

Data Center Power System Reliability Beyond the 9's: A Practical Approach

Bill Brown, P.E., Square D Critical Power Competency Center

1. Abstract

Reliability has always been the focus of mission-critical power system design. While a rigorous mathematical approach to reliability is necessary to achieve the best results, the techniques used are often misunderstood or misapplied. Further, those practical elements so crucial to successful implementation are frequently overlooked. This paper gives an overview of the theoretical underpinnings of power system reliability analysis and the limiting “real-world” factors that must be used to temper any rigorous mathematical approach.

2. Introduction

Modern data center power systems represent the “ultimate” in optimization for reliability. This is a necessity, since the computer and IT equipment which these power systems serve are very sensitive to even momentary loss of power. At the highest level are those systems which qualify as “Tier IV” systems per the Uptime Institute classification system [2], which boast a representative site availability of 99.99%, or “4 nines”. But just what does “availability” really mean? Can a single number like this really give all of the necessary information about the reliability of a power system? What about individual component statistics such as MTBF – how do they come into play? Are there other features, such as maintainability, that have an effect on just how often there will be problems with the power to critical equipment? The answers to these questions can be surprising.

Fundamentally, reliability analysis for electric power systems is similar to such analysis for any other engineered system. Such analysis is predicated upon the fact that no engineered system can operate indefinitely without failure. The failure of a component within the system is treated as a random process, for which there are mathematical tools for characterization and analysis. Using these tools and making certain assumptions yields further characterizations that can be applied to the system as a whole. However, these “characterizations” (“availability” is one of them) can be misused, or, more often, over-used. And, concentrating upon a single characterization such as availability can cause the wrong conclusions to be made concerning certain practical aspects of the system, such as maintenance.

In order to get a broad perspective of just what keeps such a system up and running with reliability that is best-in-class, one must look “beyond the 9’s” at the practical aspects of reliable power system operation and how it is achieved, while at the same time knowing how the rigorous mathematical approach that quantifies “the 9’s” works and what its limitations are. With such a perspective, decisions can be made that further the cause of “keeping the power up”, rather than diminishing it.

3. “The 9’s” – The Mathematics of Reliability

The “language” of reliability is the mathematics of probability and statistics. Unfortunately, it isn’t the simplest language to understand, as one must be fluent in calculus to grasp the intricacies involved. The approach herein will be to cover the basics without extensive mathematical equations. This will give a basic acquaintance with the concepts to most readers, while the more ambitious reader can explore further the information in the Appendix and the reference materials given in section 6 for detailed answers. In particular, reference [3] contains a sound introduction to the subject.

3.1. Component-Level Reliability

Central to reliability is the concept of probability. Central to probability is the concept of an "event." Simply stated, an event is "something which happens." Because the context herein is reliability, an example event could be "component A fails" (we won't worry about the definition of "failure" for now; more on that later). If component failure is a random process, then some probability can be assigned to it, for example the probability that the component will fail in the next instant of time. Probabilities may be given in terms of a number between 0 and 1 or in percentage form – "1%" = 0.01. If the probability that component "A" fails in the next year is 1%, for example, for every year that passes there is a 1% chance that component A will fail. Another example might be the number of years between failures of component A. Such a figure is a number that behaves according to the laws of probability, and is referred to as a "random variable".

In our case, we will concentrate on the "random variable" that represents the number of years between failures of component "A." Such a random variable has certain properties. One of these is "probability density", which for a given number is the probability that the random variable is equal to that number. For the random variable we have defined here, the value of the probability density for a given number y is the probability that y is exactly the number of years until component failure. Many types of probability density exist, among them the "Gaussian" or "normal" density which is the frequently-encountered bell curve used extensively in statistics. Using the probability density one may also define, via calculus, a "probability distribution" which is the probability that a random variable is less than or equal to a given number

Because the physical phenomenon to be dealt with in reliability analysis is that of component failure, it would seem logical to assume that the longer a given component is in service, the more likely failure would become. Stated another way, the predicted amount of time between failures for a component is more likely to be less than the time in service as the time in service increases. The *exponential probability density* and *exponential probability distribution* are one way to characterize such behavior. A detailed mathematical description of the exponential probability density and distribution is given in the appendix. For our purposes here, is it sufficient to know that the exponential probability density and distribution are commonly used when describing the behavior of component failure.

