

The Effect of Correlated Faults On Software Reliability

Kang Wu and Yashwant K. Malaiya
Computer Science Department
Colorado State University
Fort Collins, CO 80523

Abstract

The reliability models often assume random testing and statistical independence of faults to keep the analysis tractable. In practice, these assumptions do not hold. This paper presents a reliability modeling approach that considers nonrandom testing. This approach is used to calculate the fault exposure ratio, which characterizes the testing process. The analysis of the experimental data suggests that the fault exposure ratio varies differently in the early and the later stages of testing. The analysis here presents an explanation of this behavior.

1 Introduction

Software reliability is defined to be the probability of failure-free operation of a computer program in a specified environment for a specified time. One quantity used to measure software reliability is the fault exposure ratio [1 – 3]. The fault exposure ratio K is defined to be the ratio of the rate of change of the number of the faults and the number of the faults in a program. In the simplest case that there is one fault in a program and the testing is random, the fault exposure ratio is a constant. Therefore, the behavior of the fault exposure ratio as the function of time t tells us how far the number of the faults in a program deviates from an exponential function or the behavior of a single fault in a program. On the other hand, the value of the fault exposure ratio also gives us an idea about if the faults in a program is easy or difficult to be located.

From an extensive data analysis by Malaiya, et.al. [3], the fault exposure ratio indeed varies with the testing time. From the data analysis, Malaiya, et.al. also proposed that there are two distinguishable phases in testing process. They are the early stage of debugging and the later stage of debugging. In the early stage of debugging, the fault exposure ratio decreases against time t . In the later stage of debugging, the fault exposure ratio increases with time t . In this paper, we proposed a reliability modeling approach that considers deterministic testing. We used this approach to calculate the fault exposure ratio. The experimental data is compared with our results.

This paper is organized as follows: In Section 2, we will introduce a single fault model and an independent faults model. In Section 3, we will introduce a correlated faults model and obtain the exposure ratio for this model. We will compare our numerical results with the experimental data. We will present our conclusions in Section 4.

2 Independent Faults Model

In this section, we will introduce the independent faults model to model the random testing. We will start with the simplest case in which there is only one fault in a program. Then we will discuss the case in which there are multiple faults in a program. Finally, we will calculate the fault exposure ratio for the independent faults model.

2.1 Single Fault Model

In this section, we will deal with the case in which there is only one fault in a program. We will make some assumptions on locating the fault. Later we will use these assumptions for the case in which there are multiple faults in a program. Throughout the paper, we assume that when the failure appears, the fault is located.

Now let us start to explain our single faults model. In our single fault model, there are N identical testing teams, where $N \gg 1$. The term “identical teams” here means that they have the same efficiency and use the same testing strategy in the testing. In the other words, if two teams are identical, the probability that one team locates a fault in a program at time t is the same as the probability that another team locates the same fault in the same program at time t .

Suppose that we have a program in which there is only one fault. In general, a testing strategy is chosen before the program is tested. Test engineers use their knowledge about the program to decide which strategy is going to be used in the testing. After a testing strategy is chosen, there is no guarantee on how many inputs are used before the fault causes a failure. The quantity for evaluating a testing strategy is the average time \bar{t} that we spend before the failure appears.

To calculate the average time \bar{t} , we consider the following experiment: Each of the N identical teams has a copy of the program. They perform the testing independently. To simplify the calculation, we assume that in this experiment, the times at which the failure appears can only be t_1, t_2, \dots, t_m , where $t_1 < t_2 < \dots < t_m$. They are discretized and fixed. Notice that this assumption is just for simplifying the calculation. It will not affect the generality of the following calculation because we do not make any assumption on the value of $t_{i+1} - t_i$. Later we will show how the results for the continuous time can be obtained from the results for the discretised time. We call this experiment Experiment I. If we let the N identical testing teams start the testing, finally, we will find that there are \bar{n}_i teams locating the fault at time t_i , where

$i = 1, 2, \dots, m$. Considering the fact

$$\sum_{i=1}^m \tilde{n}_i = N, \quad (1)$$

we calculate the average time \bar{t} using

$$\bar{t} = \frac{1}{N} \sum_{i=1}^m \tilde{n}_i t_i. \quad (2)$$

Because for a testing strategy, it may be easy to locate one fault but hard to locate another fault, the average time \bar{t} depends on the following three parameters: the efficiency of the testing team, the testing strategy, and the feature of the fault. These three parameters are fixed in our single fault model. Therefore, the average time \bar{t} is a constant in the single fault model.

Because the details of the fault is unknown before we find it, we do not know which inputs can let the fault cause a failure. We call the input which lets the fault cause a failure the correct input for the fault. In the testing, before we choose a correct input for the fault, a number of inputs are chosen. The number of inputs that we use before a correct input is chosen will never be known before the fault is located. Thus the testing is random in the sense that there is no deterministic way to find the correct input for the fault. On the other hand, the probability that the correct input for the fault is chosen largely depends on the strategy we use in the testing. Therefore, the testing is not purely random. The testing is random under the ‘‘constrain’’ of the strategy used in the testing. For example, we suppose that the total number of inputs for a program is 4 and they are *input*₁, *input*₂, *input*₃, and *input*₄. For a purely random testing, all the four inputs are chosen with the same probability. Assume that the testing strategy *strategy*₁, in which *input*₄ is excluded in the testing, is applied in the testing. Then, in the testing, any of the inputs from *input*₁ to *input*₃ will be chosen with the same probability while *input*₄ will never be chosen. Therefore, for inputs from *input*₁ to *input*₃, the testing is purely random, but for all the possible inputs, the testing is not purely random. In our single fault model, Eq. (2) is the constrain due to the fact that the efficiency of the team, the strategy in the testing, and the feature of the fault are fixed.

Before we continue our discussion on the random testing, let us first introduce the notion of configurations. In Experiment I, we label the N teams using the numbers from 1 to N . We repeat Experiment I W times. To do this, we let each team has W copies of the program and use one of the W copies of the program in each experiment. Here we assume that there is not any correlation among the W experiments. This means that after each team locates the fault on one of the copies, it forgets all the information about the program and the fault so that in the testing on the other copies of the program, they will not take any advantage of the previous testings of the program.

To explain the configuration, let us first define the notations that we are going to use. We let $\tilde{n}_i[j]$ be the number of teams which locate the fault at time t_i in the j th experiment. We use $u_i^j[k]$ to denote the assigned number of the team which locates the fault at time t_i in the j th experiment, where $k = 1, 2, \dots, \tilde{n}_i[j]$. According to the way we

label the N teams, for a fixed value of j , all the values of $u_i^j[k]$'s are different and from 1 to N .

