in all stages of their operating life, the time-scale on which the fluctuations in the service take place will, for the vast majority of engineered systems, remain constant. Consequently, service failure will occur at random, and the failure rate will be constant, resulting in the familiar exponential distribution of the time between failures. Let the reliability be, as always, defined as the probability of not having had a failure in the time interval from $t=0$ to $t=T$ and denoted by $R(T)$, then

$$R(T) = e^{-\frac{T}{MTBF}} , \quad (2)$$

where MTBF is the Mean Time Between Failures.

After each failure, the service will be in its failed state, i.e. with $s < s_1$, for a certain time before the value of s again rises above s_1 . The times spent in the failed state will be distributed according to some distribution function (usually not exponential), and can be characterised by the mean value, the Mean Time To Restore, or MTTR. The availability, MTBF, and MTTR are related by an expression similar to that between the reliability, availability, and maintainability of a piece of equipment, i.e.

$$\text{Availability} = \frac{MTBF}{MTBF + MTTR} ,$$

However, while this expression may be seen simply as the definition of availability, expressing nothing further about the service than what is already contained in the two parameters MTBF and MTTR, the relationship between these parameters must involve the stochastic behaviour of the service in an essential way, as it is the fluctuation of s which leads to failure in the first place. In particular, while availability could be defined in terms of the service distribution function $\varphi(s)$, the MTBF cannot; it requires knowledge of the frequency distribution of the fluctuations in the stochastic variable s , whereas $\varphi(s)$ expresses only the long-term average behaviour of the fluctuations.

5 Representation of System Behaviour

The previous sections introduced the concept of a single parameter characterising the service, the Quality of Service, S , and its stochastic nature, as described through the variable s and its density distribution, $\varphi(s)$, and defined service failure and availability in terms of this density distribution. The very general behaviour of $\varphi(s)$, reflecting the balance between the random, destructive forces to which any system is exposed and the restoring effect of maintenance, was outlined in an intuitive manner; what is now required is a quantitative expression of this behaviour. Consequently, what is needed is a *representation* of that behaviour in a manner which allows relations between system parameters to be expressed conveniently.

The approach here is exactly the same as in other, more well-known cases of representations, such as the representation of the relationship between two points in space using a coordinate system or the shape of a function using a complete set of orthonormal functions. In the former case, the numbers expressing the relationship will be different depending on the choice of coordinate system, i.e. the representation will look different, but the physical reality (i.e. the distance between the points) remains the same. In the latter case, the expansion coefficients will depend on the choice of orthonormal functions, e.g. trigonometric or Bessel, but the reality of the function, such as the information contained in it, remains the same.

A framework for representing system behaviour may be obtained by considering the service, as measured by the QOS, to be produced by N identical functional elements, each contributing $1/N$ to the value of the QOS, such that if the number of elements in the operating state is denoted by n , then $s = n/N$. Each element can be in only one of two states, operating or failed, with transition rates equal to λ (failure rate) and ρ (repair rate). As a result, the mean value of n , \bar{n} , is given by

$$\bar{n} = N \frac{\rho}{\rho + \lambda},$$

and

$$S = \bar{s} = \frac{\bar{n}}{N} = \frac{\rho}{\rho + \lambda} \quad (3)$$

In this representation, the previously introduced service density distribution $\varphi(s)$ is nothing but the binomial density distribution,

$$\varphi(s) = \varphi\left(\frac{n}{N}\right) = P(n, N, s_0) ,$$

where

$$P(n, N, s_0) = \frac{N!}{n!(N-n)!} s_0^n (1-s_0)^{(N-n)}$$

is the probability of finding n of N elements in the operating state, given that the probability of finding any one element in the operating state is s_0 . The value of N is given by the expression for availability, Eq.(1), which now reads

$$A_0 = \sum_{n=n_1}^N P(n, N, s_0) , \quad (4)$$

where $n_1 < n_0 = s_0/N$, such that whenever $n < n_1$ the service is said to have failed. That is, the previously introduced parameter s_1 is now given by $s_1 = n_1/N$.

The function $A = A(s_0, s_1, N)$ is, of course, nothing but $1-B(n_1, N, s_0)$, where B is the binomial distribution function. Consequently, the function cannot be expressed as any combination of simple functions of the variables, it is most commonly given in tabular form. However, for many applications in systems design the following approximation is adequate:

$$N(s_0, s_1, A_0) = \frac{1.97(1-s_0)^{0.42}}{(s_0 - s_1)^{1.7}} (\log(1 - A_0) - 0.3)^{1.78 - 0.95s_1} \quad (5)$$

And in any case this approximation is useful as a first step in a numerical routine that calculates $N(s_0, s_1, A_0)$ exactly.