Assuming that the random variable defined above (the number of years between failures of component A behaves according to the exponential probability density and distribution leads to the assumption that there is a constant *mean time between failures* and associated *failure rate* for component A. The formal definitions for MTBF and Failure Rate are [3]:

MTBF: The mean exposure time between consecutive failures of a component. It can be estimated by dividing the exposure time by the number of failures in that period, provided that a sufficient number of failures has occurred in that period.

Failure Rate: The mean number of failures of a component per unit exposure time. Usually exposure time is expressed in years and failure rate is given in failures per year.

This creates a fundamental issue, since the use of constant MTBF and Failure Rate in reliability analysis is based upon the assumption that the probability density and distribution functions for the time between failures are exponential. Further, in applying the definitions for MTBF and Failure Rate unless a sufficiently long exposure time is used the results may be overly optimistic. Both of these limitations come into play in practical situations in the data center, as will be explained herein.

Another frequently-used characteristic of a component is the *mean time to repair* or MTTR, often given the designation μ and defined as [3]:

MTTR: The mean time to repair or replace a failed component. It can be estimated by dividing the summation of repair times by the number of repairs, and, therefore, it is practically the average repair time.

The reciprocal of the MTTR is defined as the *repair rate*.

3.2. System-Level Reliability

Engineered systems can be considered to consist of multiple components, each of which has its own set of failure characteristics. System-level reliability analysis considers how these components are “connected” to form the system, and uses this information to calculate the various aspects that describe the reliability of the system. There are a number of methods available to do this. Several methods are described in detail in [3], and the method of *minimal cut-sets* is described in the appendix.

For our purposes here, it is sufficient to recognize that, whatever method we use, we have the capability to calculate the certain “reliability indices” for the system. The most common of these is *availability*, defined as [3]:

Availability: The long-term average fraction of time that a component or system is in service and satisfactorily performing its intended function. Equivalently, this is the steady-state probability that the component or system is in service.

A second index is the frequency of system failure, defined as [3]:

Frequency of System Failure: The mean number of system failures per unit time.

A third index is the expected failure duration, defined as [3]:

Expected Failure Duration: The expected or long-term average duration of a single failure event.

4. Beyond the 9's – Practical Considerations

The results of a rigorous mathematical approach to reliability are numbers for the reliability indices described in section 3. In a data center these are commonly calculated at the point of supply to the critical computer and IT equipment – thus the term “load-point reliability.” Such calculations, when viewed in the context of the assumptions made in calculating them, are very powerful tools. However, misuse or over-use of such numbers may lead to a false sense of security. Herein are given some practical aspects that must be used to temper the results of any rigorous reliability analysis.

4.1. Absolute Reliance on “Availability”

Availability, as defined in section 3, is the most often-quoted specification regarding the reliability of a data center. The much sought-after goal is “5 nines”, or 99.999% availability. But, what does this really mean? With 99.999% availability, in a given year only 0.001% of the time, or 5 minutes, 15.36 seconds. Is this one outage of 5:15.36? 5 outages of 1 minute per outage, plus another 15.36s outage? One outage of 10:30.72 every two years? In reality, availability is too vague a figure to make such specific predictions. The frequency of system failure and expected failure duration, both described in section 3, should be used for such purposes instead of the availability.

What is even more worrisome is that availability numbers calculated via a typical reliability study only apply if the component failure rates for the system hold true. In other words, in the very best case hard reliability numbers from a system reliability study are only a “snapshot” of the system availability given the conditions assumed in the study. In reality, the system reliability decreases over time as the sum aggregate effect of component aging, as illustrated in Fig. 1. The true

average availability is as given in the figure, i.e., mean uptime divided by the sum of mean uptime and mean downtime. This will usually be lower than the availability predicted by a rigorous reliability analysis.