Now we are ready to explain the configuration. We define $G(j)$ to be

$$G(j) \equiv \left\{ \begin{array}{l} \{u_1^j[1], \dots, u_1^j[\tilde{n}_1[j]]\}, \\ \{u_2^j[1], \dots, u_2^j[\tilde{n}_2[j]]\}, \\ \vdots \\ \{u_m^j[1], \dots, u_m^j[\tilde{n}_m[j]]\} \end{array} \right\}.$$

The configuration for the j th experiment is defined to be $G(j)$. We need to remind the readers that in $G(j)$, the sequence of $u_i^j[k]$'s in $\{u_i^j[1], \dots, u_i^j[\tilde{n}_i[j]]\}$ is arbitrary and different sequence of it will not give us a new configuration. The meaning of the configuration $G(j)$ is the information about which team locates the fault at time t_i , where $i = 1, \dots, m$. Therefore, we can see that the configuration G is the most detail information about the experiment.

The probability that a configuration appears is related to the way in which the inputs are chosen. If all the possible inputs are chosen with the same probability, all the possible configurations will appear with the same probability. In the pure random testing, every possible input can be chosen with the same probability. Therefore, the probability of locating the fault at time t_{i1} is the same as the probability of locating it at time t_{i2} . It is apparently that in this case, every possible configurations appear with the same probability. For the nonpure random testing, the inputs are chosen randomly with the same probability but under the constrain of Eq. (2). Therefore, the possible configurations which are subject to Eq. (2) appear with the same probability. However, the probability of locating the fault at time t_{i1} is different than the probability of locating it at time t_{i2} .

Now let us summarize our single fault model. In our single fault model, we make the following assumptions:

- N identical teams independently test a program in which there is only one fault, where $N \gg 1$.
- The testing is random.
- The average time of locating the fault is a constant.

In our single fault mode, we need to know the probability of locating the fault at time t_i, p_i . Let \bar{n}_i be the average number of teams which locate the fault at time t_i and defined to be

$$\bar{n}_i \equiv \frac{1}{W} \sum_{j=1}^W \tilde{n}_i[j] \Big|_{W \rightarrow \infty}. \quad (3)$$

According to the definition of \bar{n}_i , we obtain

$$p_i = \frac{\bar{n}_i}{N}. \quad (4)$$

Thus if \bar{n}_i is calculated, p_i can be obtained using Eq. (4). Now we are going to take another approach to calculate \bar{n}_i instead of using Eq. (3).

To obtain \bar{n}_i , we first calculate the set of \tilde{n}_i 's which has the highest probability to appear, i.e. the most likelihood set of \tilde{n}_i 's. Then we take the limit $N \rightarrow \infty$. The most

likelihood set of \tilde{n}_i 's in the limit $N \rightarrow \infty$ is equal to the set of \bar{n}_i 's. We know that all the possible configuration appear with the same probability. The most likelihood set of \tilde{t}_i 's is such set that has the most number of configurations which are subject to Eqs. (1) and (2)

Now we are ready to calculate \bar{n}_i 's. Before doing this, we first calculate the number of configurations $\Gamma(\{\tilde{n}_i\})$ for a given set of \tilde{n}_i 's. Suppose that we first arrange the N teams in such a way that the number held by the left team is less than the number held by the right team. Then we exchange the position of any two of the N teams. There will be $N!$ configurations. All the configurations we are going to find are in the $N!$ configurations. Now let us find out which configurations in the $N!$ configurations should not be counted. For those teams which locate the fault at time t_i , if we exchange any two of the \tilde{n}_i teams, it will not give us a new configuration which we are interested in. There are $\tilde{n}_i!$ number of such exchanges for these teams. We should consider these $\tilde{n}_i!$ exchanges in the $N!$ configurations as one configuration. Doing the same consideration for all the other teams, we obtain

$$\Gamma(\{\tilde{n}_i\}) = \frac{N!}{\prod_{i=1}^m \tilde{n}_i!}.$$

Considering that the value of N is sufficiently large and using Sterling formula $\ln X! \approx X \ln X - X$, we obtain

$$\ln \Gamma(\{\tilde{n}_i\}) = - \sum_{i=1}^m \tilde{n}_i \ln \frac{\tilde{n}_i}{N}. \quad (5)$$

To calculate \bar{n}_i 's, we need to calculate the maximum of the function $\ln \Gamma(\{\tilde{n}_i\})$ in Eq. (5) under the constraints in Eqs. (1) and (2). Using Lagrangian multiplier method, we obtain

$$p_i \equiv \frac{\tilde{n}_i}{N} = \frac{\exp[-k(t_i - \bar{t})]}{\sum_{i=1}^m \exp[-k(t_i - \bar{t})]}. \quad (6)$$

Here k is the Lagrangian multiplier for Eq. (2) and can be determined using Eqs. (2) and (6) in principle. Because the equation for k cannot be solved in general, we just consider k as a constant at this moment.

So far, the result we obtained is for the discretized time. Now we are going to obtain the result for the continuous time from Eq. (6). We assume $t_1 = \Delta t$ and $t_{i+1} - t_i = \Delta t$, namely $t_n = n\Delta t$. Thus Eq. (6) becomes

$$p_n = \frac{\exp[-k(n\Delta t - \bar{t})]\Delta t}{\sum_{i=1}^m \exp[-k(i\Delta t - \bar{t})]\Delta t}. \quad (7)$$

We let $t \equiv n\Delta t$, and $p(t)\Delta t \equiv p_n$. As we take the limit $\Delta t \rightarrow 0$ and keep $m\Delta t$ to be a constant, q_0 , the summation in Eq. (7) becomes an integral. The meaning of q_0 is the maximum time needed to locate the fault. If we assume that the time spent to locate the fault can range from 0 to infinity, using Eq. (2) to determine k , we obtain

$$p(t) = \frac{1}{\bar{t}} \exp(-t/\bar{t}). \quad (8)$$

The meaning of $p(t)dt$ is the probability of locating the fault in the time interval from time t to time $t + dt$. As we

have seen that, to obtain the results for continuous time from the results for discretized time, we can just simply write the summation over time t_i into the integral on time t .

Now let $P(t)$ be the probability of locating the fault from time 0 to time t , then

$$P(t) = \int_0^t p(\tau) d\tau = 1 - \exp(-t/\bar{t}).$$

We can see that for a fixed time t , $P(t)$ is an increasing function of \bar{t} . Therefore we use the value of \bar{t} to measure the hardness of locating the fault.

2.2 Multiple Faults Model

Now let us consider the case in which there are multiple faults in a program. We will use all the assumptions in Section 2.1 and will not repeat those assumptions in this section. We assume that all the faults in the program are independent. Therefore, locating one fault will not give the test engineers any clue for locating other faults. This assumption also excludes the case in which one fault masks others. For example, in the testing process, the following situation happens sometimes: In a program, because of the presence of fault i , the failure caused by fault j will not appear in the execution of this program. The only way that fault j can affect the function of the program is that fault i is removed. The independent faults assumption excludes this situation. As long as this situation does not happen frequently, the independent faults assumption will be a fairly good approximation.