Within this representation of service behaviour, the MTBF is the mean time between transitions from $n = n_1$ to $n = n_1 - 1$. The frequency of such transitions is given by the product of the probability of $n = n_1$, which is $P(n_1, N, s_0)$, and the rate of the transition, $n_1 \lambda$, so that

$$MTBF = \frac{1}{P(n_1, N, s_0) n_1 \lambda} \quad (6)$$

Representing the stochastic behaviour of the system in terms of N binary elements is clearly an approximation, as is giving the value of a continuous variable in terms of a binary number with a finite number of digits, or a continuous function in terms of a staircase function. The number of elements in the representation, N , is determined by the nominal values of QOS and availability, s_0 and A_0 , as well as the minimum acceptable value of the QOS, s_1 ; all part of the user requirements. Therefore, once the user requirements, i.e. the requirements on the service as seen from the users, are given, the representation to be used for the service is also determined, and the support for any function of s is the set of $N+1$ discrete values, n/N , $n=0, 1, \dots, N$.

At this stage, all the introduction of this representation has resulted in is a mapping of the four original parameters, s_0 , s_1 , A_0 , and MTBF into the four parameters n_1 , N , λ , and ρ , and while this representation makes the role of the fluctuations and their underlying statistics explicit, it is not clear that it leads to anything useful. However, if one recalls the underlying cause of the fluctuations in the service - the random failures and repairs of components in the system which produces the service, it can be argued that for a given system, the value of N is a characteristic of the system, and that the two parameters λ and ρ can be thought of as reflecting the operating point of the system. That is, the usefulness of the representation is that the two parameters λ and ρ can provide a relationship between parameters of the service and parameters of the system producing the service, such as the failure and repair rates of physical components, or acquisition and maintenance costs. Such relationships are an important part of the basis for design in the functional domain, they form the link between functionality (which relates to a class of physical systems) and service (which relates to the users), and an example of the development of such a relationship is given in the next section.

6 A Measure of Maintenance Effort

From the discussion in Sec. 3, equilibrium is reached when the amount of work expended on maintenance balances the destructive effect of the random forces breaking down the order of the system, which is proportional to $-dS/dt$ in the absence of any maintenance. That is, producing the service cannot be the only interaction between the system and its environment (in equilibrium); in order to maintain the equilibrium, work has to be performed on the system by the environment. (This is, of course, just basic thermodynamics.)

However, it is not only the amount of maintenance work that is important; the speed with which each failed functional element is restored, as measured by the parameter ρ , is also an essential factor in determining the operating point of the system, and thereby the QOS delivered. Therefore, the quantity $-\rho \cdot dS/dt$ will be called the *Level of Maintenance*. (Of course, this says nothing about the *cost* of maintaining a particular QOS.)

The quantity $-dS/dt$ can be calculated within the current representation in the following manner: Let the probability of finding the system in the state where i out of N binary functional elements are in the operating state at time t be denoted by $x_i(t)$, then the system behaviour is given by the set of coupled differential equations:

$$\frac{dx_i(t)}{dt} = (i+1)\lambda x_{i+1} - i\lambda x_i, \quad i=0, \dots, N-1,$$

and

$$\frac{dx_N(t)}{dt} = N\lambda x_N,$$

with the initial condition $x_i(0) = 1$ for $i=N$ and 0 for all other values of i . Approximating time by a discrete variable, j , with a time interval of Δ , a recursion relation is obtained for the value of x_i at the end of the j -th time interval, denoted by $x(i, j)$,

$$x(i, j) = \frac{(2 - \Delta i \lambda)x(i, j - 1) + \Delta(i + 1)\lambda(x(i + 1, j - 1) + x(i + 1, j))}{2 + \Delta i \lambda}$$

and multiplying by i and summing over all values of i gives the expectation value of s for each value of j , or $S(j)$. The differential dS/dt is then obtained as $(S(j+1)-S(j))/\Delta$, and the result can be expressed in the form

$$-\frac{dS}{dt} = \Phi(S, N)\lambda \rho S,$$

so that

$$LOM = \Phi(S, N)\lambda \rho S. \tag{7}$$

The function $\Phi(S, N)$ is shown in Fig. 1.