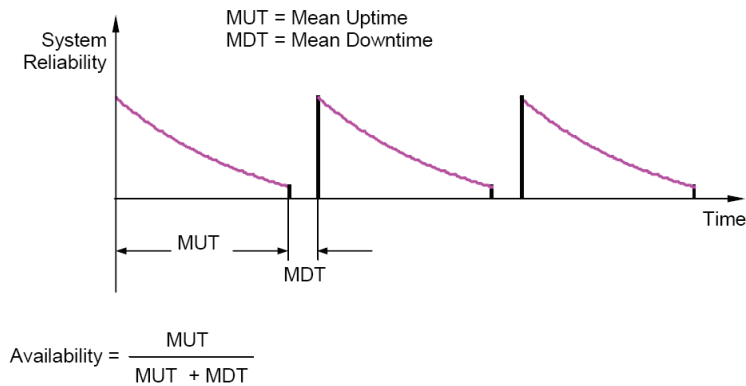


Fig. 1 Illustration of system reliability changes with time [2]

4.2. Component Failure Rates

A major assumption made in rigorous reliability analysis is the failure rate for each system component. This has a dramatic effect on the outcome of the analysis, and therefore good data is essential. In the absence of specific data from the component manufacturer survey data, such as presented in [3] and [4], can be used; indeed, the use of such data is recommended in lieu of manufacturer-specific data unless that manufacturer's component will be used or the manufacturer data contains worst-case data for all of the potential component manufacturers whose components could be used.

The first difficulty with component failure rates is the fact that, in reality, component failure rates are not constant. They change over the lifetime of the component. The nature of this change is shown graphically in the "bathtub curve" of Fig. 2. The period of heightened failure rate at the beginning of component life is known as "infant mortality". Hopefully, failures associated with infant mortality are caught during the commissioning process (more on this below). The period of heightened failure rate at the end of component life is known as "wearout". The possibility of wearout makes it very necessary to follow the manufacturer's recommended maintenance schedule in order to avoid a lowering of the reliability of the system (more on this below also). Between the infant mortality and wearout points the failure rate is approximately constant. It is this period of time, and only this period of time only, for which the results of most rigorous reliability analyses are valid.

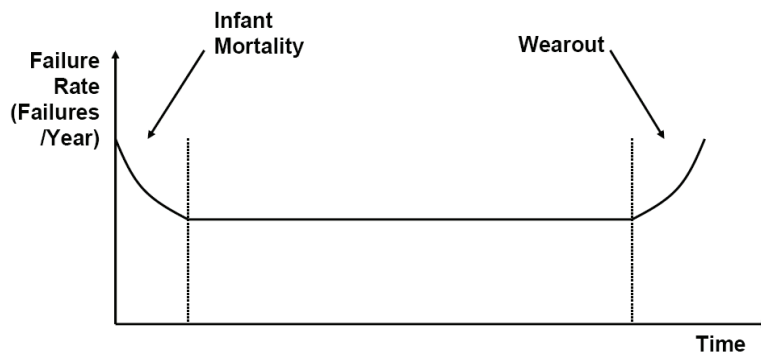


Fig. 2 Illustration of time-varying nature of component failure rate

The second issue with the use of constant failure rates is that, for most types of components, it is simply either very difficult or impossible to give such a number without survey data. Most component types have many failure modes, and these failure modes are not necessarily independent of one another. Attempts to supply such data based upon “accelerated lifetime testing”, or any method other than field history data (which takes years to build up in order to have a valid time basis for the calculation) can result in overly-optimistic results.

The third difficulty with the use of constant failure rates is their misuse if something goes wrong. Suppose “component A” in a data center power system fails, and causes an outage to the critical load. Is the component to blame for the outage? On the surface, possibly; however, shouldn't the system have been designed to take into account a single point of failure? Also, if the data center has been up and running for 1 year, and this was the only failure of component A that year, is the failure rate of component A 1 failure per year? From the end-user point of view, possibly; from a statistical point of view, no. Instead it would affect the over-all failure rate tabulated for that component by a manufacturer. In many cases, “any failure is unacceptable” is the attitude toward such component failures, however it must be pointed out that if no component ever failed (a distinct impossibility!) data center availability would always be 100%. The concern should be, instead, the historical failure rate of this component, up to and including recent data; if the failure rate is significantly higher than reference survey data such as given in [3] and [4], this could be an issue. If this is not the case, the concern should be whether or not the design takes the component failure rate in question into account.

