

Assessing the Reliability of Diverse Fault-Tolerant Systems

Bev Littlewood, Peter Popov, Lorenzo Strigini
Centre for Software Reliability, City University, London
E-mail: {B.Littlewood,L.Strigini,PTP}@csr.city.ac.uk
phone: 020 7477 8420, fax: 020 7477 8585

Abstract

Design diversity between redundant channels is a way of improving the dependability of software-based systems, but it does not alleviate the difficulties of dependability *assessment*. Assuming failure independence between channels is unrealistic. Using statistical evidence from realistic testing, standard inference procedures can estimate system reliability, but they take no advantage of a system's fault-tolerant structure. We show how to extend these techniques to take account of fault tolerance by a conceptually straightforward application of *Bayesian* inference. Unfortunately, the method is computationally complex and requires the conceptually difficult step of specifying 'prior' distributions for the parameters of interest. This paper presents the correct inference procedure, exemplifies possible pitfalls in its application and clarifies some non-intuitive issues about reliability assessment for fault-tolerant software.

1. Introduction

Design diversity between the redundant channels of a fault-tolerant architecture appears to be an effective way of improving the dependability of software-based systems [Littlewood *et al.* 2000b]. However, it does not simplify the problem of assessing the reliability or safety of a specific system, e.g. for the purposes of licensing.

Consider for instance a two-channel, 1-out-of-2, software-based diverse system, as could be for instance a protection system (Fig. 1) (we will use this example throughout our discussion).

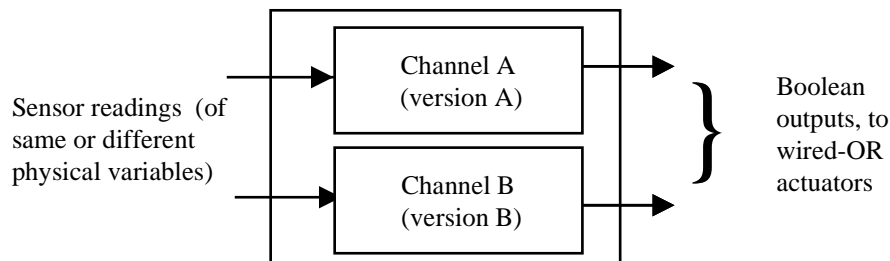


Fig. 1. Our example system

Estimating its probability of failure per demand (*pdf*) would be simplest if we could assume independence between failures of the two channels. Then, we could just assess the *pdf* of the two channels separately and multiply them together. Evidence of even modest reliability of the channels would suffice to claim much higher reliability for the system. But assuming independent failures has been shown to be completely unrealistic by both experiments [Knight & Leveson 1986] and theoretical modelling [Littlewood *et al.* 2000b]. Positive correlation between channel failures should normally be expected, essentially because, for the builders of diverse versions of a program, some demands will be more difficult - more error-prone - than others. So, even if diverse versions (channel software designs) are produced 'independently', their failures are more likely to happen on certain demands than on others, which leads to positive correlation. What is worse, research has found no simple way of setting an upper bound for the correlation between failures of the two channels. So, it is necessary actually to evaluate the *pdf* of the two-channel system as a whole.

The simplest way to assess the reliability of a system - fault tolerant or otherwise - is to observe its failure behaviour in (real or simulated) operation. If we treat the fault-tolerant system as a black box (Fig. 2a), i.e., we ignore the fact that it is indeed fault-tolerant, we can apply standard techniques of statistical inference to estimate its *pdf* on the basis of the amount of realistic testing performed and the number of failures observed. However,

this ‘black-box’ approach to reliability estimation has severe limitations [Littlewood & Strigini 1993], [Butler & Finelli 1991]: if we want to demonstrate very small upper bounds on the *pdf*, the amount of testing required becomes very expensive and then infeasible. It is then natural to ask whether we can use the additional knowledge that we are dealing with a fault-tolerant system to reduce this problem - to achieve better confidence for the same amount of testing.

We reasonably assume that we can observe whether either channel fails, so that testing produces evidence about the reliability of each channel by itself as well as of the whole system. Thus, we treat the system as a ‘white-box’ (Fig. 2b). In addition, we have a priori knowledge about the effect of the channels' failures on system failure: we know that we are dealing with a 1-out-of-2 system. In short, we have much more information than in the ‘black box’ scenario. We may hope that this additional information can be used to reduce the uncertainty about system reliability. This is the problem which we address in this paper.