In our multiple faults model, we let the N teams test the program independently. To understand our multiple faults model, let us use Experiment II and Experiment III. In Experiment II, we let each of the N teams test the program. When a team locates a fault, it will not continue to test the program but waits. After each of the N teams locates a fault in the program, we know the locations of all the located faults in the program. Because all the faults are independent, they can be located by this mean if the value of N is sufficiently large. In Experiment III, we let each of the N teams independently locate all the faults in the program. After a fault is located, it is removed from the program. We see that in Experiment III, if the located faults are not removed, it will be equivalent to Experiment II and a fault may be located multiple times by a team before the team locates all the faults. Therefore, fault removal is taken into account in Experiment III and but not in Experiment II. Fault removal is the statistic correlation between faults. Statistic correlations have nothing to do with the testing. It only affects the way in which we do the statistics. Therefore, in our multiple faults mode, Experiment II is used and the statistic correlation is not considered.

Now let r denote the number of faults in the program. The average time for locating fault j is \bar{t}_j . For any pair of \bar{t}_i and \bar{t}_j , they may be equal to or may not equal to each other. Again, we first consider the case for the discretized time. We use $t_i(j)$ to denote the time needed to locate fault j and $\tilde{n}_i(j)$ to denote the number of teams which locate fault j at time $t_i(j)$. Then we have the following

constrains on $\tilde{n}_i(j)$'s:

$$\sum_{j=1}^r \sum_{i=1}^{m_j} \tilde{n}_i(j) = N \quad (9)$$

and

$$\sum_{i=1}^{m_j} \tilde{n}_i(j)[t_i(j) - \bar{t}_j] = 0, \quad (10)$$

where m_j is the number of time tickets at which fault j is located. As the same as in Section 2.1, N , \bar{t}_j 's, and $t_i(j)$'s are constant. All $\tilde{n}_i(j)$'s are the random variables which are subject to the constrains in Eqs. (9) and (10)

Using the same method as we used in the single fault model, we calculate the number of configurations, $\Gamma(\{\tilde{n}_i(j)\})$, for a set of $\tilde{n}_i(j)$'s and obtain

$$\ln \Gamma(\{\tilde{n}_i(j)\}) = \ln \left[\frac{N!}{\prod_{j=1}^r \prod_{i=1}^{m_j} \tilde{n}_i(j)!} \right]. \quad (11)$$

Now we need to find the maximum point of the function $\Gamma(\{\tilde{n}_i(j)\})$ in Eq. (11) under the $r+1$ constrains in Eqs. (9) and (10). One can take the same approach as we used for the single fault model to obtain the maximum point of Γ . One also can take a different approach to obtain the maximum point $\{\tilde{n}_i(j)\}$. It can be shown that finding the maximum point of the function Γ with the constrains in Eqs. (9) and (10) is equivalent to finding the minimum point of the following function, E_I , without any constrains when the time becomes continuous:

$$E_I \equiv \sum_{j=1}^r \sum_{i=1}^{m_j} \left\{ \bar{t}_j^{-1} \frac{n_i(j)}{N} [t_i(j) - \bar{t}_j] + \frac{n_i(j)}{N} \ln \left[\frac{n_i(j)}{N} \right] + c \frac{n_i(j)}{N} \right\}, \quad (12)$$

where $n_i(j)$ is the number of teams which locate fault j at time $t_i(j)$ and c is to be determined by Eq. (9). Notice that we use the different symbols in Eq. (11) and (12). This is because we would like to emphasize that the variables $\tilde{n}_i(j)$'s in Eq. (11) are the solution of Eq. (10) and the variables $n_i(j)$'s in Eq. (12) are not necessarily the solution of Eq. (10), but they are still the solution of Eq. (9). Calculating the derivatives of the function E_I with respect to $n_i(j)$'s respectively and letting them be zero, we obtain

$$p_i(j) \equiv \frac{\tilde{n}_i(j)}{N} = \left\{ \sum_{j=1}^r \sum_{i=1}^{m_j} \exp\{-[t_i(j) - \bar{t}_j]/\bar{t}_j\} \right\}^{-1} \times \exp\{-[t_i(j) - \bar{t}_j]/\bar{t}_j\}. \quad (13)$$

In Eq. (13), $p_i(j)$ is the probability of locating fault j at time $t_i(j)$. If we substitute Eq. (13) into Eq. (10), we will find that $\tilde{n}_i(j)$'s are the solution of Eq. (10) in the continuous time limit. This is because we choose the correct factor $\bar{t}_i(j)^{-1}$ for the first term in Eq. (12).

Eq. (13) is the result for the discretized time. As we mentioned in the discussion for the single fault model, to obtain the results for the continuous time from the results for the discretized time, we can just simply change the summation over $t_i(j)$'s into the integral on time t . Performing this substitution and assuming that the faults can be located at any time between 0 to infinity, we obtain the result for the continuous time:

$$p_j(t) = \left(\sum_{j=1}^r \bar{t}_j \right)^{-1} \exp(-t/\bar{t}_j). \quad (14)$$

where $p_j(t)dt$ is the probability of locating the fault j during the time from t to $t+dt$. If $r=1$, Eq. (14) recovers to Eq. (8).

The probability of locating fault j during the time from time 0 to time t , $P_j(t)$, is obtained using Eq. (14)

$$P_j(t) = \int_0^t p_j(\tau) d\tau = \left(\sum_{j=1}^r \bar{t}_j \right)^{-1} \bar{t}_j [1 - \exp(-t/\bar{t}_j)]. \quad (15)$$

If we set t to be ∞ , we obtain the probability of locating fault j during the time from 0 to ∞ :

$$P_j(\infty) = \frac{\bar{t}_j}{\sum_{j=1}^r \bar{t}_j}. \quad (16)$$

The value of $NP_j(\infty)$ gives us the number of teams which locate fault j . From Eq. (16), we can see that the harder fault j is to be located, the bigger the value of $P_j(\infty)$ is. This result tells us that in our theory, the value of P_j does not necessarily reflect the hardness of locating a fault. Only the average time of locating a fault gives us the hardness of locating the fault.

The ratios t/\bar{t}_j 's are very important in evaluating how much effort is used to locate the faults. Let us first calculate the average value of the ratio t/\bar{t}_j for fault j . We use $\langle A \rangle_j$ to denote the average value of a variable A for fault j . Using Eqs. (14) and (16) we obtain

$$\langle t/\bar{t}_j \rangle_j = \frac{\langle t \rangle_j}{\bar{t}_j} = \frac{\bar{t}_j}{\sum_{j=1}^r \bar{t}_j} = P_j(\infty). \quad (17)$$

Eq. (17) tells us that if we scale the quantity $\langle t \rangle_j$ by the average time \bar{t}_j , then we will obtain the fraction of the teams which take part in locating fault j .