Another aspect of component failure rates to take into account is the definition of “failure.” In establishing failure rates for components, the manufacturer and the end-user must agree on this point, otherwise component failure rates have no common basis of meaning. Per [3], the definition is:

Failure: Any trouble with a power system component that causes the following to occur:

- Partial or complete plant shutdown, or below-standard plant operation
- Unacceptable performance of user's equipment
- Operation of the electrical protective relaying or emergency operation of the plant electrical system
- De-energization of any electric circuit or equipment.

A failure on a public utility supply system may cause the user to have either of the following:

- A power interruption or loss of service
- A deviation from normal voltage or frequency outside the normal utility profile

A failure on an in-plant component causes a forced outage of the component; that is, the component is unable to perform its intended function until it is repaired or replaced. The terms “failure” and “forced outage” are often used synonymously.

Predictably, disagreements can arise regarding the precise meaning of this definition. For example, suppose “component A” acts in a manner that causes the power system to go from a higher reliability state to a lower reliability state, but does not cause an outage. Is this a “failure?” Arguments exist for both “yes” and “no” answers to this question. One potential reason is the stigma associated with the word “failure.” Component failure does not equal a defective component! As has been discussed herein, failure is a phenomenon associated with any engineered component, and is a function of usage history and maintenance as well as quality of manufacture.

4.3. Component Maintenance

A hallmark of well-designed data center power systems is *maintainability*. The necessity of maintainability is unavoidable: In order to keep component failure rates to their normal, nearly-constant values, maintenance must be performed. Otherwise, the wearout line of the bathtub curve in Fig. 2 may be crossed, with a resulting increase in failure rate.

One important point regarding system maintainability comes from the fact that *the reliability of the system is not constant*. This fact is illustrated in Fig. 1. Reliability decreases over time. Ideally, when reliability reaches the lowest acceptable level, maintenance is performed to bring the system back to an acceptable level, and the process repeats. In reality, maintenance is rarely managed to this degree of detail, but according to the component manufacturer's maintenance schedule.

The time during repair is shown as "mean downtime" in Fig. 1. For a "maintainable" system, the reliability of the system during the repair time does not go to zero as shown in the figure, but rather to some lower but non-zero level. The better the maintainability of the system, the higher the reliability level during maintenance. Such maintainability is achieved by parallel power paths and proper switching devices to allow component isolation, sufficient working space, etc.

Another aspect of component maintenance is the need for benchmarking. During maintenance, tests are performed, such as the high-potential tests, thermal scanning, etc. The results of such tests are the most meaningful when they are tracked over time, rather than simply counted as "pass/fail". An abrupt change in a test result usually signals a problem, and the best way to notice such a change is by comparing with test results from previous maintenance periods. Computerized storage of such records can facilitate this process.

Maintenance can also help pinpoint abnormal sources of component deterioration, such as overloaded circuits, improperly set protective devices, changing voltage conditions, etc.

To avoid building maintainability into a data center power system will allow reduction of construction costs. However, without maintainability outages are unavoidable – either scheduled, during maintenance, or unscheduled, due to component failure. This is true regardless of the results of a rigorous reliability study.

4.4. Commissioning

The importance of commissioning cannot be over-emphasized. True commissioning is different from simple startup of components, which only tests the component in question, usually per pre-defined standard procedures. At the highest level, commissioning tests whole systems and across systems to make sure all components work together properly. The manufacturer's startup procedures for a given type of component are simply designed to bring that component up to the point that it can be energized. True commissioning takes a system-level approach, with the goal being to ensure that the facility is functioning according to its intended purpose.

Commissioning is an engineering-based activity because it requires custom procedures to be developed based upon engineering knowledge of the system. It should test real-world conditions. Often, interoperability problems can only be found through commissioning. No system can perform to its true potential unless this important process is performed, and performed correctly. Failure to perform commissioning can lead to availability numbers that are significantly lower than those predicted by a reliability analysis.

4.5. Emergency Procedures and Trained Maintenance Staff

Emergency contingency procedures are a must to allow speedy resolution of power system issues while minimizing the impact on critical loads. Such procedures should list, step by step, actions to be taken for a given type of emergency. Unfortunately, such procedures are not the

norm for data centers, and further, even if they are their use is dependent upon trained maintenance staff.