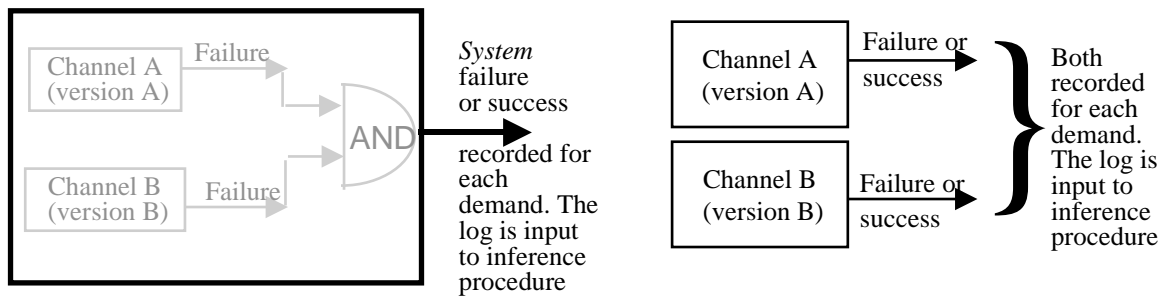


Fig. 2. Black-box vs. white-box inference

We first briefly introduce Bayesian inference: the more widely known approach of "classical" inference, which typically employs different ad hoc methods for different inference problems, does not seem suitable in this case in which we wish to perform inference in a consistent way on various aspects of a system. We then describe the procedure for applying Bayesian inference to our 2-channel system. Bayesian inference presents two kinds of difficulty: conceptually, it depends on the user specifying "prior" probability distributions, which many people find difficult to specify; computationally, it can be very demanding, requiring numerical computation of complex integrals. Various methods are commonly applied to reduce both difficulties. In the rest of the paper we proceed to discuss both some standard method and some apparently promising ad hoc methods. It turns out that none of these will be useful in all cases.

2 Bayesian inference

In our scenario, we count the demands to the system and the failures (of one or both channels) observed, and from this information try to predict the probability of failures on a future demand. This is a problem of statistical inference. Standard techniques for statistical inference are divided into "classical" and "Bayesian". Their applications to estimating the reliability of a system as a black box (i.e., ignoring how it behaves internally - in our case, ignoring that one channel may fail without the whole system failing) is standard textbook material. The classical methods produce "confidence" statements, like "we have 95% confidence that the probability of failure per demand is less than 10^{-3} ". Classical inference is the more widely known approach, but it has drawbacks. The meaning of a "confidence level" is defined in terms of the experiment that produced it, and cannot be translated into probabilities for events of actual interest, e.g. system failure over a pre-specified duration of operation. It is difficult or meaningless to compare values of confidence bounds and confidence levels obtained for different systems or under different regimes of observation, and to devise a classical inference procedure for a system described by multiple parameters, as is our case.

The Bayesian approach, on the other hand, produces probability statements, which we can combine to derive probabilities of other events of interest. Furthermore, inference procedures for any situations can be easily derived from the general approach.

In our case, the Bayesian approach considers that the actual *pdf* of the system is unknown, and thus treats it as a random variable. In a sense, the system that one is trying to evaluate was extracted at random from a population of possible systems, with different reliabilities and different probabilities of being actually produced. Any one of these *could* have been delivered, as far as the observer can tell from the available information. Reliability estimation consists, roughly speaking, in deciding whether the actual system is, among this population, one of those with a high *pdf* or with low *pdf*. This population is described by a *prior* probability distribution: for each

possible value of the *pdf*, a probability is stated that the system has that value of *pdf* (more precisely, a *probability density function* is specified). This prior distribution must describe the knowledge available before testing. Then, the frequency with which we observe failures gives us reason to alter this probability distribution. For instance, passing a certain number of tests shows that the system is less likely to be one with very high *pdf*. Bayes's theorem completely specifies the changes in probabilities as a function of the observations. A *posterior* distribution for the *pdf* is thus obtained, which takes account of the knowledge derived from observation.

With Bayesian inference, one can answer the question 'How likely is it that this software has $pdf \leq 10^{-4}$?' with an actual probability. This can be used in all kinds of reliability calculations. One can also, given the probability distribution for the *pdf*, calculate the probability that the software will survive a given number of demands without failures, i.e., the probabilities of events of actual interest.