Now let us discuss the meaning of the ratio t/\bar{t}_j in more details using the discretized time. We assume that a set of $n_i(j)$'s is assigned for locating fault j . Let T_j be the average time used to locate fault j for the set of $n_i(j)$'s. T_j is calculated using the following equation:

$$T_j = \frac{\sum_i n_i(j)t_i(j)}{N_j},$$

where N_j is the total number of the teams participating in locating fault j , namely

$$N_j = \sum_i n_i(j).$$

Let ΔT_j denote the value of $T_j - \bar{t}_j$. The value of ΔT_j tells us that on the average, how much more time a team needs to locate fault j if $\Delta T_j < 0$ and how much more time a team spends to locate fault j than it should spend if $\Delta T_j > 0$. In the other words, if $\Delta T_j < 0$, then more effort is needed to locate fault j ; if $\Delta T_j > 0$, then too much effort is used to locate fault j . In order to compare how much more or less effort is used for locating fault j with the efforts for locating other faults, we need to use the relative deviation $\Delta T_j / \bar{t}_j$ multiplied by the total number of people who locate fault j , e_j , namely,

$$e_j = N_j \bar{t}_j^{-1} \Delta T_j.$$

The first term of function E_I in Eq. (12) is e_j/N . Now we understand that the first term of E_I is actually the time related effort for locating fault j . Later, we will discuss more about the effort on the testing.

2.3 The Fault Exposure Ratio For The Independent Fault Model

In this section, we will use the results from the previous sections to calculate the fault exposure ratio for the independent faults model. The fault exposure ratio, K , is defined to be

$$K = -M(t)^{-1} \frac{dM(t)}{dt} T_L,$$

where $M(t)$ is the number of the faults left in the program at time t and T_L is the average time needed for a single execution. In Experiment II, after a fault is located, the fault can be considered to be removed from the program. Therefore, we have $M(t) = M_0 R(t)$, where M_0 is the total number of faults at $t = 0$ and $R(t)$ is the probability that a fault is in the program after time t , i.e. $R(t) = 1 - P(t)$. Thus we obtain

$$K = -R(t)^{-1} \frac{dR(t)}{dt} T_L. \quad (18)$$

Eq. (18) gives us the relation between the fault exposure ratio K and the quantity $R(t)$ which can be calculated analytically in our independent faults model.

Now let us turn to the multiple faults model. $P_j(t)$ in Eq. (15) is the probability of locating fault j during time t . The probability of locating any fault is the sum of all the $P_j(t)$'s. Doing this using Eq. (18), we obtain

$$K(t) = \frac{\sum_{j=1}^r \exp(-t/\bar{t}_j)}{\sum_{j=1}^r \bar{t}_j \exp(-t/\bar{t}_j)} T_L. \quad (19)$$

Eq. (19) is different than the result obtained by Y.K.Malaiya, et.al. [3]. This is because they are two different models. For large value of t , the terms that contain

the largest value of \bar{t}_j 's in the summations, \bar{t}_{max} , dominates the values of the summations. Thus, for large value of t , we have $K(t) \approx \bar{t}_{max}^{-1} T_L$. We can see that for large value of t , the fault exposure ratio for the independent faults model recovers to the fault exposure ratio given by Y.K.Malaiya, et.al. [3].

It can be shown that the function $K(t)$ in Eq. (19) is always a decreasing function of time t for any set of \bar{t}_j 's. This tells us that the fault exposure ratio calculated from the independent faults model is a decreasing function of time t . As the time passes, the value of the fault exposure ratio becomes smaller and smaller and the smallest value of the fault exposure ratio is the value of $\bar{t}_{max}^{-1} T_L$. This behavior does not agree with what we observed in reality. This deviation results from the assumption of our independent faults model.

3 Correlated Faults Model

In this section, we will take the nonrandom testing into account in the calculation. We will see that the nonrandom testing can be treated as the dynamic correlations between faults. The correlated faults model will be introduced to model the dynamic correlations. As in the independent faults model, because statistic correlations are not related to the testing, it is not considered in the correlated faults model.

Now let us first explain the dynamic correlations between faults. In programs, we can divide the faults into three types according to the way in which the testing is performed and they are found. The faults of the first type are the independent faults which we have considered in the previous section. The independent faults are located due to random testing. Most of this type of faults are found in the early period of testing. In this period, the testing engineers randomly choose the inputs to test the program using a testing strategy. Therefore there is no correlation among the faults except the statistic correlation. This type of faults can be treated independently.

In reality, the randomness of choosing the inputs is not always the same during the testing. After a period of testing, test engineers may gain some experience from the previous testing. Then they actually have a fairly good idea about what types of inputs have a high probability of exposing the faults and ought to be chosen and what types of inputs are unlikely to expose the faults and should be chosen later. Therefore, after a period of time, the testing is not as random as the testing at the beginning. The effect of this kind of testing on locating faults is the following: After finding fault i , we will take less time to find fault j than the time taken to find fault j without finding fault i before. Such faults like fault i and fault j are the faults of the second type.

The faults of the last type are those which are masked by other faults. As in the hardware fault testing, one fault may mask another in software fault testing. If fault j is masked by fault i , then in order to find fault j , we have to find fault i first.

We notice that in both the second type and third type of faults, if faults B is correlated to fault A, finding fault B is related to the time at which fault A is found. Therefore, we call the correlations among the faults of the second type and the third type are dynamic correlations. As pointed

out in reference [3], the dynamic correlations have a major effect on the fault exposure ratio especially in the late stage of testing.

Now let us to explain how we model these three types of faults in our theory. In the previous section, we have shown how we model the faults of the first type using the independent fault model. For the faults of second and third types, we will use our correlated faults model to take the dynamic correlation into account. In the correlated faults model, we consider the following correlations: Suppose that we have two correlated faults of the second type, faults i and j . If we locate fault i before fault j is located, then the average time used to locate fault j will be smaller than the average time used to locate fault j without locating fault i first, and vice versa. For the faults of third type, if fault i masks fault j , then fault j can not be located until fault i is located. In the correlated faults model, for the faults of both the second type and the third type, we only consider the one-step correlations between faults. For example, for the two correlated faults of the second type, fault i and fault j , after finding fault i , we find fault j and finding fault j will not affect on finding other faults. For the faults of the third type, we only consider the case in which a masked fault does not mask other faults and is masked by only one fault.

In the correlated faults model, because some of the faults are correlated to each other, we need to use correlation functions to describe this system. In order to state the correlation functions clearly, we use the following conventions to specify the faults and the times. We use indices i , j , and k to specify the faults. If indices i and j appear in the same correlation function, then faults i and j are independent of each other. We use indices i and k_i to denote that fault k_i is correlated to fault i and is located after fault i is found. We use indices n and m for the time tickets, where $n = 1, 2, \dots$ and $m = 1, 2, \dots$. If fault i is located independently, then $t_n(i)$ is used to specify the time when it is located. We use $t_{n,m}(i, k_i)$ to denote the additional time needed to locate fault k_i after fault i is located at time $t_n(i)$. In the other words, $t_n(i) + t_{n,m}(i, k_i)$ is the time when fault k_i is located. Notice that the time interval $t_{n,m}(i, k_i)$ is always greater than or equal to zero.

Two types of correlation function are used. The function $\rho_{ij}[t_n(i), t_m(j)]$ is defined to be the probability of locating both fault i at time $t_n(i)$ and fault j at time $t_m(j)$. The function $\rho_{ik_i}[t_n(i), t_{n,m}(i, k_i)]$ is defined to be the probability that fault i is located at the time $t_n(i)$ and after the time interval $t_{n,m}(i, k_i)$, fault k_i is located. Once we obtain all the correlation functions, all the quantities that we need can be calculated.