Given the complexity of most data center power systems, it would be logical that the typical data center would have a maintenance staff that is extremely familiar with the system and large enough to cope with major emergencies. Surprisingly, this is not the case. In many instances, maintenance is contracted, with minimal on-site staff to cope with emergencies. In some cases, even the on-site staff are not familiar with system operation, instead relying upon component manufacturers to supply this familiarity. The result is that when an emergency does occur, even if emergency procedures are in place there is no one familiar enough with the system to properly implement them.

Having adequate system documentation is also a must. Data centers are dynamic environments, and as critical load components and the additional infrastructure to support it are added this documentation must be maintained. Unfortunately, in some cases up-to-date documentation, such as a single-line diagram, does not exist. If required, the services of the original engineer of record for the facility should be retained to keep system documentation up to date. Having such documentation available in key easily-accessible locations (such as posting single-line diagrams for easy reference) is also a must.

5. Summary

Data centers are dynamic systems, and therefore any type of analysis that provides a “snapshot” picture of the performance of such a system may give overly-optimistic results. Rigorous reliability analyses, when used properly to compare alternate system designs or rank the relative performance of different facilities, are powerful tools. However, other uses of the results of such analysis, such as attempting to “guarantee” the availability of a given facility over time, can lead to overly-optimistic results.

In order to maximize the reliability of any data center power system, recognition of “real world” factors such as the over-reliance on availability figures, complexities associated with component failure rates, the importance of commissioning and maintenance, and the need for adequate emergency procedures and maintenance staff is a must. Being cognizant of these real-world factors, going “beyond the 9's” and implementing accordingly, gives the best chance for success in keeping data center power systems up and running.

6. References

- [1] Tajali, R., “Maintaining the Long-Term Reliability of Critical Power Systems”, available at www.criticalpowernow.com
- [2] Turner, P.W., Seader, J.H., Brill, K.G., “Tier Classifications Define Site Infrastructure Performance, 2006, available at <http://uptimeinstitute.org>
- [3] IEEE Std.493-1997, *IEEE Recommended Practice for the Design of Reliable Industrial and Commercial Power Systems*
- [4] Hale, P.S., Jr., Arno, R.G., “Survey of Reliability and Availability Information for Power Distribution, Power Generation, and HVAC Components for Commercial, Industrial, and Utility Installations”, *Industrial and Commercial Power Systems Technical Conference, 2000, Conference Record*, 7-11 May 2000, pp. 31-54

7. Appendix – Reliability Mathematics

As stated in section 3, the “language” of reliability is the mathematics of probability and statistics. This section will give a deeper look into the concepts presented in section 3 for the reader who is interested in learning more about this subject. The approach herein will be to cover the basics without using calculus; this means that certain results will be presented without showing how they were derived. For further detail, it is recommended that the reader consult reference [3].

7.1. Component-Level Reliability

We will define the random variable X as the time between failure of component “A”, in years.

The “probability density” of a random variable is the probability that the random variable is equal to that number, written for X mathematically as $P(X=y)$ where y is the number. For the random variable X , the value of the probability density, for a given number y , is the probability that y is exactly the number of years until component failure. Using the probability density one may also define, via calculus, a “probability distribution” which is the probability that a random variable is less than or equal to a given number, written mathematically as $P(X \leq y)$.

The *exponential probability density* is defined as:

$$P(X = y) = \lambda e^{-\lambda y} \quad (1)$$

The associated *exponential probability distribution* is found using calculus and is:

$$P(X \leq y) = 1 - e^{-\lambda y} \quad (2)$$

The number λ is simply that – a number. It is interesting to note the behavior of the exponential probability density and distribution as y increases – the probability density starts at a value of λ and decreases to zero, while the probability distribution starts at a value of 0 and increases to 1. This is illustrated graphically in Fig. 3. If the exponential probability density/distribution were applied to the random variable X (which we have defined above as the number of years until failure of component A), the probability density of X for a given number y would be the probability that the number of years until the failure of component A is equal to y . Similarly, the probability distribution of X for a given number y would be the probability that the number of years until failure of component A is less than or equal to y . If the predicted amount of time between failures is more likely to be less than the time in service as the time in service increases, then the form of the exponential probability distribution per Fig. 3 b.) makes intuitive sense, since the larger the value of y (the number of years in service), the higher the probability that the time between failures is less than y . The physical meaning of the exponential probability density is harder to visualize intuitively in this case; it is sufficient to know that the two are related through calculus.