The Bayesian approach has the advantage of a consistent and rigorous treatment of all inference problems, but in our case we have additional reasons for preferring it over the "classical" approach: we need to produce an inference procedure for a new, non-textbook scenario - a fault-tolerant system; and we need inference about multiple variables (the *pdfs* of the individual channels and of the system) linked by mutual constraints.

However, Bayesian methods present two difficulties. First, although the formulae for the inference are straightforward to derive, the calculations which they require may be very complex, often with no closed-form solution. Numerical solutions may be time-consuming and vulnerable to numerical errors. Fast computers help, but one may need to write ad hoc software.

The second difficulty is more basic. Bayesian inference always requires one to start with prior probability distributions for the variables of interest: it (rightly) compels us to state the assumptions that we bring to the problem. But formulating the prior distributions may require somewhat subtle probabilistic reasoning. The prior distribution must be one that the assessor does consider a fair description of the uncertainty about the system before the system is tested. Even experts in a domain may find it very difficult to specify their prior beliefs in a mathematically rigorous format. In some cases, if undecided between alternative priors, the only practical solution may be to adopt the more *pessimistic* one. The difficulty may be alleviated by checking how sensitive the predictions are to the variation between the different priors that appear plausible. As observations accumulate, they may start to "speak for themselves", making the differences in the priors irrelevant. Statisticians have developed various ways for simplifying both problems (computational complexity and difficulty in specifying priors). In our discussion we will consider the most popular among such general "tricks", as well as some ad hoc ones.

3 Problem statement and Bayesian inference procedure

We consider the system of Fig. 1, subjected to a sequence of n independent demands.

If we treat the system as a black box, i.e. we can only observe *system* failure or success (Fig. 2a), the inference proceeds as follows. Denoting the probability of failure on demand for the system as p , the posterior distribution of p after seeing r failures in n demands is:

$$f_p(x|r, n) \propto L(n, r|x) f_p(x), \tag{1}$$

where $L(n, r | x)$ is the *likelihood* of observing r failures in n demands if the *pdf* were exactly x . This is given in

this case by the *binomial* distribution, $L(n, r|x) = \binom{n}{r} x^r (1-x)^{n-r}$. $f_p(\bullet)$ is the prior distribution of p , which

represents the assessor's beliefs about p , before seeing the result of the test on n demands.

(1) is the general form of Bayes's formula, applicable to any form of the likelihood and any prior.

In the white-box scenario, instead, we can discriminate among four different possible outcomes for each demand: We use these notations:

Event	Version A	Version B	Number of occurrence in n tests	Probability
α	fails	fails	R_1	P_{AB}
β	fails	succeeds	R_2	$P_B - P_{AB}$
γ	succeeds	fails	R_3	$P_A - P_{AB}$
δ	succeeds	succeeds	R_4	$1 - P_A - P_B + P_{AB}$

The probability model now has the four parameters shown in the last column of the table, but since these four probabilities sum to unity, there are only three degrees of freedom: the triplet P_A, P_B and P_{AB} completely specifies the model. An assessor will need to specify a *joint* prior distribution for these three parameters,

$$f_{P_{AB}, P_A, P_B}(x, y, z).$$

The likelihood of observing r_1 common failures of both channels, r_2 failures of channel A only and r_3 failures of channel B only in n tests is now given by a *multinomial function*:

$$L(r_1, r_2, r_3, n | P_{AB}, P_A, P_B) = \frac{n!}{r_1! r_2! r_3! (n - r_1 - r_2 - r_3)!} P_{AB}^{r_1} (P_B - P_{AB})^{r_2} (P_A - P_{AB})^{r_3} (1 + P_{AB} - P_A - P_B)^{n - r_1 - r_2 - r_3} \quad (2)$$

The posterior distribution, similarly to (1), is:

$$f_{P_{AB}, P_A, P_B}(x, y, z | r_1, r_2, r_3, n) \propto L(r_1, r_2, r_3, n | P_{AB}, P_A, P_B) f_{P_{AB}, P_A, P_B}(x, y, z) \quad (3)$$

Given a joint distribution for P_A, P_B, P_{AB} , we can always deduce the distribution $f_{P_{AB}}$ of the system *pdf*, by integrating out P_A and P_B . So, for a given *prior* joint distribution, there are two options for inferring system reliability from the test results. In the white-box method, we obtain the posterior joint distribution via (3) and then deduce the posterior $f_{P_{AB}}$ from this. We can also apply the black-box method: we first derive the prior $f_{P_{AB}}$ and then update it to obtain a posterior via (1). Comparing the two results will be for us a way of comparing the two methods.