In the correlated faults model, we will take another approach to calculate the correlation function. We use Zipf's least effort principle [4, 5] to solve the problem, which states that people always tend to spend least effort to achieve an object. In order to find out how the faults are found in our model, we need to construct the effort function, E . Then the minimum point of the effort function E will give us the behavior of the testing.

We use T to denote the effort function which is proportional to the time spent to locate the faults. From the independent fault model, we know that the average value of t for a fault alone does not reflect how much effort a team spends on locating the fault. Only the ratio of $\langle t \rangle_j$ and \bar{t}_j is proportional to the effort on locating fault j . Therefore, T must be the sum of all the ratios for the faults in

the program. Actually, the function E_I is the effort function for the independent faults model. We can use the form of E_I to obtain the effort function E for our correlated faults model. Let \bar{t}_i denote the average time used to independently locate fault i . Let \bar{t}_{i,k_i} be the average time used to locating fault k_i after fault i is located. Considering the form of E_I in Eq. (12) and the above discussion, we obtain

$$\begin{aligned}
T = & \sum_{i,j,n,m} \bar{t}_i^{-1} \rho_{ij}[t_n(i), t_m(j)][t_n(i) - \bar{t}_i] \\
& + \sum_{i,j,n,m} \bar{t}_i^{-1} \rho_{ij}[t_n(i), t_m(j)][t_m(j) - \bar{t}_j] \\
& + \sum_{\{i,k_i\},n,m} \bar{t}_i^{-1} \rho_{ik_i}[t_n(i), t_{n,m}(i, k_i)] \\
& \quad \times [t_n(i) - \bar{t}_i] \\
& + \sum_{\{i,k_i\},n,m} \bar{t}_{i,k_i}^{-1} \rho_{ik_i}[t_n(i), t_{n,m}(i, k_i)] \\
& \quad \times [t_{n,m}(i, k_i) - \bar{t}_{i,k_i}].
\end{aligned} \tag{20}$$

Here $\{i, k_i\}$ means that the summation is taken over all the possible pairs of correlated faults.

Now let us turn to another aspect of the effort for testing. We understand that if we want to choose a test which can locate a specific fault, we need to expend a lot of effort to do so. On the other hand, if we randomly choose one of all the possible inputs to see if the input exposes any fault, we will spend much less effort to do so. Therefore, the effort function also depends on the randomness of the testing. The more random the testing is, the less effort we spend in the testing. A quantity of measuring the randomness for a system is "entropy". We use S to denote the entropy. In our case, the entropy can be easily calculated and is

$$\begin{aligned}
S = & S_0 - \sum_{i,j,n,m} \rho_{ij}[t_n(i), t_m(i)] \\
& \quad \times \ln \{ \rho_{ij}[t_n(i), t_m(i)] \} \\
& - \sum_{\{i,k_i\},n,m} \rho_{ik_i}[t_n(i), t_{n,m}(i, k_i)] \\
& \quad \times \ln \{ \rho_{ik_i}[t_n(i), t_{n,m}(i, k_i)] \},
\end{aligned} \tag{21}$$

where S_0 is a constant.

From the function E_I in Eq. (12), we see that function E_I contains the two parts we discussed above. Here we assume that the effort function E only contains these two parts. One part is T , which is the effort related to the time spent on the testing and defined in Eq. (20). The other part is the entropy S , which is the effort related to the randomness of the testing and defined in Eq. (21). Then the effort function is

$$\begin{aligned}
E = & T - S + c \sum_{i,j,n,m} \rho_{ij}[t_n(i), t_m(j)] \\
& + c \sum_{\{i,k_i\},n,m} \rho_{ik_i}[t_n(i), t_{n,m}(i, k_i)]
\end{aligned} \tag{22}$$

Here the two terms with a factor c are used to assure

$$\begin{aligned} & \sum_{i,j,n,m} \rho_{ij}[t_n(i), t_m(j)] \\ & + \sum_{\{i,k_i\},n,m} \rho_{ik_i}[t_n(i), t_{n,m}(i, k_i)] = 1. \end{aligned} \quad (23)$$

and Eq. (23) is used to determine the constant c later.

In Eq. (22), the effort function E is the function of $\rho_{ij}[t_n(i), t_m(j)]$'s and $\rho_{ik_i}[t_n(i), t_{n,m}(i, k_i)]$'s. In the testing, the least effort is used. Therefore, we need to find the sets of $p_{ij}[t_n(i), t_m(j)]$'s and $p_{ik_i}[t_n(i), t_{n,m}(i, k_i)]$'s such that the effort function E reaches its minimum. These sets can be obtained by starting with calculating the derivatives of the function E with respect to $\rho_{ij}[t_n(i), t_m(j)]$'s and $\rho_{ik_i}[t_n(i), t_{n,m}(i, k_i)]$'s and let each of them be zero. Then Eq. (23) is used to determine the value of c . This procedure is similar to what we did in the multiple faults mode. Therefore, we skip all the calculations and give the result for the continue time directly:

$$p_{ij}(t, \tau) = Z^{-1} \exp(-t/\bar{t}_i - \tau/\bar{t}_j), \quad (24)$$

$$p_{ik_i}(t, \Delta t) = Z^{-1} \exp(-t/\bar{t}_i - \Delta t/\bar{t}_{i,k_i}), \quad (25)$$

where

$$Z = \sum_{i,j} \bar{t}_i \bar{t}_j + \sum_{\{i,k_i\}} \bar{t}_i \bar{t}_{i,k_i}.$$

The meaning of $p_{ij}(t, \tau) dt d\tau$ is the probability that fault i is located during the time from time t to time $t + dt$ while fault j is located during the time from time τ to time $\tau + d\tau$. $p_{ik_i}(t, \Delta t) dt d\Delta t$ is the probability that fault i is located during the time from time t to time $t + dt$ while fault k_i is located during the time from time $t + \Delta t$ to time $t + \Delta t + d\Delta t$.

Now let us start to calculate the fault exposure ratio using Eq. (18). Let $P(t)$ be the probability of locating any of the faults from time 0 to time t . We use $P_i(t)$ to denote the probability of fault i being found from time 0 to time t and $P_{ik_i}(t)$ to denote the probability of locating both fault i and fault k_i from time 0 to time t . Then the function $P(t)$ can be calculated using the following equation:

$$P(t) = \sum_i P_i(t) + \sum_{\{i,k_i\}} P_{ik_i}(t). \quad (26)$$

According to the definition for the functions $P_i(t)$ and $P_{ik_i}(t)$, we calculate the functions $P_i(t)$ and $P_{ik_i}(t)$ using the correlation functions $p_{ij}(t, \tau)$ and $p_{ik_i}(t, \Delta t)$ in Eqs. (24) and (25) respectively and obtain

$$\begin{aligned} P_i(t) &= \sum_j \int_0^t dt' \int_0^\infty d\tau p_{ij}(t', \tau) \\ &= Z^{-1} \left(\sum_j \bar{t}_j \bar{t}_i [1 - \exp(-t/\bar{t}_i)] \right) \end{aligned} \quad (27)$$

and

$$\begin{aligned} P_{ik_i}(t) &= \frac{1}{2} \int_0^t dx \int_{-x}^x dy p_{ik_i}[(x+y)/2, (x-y)/2] \\ &= Z^{-1} \bar{t}_i \bar{t}_{i,k_i} \\ &\quad - Z^{-1} \frac{\bar{t}_i \bar{t}_{i,k_i}}{\bar{t}_{i,k_i} - \bar{t}_i} [\bar{t}_{i,k_i} \exp(-t/\bar{t}_{i,k_i}) - \bar{t}_i \exp(-t/\bar{t}_i)]. \end{aligned}$$