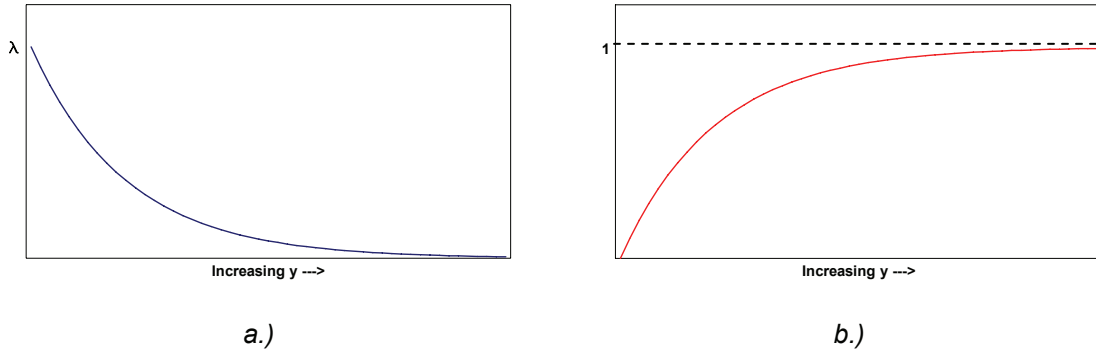


Fig. 3 Graphical representation of a.) Exponential Probability Density b.) Exponential Probability Distribution

One more thing may be determined from the exponential probability density/distribution as applied to X . Via calculus, the mean or *expected* value of X may be determined. This represents the value of X that would be determined over many observations of the failure of component A. The result: The mean value of $X = 1/\lambda$, which a constant. And since X was defined as the time between failures of component A, the number $1/\lambda$ represents the *Mean Time Between Failures*, or MTBF, for component A. The number λ itself is known as the *failure rate* of component A.

The mean time to repair (MTTR) and its reciprocal, the repair rate, were introduced in section 3. Herein, the MTBF is referred to as the variable d , the failure rate as λ , MTTR as r , and the repair rate as μ .

7.2. System-Level Reliability

As stated in section 3, there are several methods to calculate the reliability indices for a power system. Herein we introduce the method of *minimal cut-sets* as it is the simplest to understand and allows a good introduction to how such analysis is performed. This method, and several others, are described in detail in [3].

In the method of minimal cut-sets, a *cut-set* is a “set whose failure alone will cause system failure” and a *minimal cut-set* is “a cut-set with no sub-set of components whose failure alone will cause system failure.” In a power distribution system, a minimal cut-set can be considered to be a portion of the power path from the source to the load which, if removed, would cause power to the load to be lost. Inside a single minimal cut-set there may be a number of parallel paths, all of which must fail for the entire cut-set to fail. This is illustrated in Fig. 4:

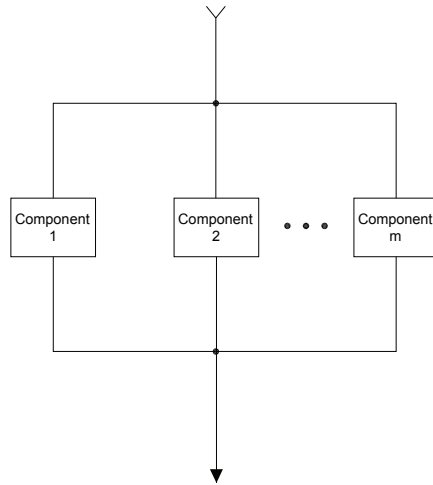


Fig. 4 Visualization of a minimal cut-set

The minimal cut-set in Fig. 4 consists of m components, all of which must fail for the cut-set to fail. The probability of failure of the cut-set is the probability that all components fail. Consider then, a single component i . Assuming that the component can be characterized with a constant MTBF and MTTR (i.e., assuming an exponential probability density for the time between failures), the component can be considered to be working for an amount of time equal to the MTBF, and not working (i.e., failed) for an amount of time equal to the MTTR. The probability that component i is in a failed state, P_i , can be calculated as:

$$P_i = \frac{\text{MTTR}_i}{\text{MTTR}_i + \text{MTBF}_i} = \frac{r_i}{r_i + d_i} = \frac{\lambda_i}{\lambda_i + \mu_i} \quad (3)$$