How to solve these formulas is clear even though it may be computationally expensive. There remains the problem of specifying prior distributions, which we address in the next section.

4 Prior distributions

Here we study ways of specifying prior distributions. Our main concern is to help assessors to specify priors, by imposing a useful structure for their interrelated beliefs about the *pdfs* of the channels and of the system. A useful side effect is often a simplification of the calculations. We omit the mathematical details and concentrate on the practical conclusions; a more mathematical and more detailed discussion is available in [Littlewood *et al.* 2000a].

4.1 Dirichlet distribution

It is common in Bayesian statistics to use a *conjugate family* of distributions to represent prior beliefs. This term denotes a parametric family of distributions that has the property for a particular problem (i.e. likelihood function) that if an assessor uses a member of the family to represent his/her prior beliefs, then the posterior will automatically also be a member of the family. If a conjugate family exists for a certain likelihood function (this is not always the case), it is unique. For our white-box scenario, the conjugate family is that of *Dirichlet* distributions.

It turns out that, with a Dirichlet prior, the posteriors for the probability of system failure derived via the ‘white-box’ and via the ‘black-box’ methods are identical, no matter what we observed. In other words, whatever the detailed failure behaviour of the two channels, there is no benefit from taking this extra information into account in assessing the reliability of the system. So if an assessor’s prior belief is indeed a Dirichlet distribution, there is no advantage in using ‘white-box’ inference. On the other hand, if the assessor’s belief are *not* represented by a Dirichlet distribution, choosing this distribution as a convenient simplification would make it impossible to exploit any potential gain from the white-box inference.

4.2 Prior distributions with known failure probabilities of the versions

Another form of simplification of the prior distribution may be possible if there is a very great deal of data from past operational use for each version (e.g. if they are commercial-off-the-shelf - COTS - items), so that each channel’s probability of failure on demand can be estimated with great accuracy. We can then approximate this situation by assuming that the *pdfs* of the versions are known *with certainty* and are P_{Atrue} and P_{Btrue} . In other words, the uncertainty of the assessor concerns only the probability of *system* failure.

We illustrate this set-up with a few numerical examples, shown in Table 1. In each case we assume that $P_{Atrue} = 0.001$, $P_{Btrue} = 0.0005$. Clearly, a 1-out-of-2 system will be at least as reliable as the more reliable of the two versions, so the prior distribution of the system *pdf* is zero outside the interval $[0, 0.0005]$. We consider two examples of this distribution: a uniform distribution and a Beta($x, 10, 10$) both constrained to lie within this interval.

We make no claims for ‘plausibility’ for these choices of priors. However, it should be noted that each is quite pessimistic: both priors, for example, have mean 0.00025, suggesting a prior belief that about half channel B failures will also result in channel A failure.

Table 1

Uniform prior $P_{ab} P_a,P_b$		Percentiles				
		10%	50%	75%	90%	95%
Prior		0.00005	0.00025	0.000375	0.00045	0.000475
$r_1=0$	Black Box	0.000011	0.00007	0.000137	0.000225	0.000286
	White Box (version failures)	0.000008	0.00005	0.000095	0.000148	0.00018
	White Box (no version failures)	0.000268	0.00042	0.000462	0.00048	0.000485
$r_1 = 1$	Black Box	0.000045	0.000155	0.000246	0.000342	0.000396
	White Box (version failures)	0.00004	0.00012	0.000179	0.000238	0.000271
$r_1 = 3$	Black Box	0.00015	0.0003	0.000384	0.000443	0.000465
	White Box (version failures)	0.000165	0.000283	0.000343	0.00039	0.000413
$r_1 = 5$	Black Box	0.000235	0.000375	0.000435	0.00047	0.00048
	White Box (version failures)	0.000345	0.00044	0.000469	0.000482	0.000485
Non-uniform prior $P_{ab} P_a,P_b$		Percentiles				
		10%	50%	75%	90%	95%
Prior		0.000175	0.000245	0.000283	0.000317	0.000335
$r_1=0$	Black Box	0.000146	0.000215	0.000253	0.000286	0.000306
	White Box (version failures)	0.00013	0.000188	0.00022	0.00025	0.000269
	White Box (no version failures)	0.000205	0.000278	0.000313	0.000345	0.00036
$r_1 = 1$	Black Box	0.000161	0.000228	0.000265	0.0003	0.000318
	White Box (version failures)	0.00015	0.00021	0.000244	0.000275	0.000291
$r_1 = 3$	Black Box	0.000185	0.000251	0.000287	0.00032	0.000336
	White Box (version failures)	0.000195	0.000255	0.00029	0.00032	0.000335
$r_1 = 5$	Black Box	0.000205	0.000271	0.000305	0.000335	0.000353
	White Box (version failures)	0.00024	0.000304	0.000335	0.00036	0.000375