(28)

Substituting Eqs. (27) and (28) into Eq. (26) and using Eq. (18), we obtain

$$\begin{aligned} K &= B^{-1} T_L \left\{ \sum_{j,i} \bar{t}_j \exp(-t/\bar{t}_i) \right. \\ &\quad \left. + \sum_{\{i,k_i\}} \frac{\bar{t}_i \bar{t}_{i,k_i}}{\bar{t}_{i,k_i} - \bar{t}_i} [\exp(-t/\bar{t}_{i,k_i}) - \exp(-t/\bar{t}_i)] \right\}, \end{aligned} \quad (29)$$

where

$$\begin{aligned} B &= \sum_{j,i} \bar{t}_j \bar{t}_i \exp(-t/\bar{t}_i) \\ &\quad + \sum_{\{i,k_i\}} \frac{\bar{t}_i \bar{t}_{i,k_i}}{\bar{t}_{i,k_i} - \bar{t}_i} [\bar{t}_{i,k_i} \exp(-t/\bar{t}_{i,k_i}) - \bar{t}_i \exp(-t/\bar{t}_i)]. \end{aligned}$$

Comparing Eqs. (19) and (29), we can see that if there is no correlations, Eq. (29) will recover to the fault exposure ratio for independent faults model. Let $\alpha^{-1} = \max\{\bar{t}_j\}$ and $\beta^{-1} = \max\{\bar{t}_{i,k_i}\}$. For large t , we have

$$K \Big|_{t=\infty} = \min\{\alpha, \beta\} T_L.$$

We see that the asymptotic behaviors of the fault exposure ratio for both the correlated faults model and the independent faults model are the same. Both of them are equal to the largest average time of the faults in the program.

In order to see what role the dynamic correlations between faults play in the fault exposure ratio, we rewrite Eq. (29) in another form. Let us first define $R_I(t)$, $R_C(t)$, and $K_C(t)$ as followings:

$$R_I \equiv \sum_{j,i} \bar{t}_j \bar{t}_i \exp(-t/\bar{t}_i),$$

$$\begin{aligned} R_I(t) &\equiv -[R_I(t)]^{-1} \frac{dR_I(t)}{dt} T_L \\ &= \frac{\sum_{j,i} \bar{t}_j \exp(-t/\bar{t}_i)}{\sum_{j,i} \bar{t}_j \bar{t}_i \exp(-t/\bar{t}_i)} T_L \end{aligned} \quad (30)$$

$$R_C(t) \equiv \sum_{\{i,k_i\}} \frac{\bar{t}_i \bar{t}_{i,k_i}}{\bar{t}_{i,k_i} - \bar{t}_i} [\bar{t}_{i,k_i} \exp(-t/\bar{t}_{i,k_i}) - \bar{t}_i \exp(-t/\bar{t}_i)],$$

and

$$\begin{aligned} K_C(t) &\equiv -[R_C(t)]^{-1} \frac{dR_C(t)}{dt} T_L \\ &= T_L \frac{\sum_{\{i,k_i\}} \frac{\bar{t}_i \bar{t}_{i,k_i}}{\bar{t}_{i,k_i} - \bar{t}_i} [\exp(-t/\bar{t}_{i,k_i}) - \exp(-t/\bar{t}_i)]}{\sum_{\{i,k_i\}} \frac{\bar{t}_i \bar{t}_{i,k_i}}{\bar{t}_{i,k_i} - \bar{t}_i} [\bar{t}_{i,k_i} \exp(-t/\bar{t}_{i,k_i}) - \bar{t}_i \exp(-t/\bar{t}_i)]}. \end{aligned} \quad (31)$$

We express Eq. (19) in terms of $R_I(t)$, $K_I(t)$, $R_C(t)$, and $K_C(t)$:

$$K = \frac{R_I(t)K_I(t) + R_C(t)K_C(t)}{R_I(t) + R_C(t)}. \quad (32)$$

To understand Eq. (32), we first have a look at the meanings of $R_I(t)$, $K_I(t)$, $R_C(t)$, and $K_C(t)$. $R_I(t)$ is proportional to the probability that any of the independent faults is not found during the time from 0 to t . All the faults in $R_I(t)$ are independent of each other. Therefore, the function $R_I(t)$ behaves the same as the function $R(t)$ for the independent faults model. According to the definition of the function $K_I(t)$, we know that $K_I(t)$ is the corresponding fault exposure ratio for $R_I(t)$ and is the fault exposure ratio for the independent faults. Comparing Eqs. (19) and (30), we find that $K_I(t)$ shares the same formula with the fault exposure ratio for the independent faults model. This is expected because these two equations actually describe the same thing and both are the fault exposure ratio for the independent faults. As we discussed in the independent faults model, the function $K_I(t)$ is a decreasing function of time t . This means that in the correlated faults model, the independent faults are more and more difficult to be found as in the independent faults model.

In Eq. (31), the summations are over all the pairs of correlated faults. It is clear that $R_C(t)$ is proportional to the probability of not locating any correlated faults during the time from 0 to t . Eq. (31) defines the fault exposure ratio, $K_C(t)$, for the correlated faults. In the fault exposure ratio $K_C(t)$, only the events of two correlated faults being found in certain periods are taken into account. To locate a correlated fault k_i , we must locate the independent fault i first. Thus, the function $R_C(t)$ decreases much slower than the function $R_I(t)$ does at the beginning of the testing. This can be seen from the values of $K_C(0)$ and $K_I(0)$. Letting $t = 0$ in Eqs. (30) and (31), we obtain that $K_C(0)$ is zero and $K_I(0)$ takes a nonzero value. The zero value of $K_C(0)$ means that at the beginning of the testing, a correlated fault is much more difficult to be found than an independent fault. Because $K_C(0) = 0$ and $K_C(t)$ is always greater or equal to zero, $K_C(t)$ must be an increasing function of t at least for small values of t . This is the major difference between the independent faults and the correlated faults. The consequence of this behavior is that at the beginning of the testing, the correlated faults are very difficult to be located. Later, when testing engineers have more experience, the correlated faults are easier and easier to be located. This behavior is just the opposite to the behavior of independent faults.