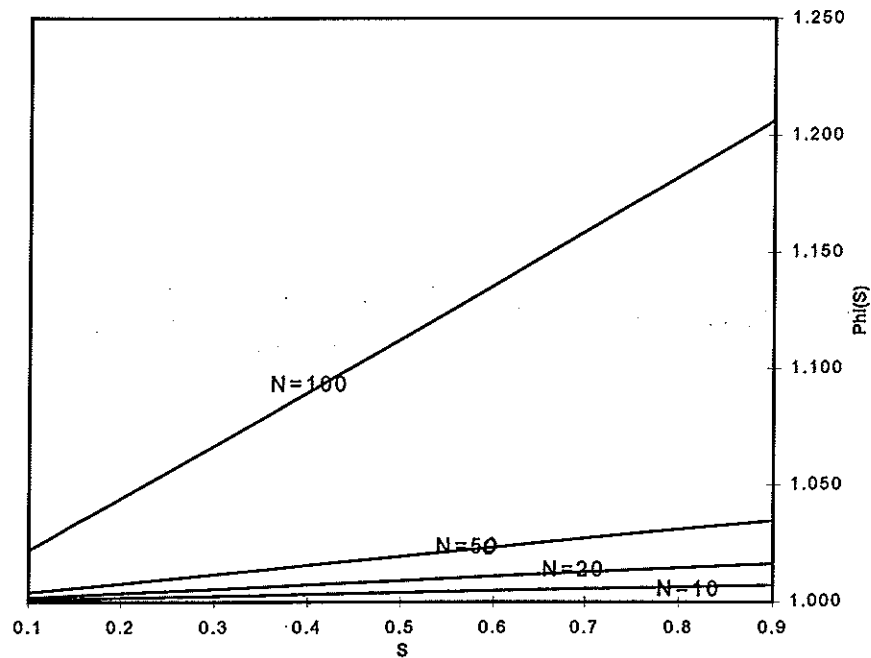


Fig.1 The function $\Phi(S, N)$, for $N=10, 20, 50, 100$.

If the decay of S were exponential, the function $\Phi(S, N)$ would be a constant. As can be seen from Fig.1, this is not the case (nor would one expect it to be, for the same reason the binomial distribution function is not a simple function); a reasonable approximation is given by

$$\Phi(S, N) \approx 1 + (1.7 \cdot 10^{-3} N + 6 \cdot 10^{-6} N^2) S .$$

Note the difference between this system parameter, LOM, and the concept of *maintainability*. The latter is a characteristic of the system design, and does not reflect the level of maintenance which is, or should be, applied to the system during its operational lifetime.

7 Conclusion

A representation of the service produced by a system has been developed in terms of a set of binary functional elements with failure and repair rates λ and ρ , respectively, with the representation being parametrised by the number of binary elements, N , and the minimum acceptable value of the Quality of Service, s_1 . The service is characterised by the three parameters Quality of Service, S , availability, A , and Mean Time Between Failure, MTBF; in addition, for the system to continue to provide this service, work must be performed on it by its environment, and the level of this work is parametrised by the Level of Maintenance, LOM. These four parameters can be expressed as functions of λ and ρ , and from the expressions for the parameters given in Eqs.(3), (4), (6), and (7), it is obvious that the curves of constant values of the entities S , A , MTBF, and LOM are straight, or almost straight lines in the $(\log(\lambda), \log(\rho))$ -plane, so that the relationships between the parameters and the two variables λ and ρ will be most conveniently displayed in this logarithmic coordinate system.

In the case of the parameter S , the function is immediately given by Eq.(3). In the case of the parameter A the function is considerably more complex, as it involves the binomial distribution function, but it is easily expressed as a small numerical routine, e.g. in Visual Basic.

In both of these cases, the lines of constant values of the parameters are straight, diagonal lines in the $(\log(\lambda), \log(\rho))$ -plane. In the case of the other two parameters, the graphic representation is not so simple but again, the functions can be expressed as small numerical routines so that, for a given representation, i.e. for given values of s_1 and N , each of the four parameters S , A , MTBF, and LOM is an easily evaluated function of λ and ρ . This now provides a basis for design in the functional domain in that, when the cost and value of having a particular QOS and availability and the cost of failure have all been determined, as well as an estimate of the cost of a given LOM, then an optimum (in terms of Return on Investment) can be found in the (λ, ρ) -plane.

By introducing the binary representation it was possible to give rigorous definitions of reliability and availability of a service without reference to any particular features of the system producing the service, thus making this approach to system design applicable to all systems.

References

- 1 Aslaksen, E.W., *A Leadership Role for INCOSE*, Proc. Fourth Annual International Symposium of the Council on Systems engineering, San Jose, CA, August 10-12, 1994, pp. 225-229.
- 2 - , *The Changing Nature of Engineering*, McGraw-Hill, 1996.
3. - , *Engineering Management - Ergonomics of the Mind*, Proc. Portland Int'l Conf. On the Management of Engineering and Technology, Portland, WA, July 1997, pp.303-306.