Table 1: Two groups of results are summarised: with uniform prior and non-uniform prior, $P_{ab}|P_a,P_b=\text{Beta}(x,10,10)$ on the interval $[0, 0.0005]$. The percentiles illustrate the cumulative distribution $P(\theta \leq X) = Y$, where X are the values shown in the table and Y are the chosen percentiles, 10%, 50%, 75%, etc. Rows labelled 'Black box' represent the percentiles, calculated with the black-box inference, those labelled 'White box' show the percentiles calculated for a posterior derived with (3). '(no version failures)' and '(version failures)' refer to two different observations, in which no individual failures of channels and individual channel failures were observed, respectively.

We assume that $n=10,000$ demands are executed in an operational test environment. The rows in Table 1, for each of the two prior distributions studied, differ in the numbers of failures (of each channel and of both together) observed over the 10,000 demands. The rows marked "version failures" describe cases in which the observed numbers of channel A and of channel B failures take their (marginal) expected values, i.e. 10 and 5 respectively. The other case is the extreme one where there are no failures of either channel.

In each case our main interest is in how our assessment of the system reliability based upon the full information, r_1, r_2, r_3 , ("white-box") differs from the assessment based upon the black-box evidence, r_1 , alone.

In Table 1, the first row with $r_1=0$ shows the increased confidence that comes when extensive testing reveals *no system failures*. The black-box posterior belief in the system *pdf* is better (all percentile values are lower) than the prior belief. More importantly, the posterior belief in the ' $r_1=0$, White box (version failures)' rows, based on observing version failures but no system failures, is more optimistic than the black-box posterior. Here the extra information of *version* failure behaviour allows greater confidence to be placed in the system reliability, compared with what could be claimed from the *system* behaviour alone.

The result is in accord with intuition. Seeing no system failures in 10,000 demands, when there have been 10 channel A failures and 5 channel B failures suggests that there is some *negative correlation* between failures of the channels: even if the channels were failing *independently* we would expect to see some common failures (the expected number of common failures, conditional on 10 A s and 5 B s, is 2.5).

The rows where ($r_1 \neq 0$) show what happens when there are system failures (common failures of the versions), with the same numbers of version failures (10 A s, 5 B s). As would be expected, the more system failures there are on test, the more pessimistic are the posterior beliefs about system reliability. More interesting, however, is the relationship between the black-box and white-box results. Consider the rows with ($r_1=5$). These rows of Table 1 represent the most extreme case, in which all demands that are channel B failures are also channel A failures. This would suggest strongly that there is *positive correlation* between the failures of the two versions. Here the black-

box method gives results that are too optimistic compared with those based on the complete (white-box) failure data.

These results show that ‘white-box’ inference can produce advantages (albeit small ones in this example).

However, this table also shows a consequence of our simplifying assumption (perfectly known channel *pdfs*) that is clearly wrong. When $r_1=0$, that is there have been no failures of either version (and hence no system failures), the posterior distribution of the *pdf* is worse than it was *a priori*. How can the observation of such ‘good news’ make us lose confidence in the system?