Now we are ready to understand Eq. (32). K in Eq. (32) is the total fault exposure ratio. Eq. (32) tells us that the total fault exposure ratio is simply calculated from the average of the fault exposure ratios for the independent faults $K_I(t)$ and the correlated faults $K_C(t)$ weighted by $R_I(t)$ and $R_C(t)$, respectively. Because $K_I(t)$ is a decreasing function of t and $K_C(t)$ is an increasing function of t in a certain area of t , unlike the fault exposure ratio for the independent faults model, a complicated behavior of K is expected.

Now let us use an example to show the above discussions graphically. In this example, we consider that there are four faults in a program. We label them using numbers 1, 2, 3, and 4, respectively. In the numerical calculations, we set T_L to be the unit time.

In the correlated faults model, we let faults 1 and 2 be independent faults and faults 3 and 4 be the correlated faults with fault 2. Faults 3 and 4 cannot be located until fault 2 is located. The average times are $\bar{t}_1 = 1.0$, $\bar{t}_2 = 4.0$, $\bar{t}_{2,3} = 5.0$, and $\bar{t}_{2,4} = 10.0$, respectively.

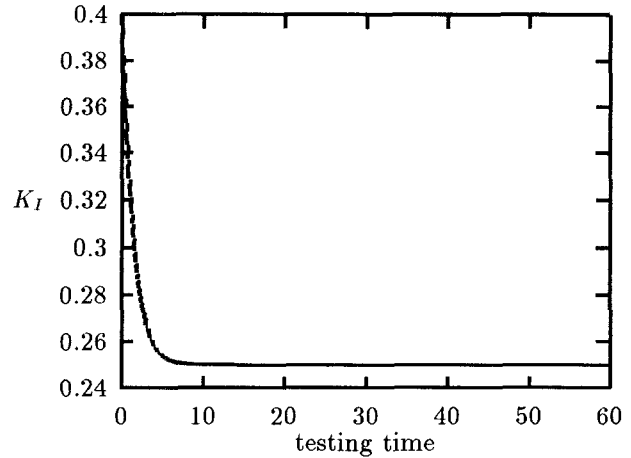


Figure 1. The vertical axis is the fault exposure ratio for faults 1 and 2.

Figure 1 shows the curve of K_I v.s. t . Only the two independent faults, faults 1 and 2, contribute to K_I . K_I decays very fast to a nonzero value. This is because the largest average time in these two faults is 4.0. For large values of t , the value of K_I approaches to the inverse of the largest average time, $1/4.0$.

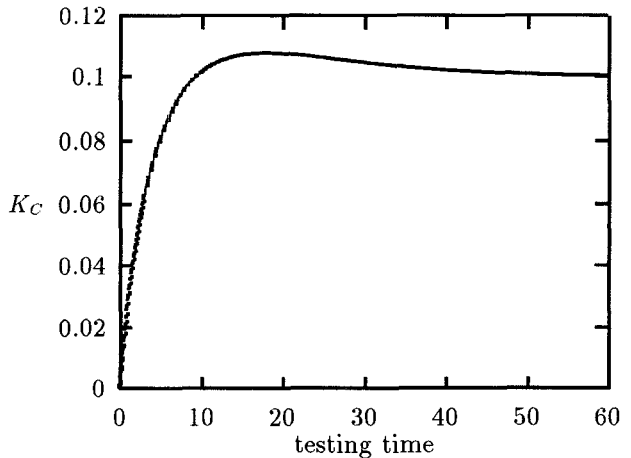


Figure 2. The vertical axis is the fault exposure ratio for faults 3 and 4 correlated to fault 2.

Figure 2 shows the curve of K_C v.s. t . In K_C , two pairs of correlated faults contribute to K_C . One pair is faults 2 and 3, and the other pair is faults 2 and 4. Comparing figures 1 and 2, we can see that there is a big difference between K_I and K_C . The value of K_C at $t = 0$ is zero and then increases as t until $t \approx 16.0$. At the beginning of the testing, because fault 2 has not been located yet, it is impossible to locate either faults 3 or 4. The fault exposure ratio for locating faults 3 and 4 is zero. As the testing continues, the probability of locating fault 2 increases. Then

locating faults 3 or 4 becomes possible and the fault exposure ratio K_C takes a nonzero value. The longer the testing time is, the bigger the probability of locating fault 2 is, and the faster faults 3 and 4 are located. So the fault exposure ratio K_C increases with time t . In the region of $t > 16.0$, K_C decays to the value of the inverse of the biggest average time, $1/10.0$. At this time, fault 2 has already been located. Thus locating faults 3 and 4 is similar to locating the independent faults. The fault exposure ratio K_C approaches the value of the inverse of the biggest average time in faults 3 and 4 as the fault exposure ratio for independent faults K_I behaves.

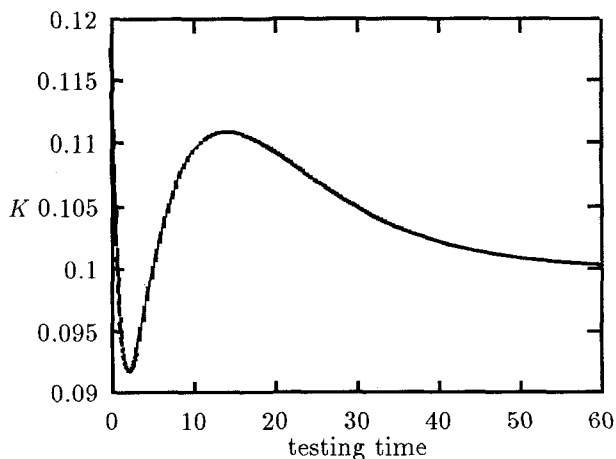


Figure 3.1. The vertical axis is the fault exposure ratio for all the four faults.

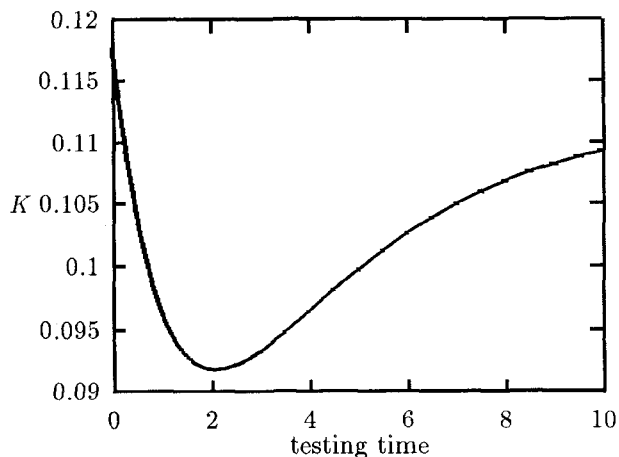


Figure 3.2. The vertical axis is the K for all the four faults. It is the same as Figure 3.1 but with a different scale in the horizontal axis.