Now, assuming component independence the probability that the entire cut-set fails, C_i , may be calculated as product of all of the probabilities of failure of the components, i.e.,

$$C_i = \left(\frac{r_1}{r_1 + d_1} \right) \times \left(\frac{r_2}{r_2 + d_2} \right) \times \dots \times \left(\frac{r_m}{r_m + d_m} \right) = \left(\frac{\lambda_1}{\lambda_1 + \mu_1} \right) \times \left(\frac{\lambda_2}{\lambda_2 + \mu_2} \right) \times \dots \times \left(\frac{\lambda_m}{\lambda_m + \mu_m} \right) \quad (4)$$

For a given cut-set, suppose that the probability of failure C_i has been calculated. This can be thought of in terms of a “master” cut-set failure rate λ_{cs} and a “master” repair rate μ_{cs} :

$$C_i = \frac{\lambda_{cs}}{\lambda_{cs} + \mu_{cs}} \quad (5)$$

Of interest is the failure rate for the cut-set, λ_{cs} . It should be noted that the “master” repair rate for the cut-set is simply the sum of the repair rates for the components, i.e.,

$$\mu_{cs} = \mu_1 + \mu_2 + \dots + \mu_m \quad (6)$$

Equation (5) may now be used to solve for λ_{cs} :

$$\lambda_{cs} = \frac{\mu_{cs} C_i}{1 - C_i} \quad (7)$$

If the probability of failure of the cut-set is small, this can be approximated as:

$$\lambda_{cs} \approx \mu_{cs} C_i = (\mu_1 + \mu_2 + \dots + \mu_m) \left[\left(\frac{\lambda_1}{\lambda_1 + \mu_1} \right) \times \left(\frac{\lambda_2}{\lambda_2 + \mu_2} \right) \times \dots \times \left(\frac{\lambda_m}{\lambda_m + \mu_m} \right) \right] \quad (8)$$

A given minimal cut-set is described as “first order” if it has one component, “second order” if it consists of two components, etc. The entire power distribution system may be visualized as the total number of minimal cut-sets connected in series. The identification of the minimal cut-sets is via a *Failure Mode and Effects Analysis* (FMEA). In general, however, higher-order cut-sets may be neglected since their probability of occurrence is low. A visualization of a power system as series-connected minimal cut-sets is given in Fig. 5:

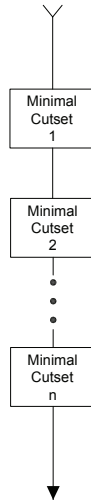


Fig. 5 Visualization of a power system as minimal cut-sets in series

In Fig. 5 the system consists of n minimal cut-sets in series, each with an associated probability of failure. The probability of failure of the system, then, is the probability that any one of the minimal cut-sets fails. This is, in general, the sum of all of the individual cut-set probabilities, minus the sum of the probabilities of any combination of simultaneous cut-set failures. If the probabilities of simultaneous cut-set failures are low, then the probability of failure of the system, S_f , can be approximated as:

$$S_f \approx C_1 + C_2 + \dots + C_n \quad (9)$$

The *frequency of failure* would provide useful insight. This is approximately the sum of the cut-set failure rates calculated per (8). This will be denoted as f_f and is given by:

$$f_f \approx \lambda_{cs1} + \lambda_{cs2} + \dots + \lambda_{csn} \quad (10)$$

The *expected failure duration* is also of interest. This will be denoted as d_f and is calculated as:

$$d_f \approx \frac{S_f}{f_f} \quad (11)$$

The reliability indices may now be calculated as:

$$\text{Availability} = 1 - S_f \quad (12)$$

$$\text{Frequency of system failure} = f_f \quad (13)$$

$$\text{Expected failure duration} = d_f \quad (14)$$

Such calculations can be automated, and this type of analysis is typically done via reliability software rather than by hand due to the number of cut-sets involved. In performing such analysis, careful consideration must be given to just what constitutes a “failure”, as this will affect the values of the component failure rates λ_i (in most cases the value used for is that for the rate of permanent failures, thus the term “permanent forced outage rate”).