We plot the fault exposure ratio K for the correlated faults model against time t in figure 3. Because the fault exposure ratio K contains all the four faults, it is combined from the feature of K_I and the feature of K_C . In the earlier stage of the testing, the independent faults are much easier to be located than the correlated faults are. The fault exposure ratio for the independent faults dominates the fault exposure ratio K . In this period of time, the fault exposure ratio K behaves the same as the fault exposure

ratio K_I for the independent faults does. After this period of time, the testing process enters the later stage of the testing. In the later stage of the testing, the independent faults have been located and some experience from the earlier stage of the testing has been obtained. Then the correlated faults start to be located. Therefore, in the later stage of the testing, the fault exposure ratio K_C for the correlated faults dominates the fault exposure ratio K . The fault exposure ratio K starts to increase with time t as we observed in figure 2.

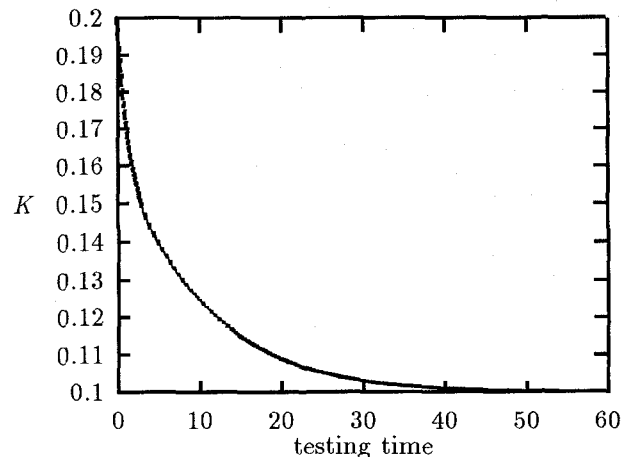


Figure 4. The vertical axis is the fault exposure ratio for the four independent faults.

Figure 4 shows the fault exposure ratio K for the independent faults model. In figure 4, there are four independent faults in a program. The average times of locating faults 1, 2, 3, and 4 are 1.0, 4.0, 5.0, and 10.0, respectively. Comparing figures 1 and 4, we can see that the fault exposure ratio K for the independent faults model and the fault exposure ratio K_I behave the same. Unlike the fault exposure ratio for correlated faults model, the fault exposure ratio for independent faults model is always a decreasing function of t . The faults are more and more difficult to be located as the testing continues. Again, for large values of t , the fault exposure ratio approaches to the value of the inverse of the largest average time, $1/10.0$. This is due to the fact that in independent faults model, the testing is random. At the beginning of the testing, there are more faults in the program. At this time, it is relatively easy to locate any of the faults. Later, some of the faults have been located and most of the rest of the faults have larger average times of being located or are more difficult to be located. Thus, the fault exposure ratio always decreases against time t .

In order to see how well the fault exposure ratio for the correlated faults model describes the reality, we use Eq. (29) to fit the experimental data using the least square fit. The experimental data comes from reference [3]. In the fitting, we assume that N_r faults with the same average time are masked by one fault and there are N_i independent faults with the same average time. Both N_r and N_i are much greater than 1. In the fitting function, there are five parameters. Figure 5 shows both the experimental data and the curve from the fitting. We can see that the experimental data can be described very well by the curve.

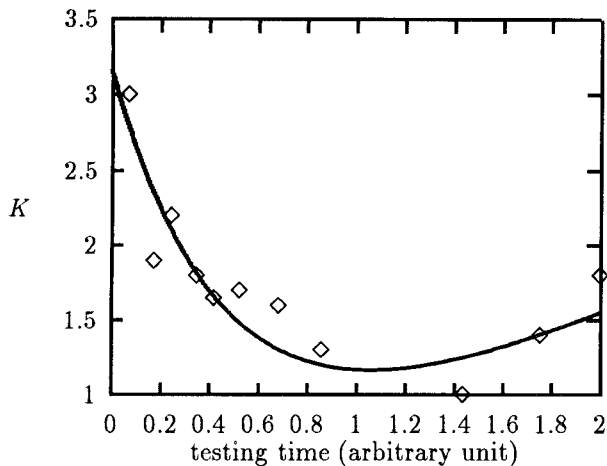


Figure 5. The vertical axis is the relative fault exposure ratio. \diamond denotes the experimental data and the solid line is from the correlated faults model.

We conclude that the correlated faults model contains the main factor of affecting the fault exposure ratio. The results from the correlated faults model show that there are two stages in the testing process. We call them the earlier stage and the later stage. These two stages are distinguished by which type of faults dominates the fault exposure ratio. In the earlier stage, the independent faults dominate the fault exposure ratio and the fault exposure ratio decreases against time t . In the later stage, the correlated faults dominate the fault exposure ratio and the fault exposure ratio increases with time t for a certain period of time.

4 Conclusions

In this paper, we discussed the effect of both random testing and nonrandom testing on software reliability analytically. The independent faults model is used to model the random testing. In this fault model, the fault exposure ratio of the independent faults model approaches to a constant as the time approaches to infinity. This behavior characterizes the testing process in the earlier stage of testing. The same conclusion has been made in reference [3] using a different fault model.

We analyzed the effect of the nonrandom testing on locating faults. According to our analysis, the nonrandom testing can be considered as the correlations between faults. Based on the analysis, the correlated faults model is proposed. The independent faults model is a special case of the correlated faults model.

We find that the behaviors of the fault exposure ratios for the independent faults and the correlated faults are significantly different. We find that there are two phases in the testing process. This agrees with the observation from the experimental data [3].

We also used the fault exposure ratio for the correlated faults model to fit the experimental data. It fits the experimental data very well. We conclude that the dynamic correlations between faults indeed play an important role in the testing as Malaiya, et.al. pointed out [3]. In the early stage of testing process, the independent faults dominate the fault exposure ratio; In the later stage of testing

process, the correlated faults dominate the fault exposure ratio.

From the theoretical aspect, the correlated faults model has two advantages. The first is that after the correlations between faults are taken into account, the testing can be treated as random testing. This can largely simplify the calculation. The second is that the fault maskings are automatically taken into account in the correlated faults model.

This paper also provides an approach to calculate the probability of locating faults. For example, if the statistic correlation is considered, it is still possible to calculate the probabilities using this approach. We speculate that this approach will be used in solving more complicated fault models in the future.

Acknowledgements

We thank Prof. R. Mark Bradley for his support of this research and Mr. Naixin Li for helpful discussions. This work was partially supported by an SDIO/IST funded project monitored by ONR and the NSF Grant No. DMR-9100257.

References

- [1] J.D. Musa, A. Iannino and K. Okumoto, *Software Reliability: Measurement, Prediction, Application*, McGraw-Hill, 1987.
- [2] J.D. Musa, "Rationale for Fault Exposure Ratio," *ACM SIGSOFT Software Engineering News*, pp. 79, July 1991.
- [3] Y. K. Malaiya, A. von Mayrhauser and P. K. Srimani, "The Nature of Fault Exposure Ratio," *Proc. of Symposium on Software Reliability Engineering*, pp. 23-32, May 1992.
- [4] G. K. Zipf, *Human Behavior and the Principle of Least Effort*, Addison-Wesley, 1949.
- [5] M. Trachtenberg, "Why Failure Rates Observe Zipf's Law in Operational Software," *IEEE Trans. Reliability*, Vol. 41, pp. 383-389, September 1981.