The reason for this paradox lies in the constraints on the parameters of the model that are imposed by assuming the versions reliabilities are known with certainty. Consider Table2:

Table 2

	A fails	A succeeds	<i>total</i>
B fails	θ	$P_B - \theta$	P_B
B succeeds	$P_A - \theta$	$1 - P_A - P_B + \theta$	$1 - P_B$
<i>total</i>	P_A	$1 - P_A$	1

There is only one unknown parameter, θ , the system *pdf*, which appears in all the cells above representing the four possible outcomes of a test. If we observe no failures in the test, this makes us believe that the entry in the (A succeeds, B succeeds) cell, $1 - P_A - P_B + \theta$, is large. Since P_A, P_B are known, this makes us believe that θ is large.

Of course, it could be argued that observing no version failures in 10,000 demands, with the known version *pdfs* 0.001, 0.0005, is extremely improbable - i.e. observing this result in practice is essentially impossible. This does not remove the difficulty, however: it can be shown that *whatever the value of n*, the ‘no failures’ posterior will be more pessimistic than the prior.

The practical conclusion seems to be that this particular simplified prior is only useful if the number of demands in test is great enough to ensure that at least some version failures are observed.

4.3 Priors allowing conservative claims for system reliability

Here we show that even if P_A and P_B are not known with certainty, assuming that they are can be used to obtain conservative estimates in many cases, and is therefore useful despite the problems described in section 4.2.

Clearly for every prior $f_{P_{AB}, P_A, P_B}(\bullet, \bullet, \bullet)$ (with its corresponding marginal distribution of the probability of system failure, $f_{P_{AB}}(\bullet)$), if we have upper bounds on the probabilities of channel failures, P_{Amax} and P_{Bmax} , we could define a new prior, $f_{P_{AB}, P_A, P_B}^*(\bullet, \bullet, \bullet)$, such that $P_A = P_{Amax}$, with certainty, $P_B = P_{Bmax}$ with certainty, and the probability of system failure is as in the true prior, $f_{P_{AB}}(\bullet)$.

Now we compare the posterior marginal distributions, $f_{P_{AB}}(\bullet | n, r_1, r_2, r_3)$ and $f_{P_{AB}}^*(\bullet | n, r_1, r_2, r_3)$, derived from the same observation ($n : r_1, r_2, r_3$), respectively with the true and the approximated priors, $f_{P_{AB}, P_A, P_B}(\bullet, \bullet, \bullet)$ and $f_{P_{AB}, P_A, P_B}^*(\bullet, \bullet, \bullet)$. We illustrate the relationship between the two posterior distributions in Table 3.

The prior $f_{P_{AB}, P_A, P_B}(\bullet, \bullet, \bullet)$ used in Table 3 is defined as follows:

- $f_{P_A, P_B}(\bullet, \bullet) = f_{P_A}(\bullet) f_{P_B}(\bullet)$, i.e. the prior distributions of P_A and P_B are independent.¹
- The marginal distributions $f_{P_A}(\bullet)$ and $f_{P_B}(\bullet)$ are Beta distributions, $f_{P_A}(\bullet) = \text{Beta}(x, 20, 10)$ and $f_{P_B}(\bullet) = \text{Beta}(x, 20, 20)$ within the interval $[0, 0.01]$: $P_{Amax} = P_{Bmax} = 0.01$.
- The assessor is "indifferent" among the possible values of P_{AB} , i.e.:

$$f_{P_{AB}}(\bullet | P_A, P_B) = \frac{1}{\min(P_A, P_B)} \text{ within } [0, \min(P_A, P_B)] \text{ and } 0 \text{ elsewhere.}$$

The system was subjected to $n = 4000$ tests, and the failures of the system, channel A and B, represented by r_1, r_2 and r_3 , respectively, are shown in the table. The selected examples cover a range of interesting testing results: no failure, no system failure but channel failures, system failure only, a combination of system and channel failures.

¹ The assumption we make can be spelled out as: “Even if I were told the value of P_A , this knowledge would not change my uncertainty about P_B (and vice versa)”. Notice that this assumption is not equivalent to assuming independence between the failures of the two channels, which is well known to be unreasonable. In fact, our assumption says *nothing* about the probability of common failure, P_{AB} .

The percentiles reveal that the simplified prior always gives more pessimistic predictions than the true prior: the probability that the system reliability will be better than any reliability target will be greater with the true prior, $f_{P_{AB}}(\bullet|data)$, than with the simplified one, $f^*_{P_{AB}}(\bullet|data)$.

This observation, if universally true, suggests a relatively easy way of avoiding the difficulty in defining the full $f_{P_{AB},P_A,P_B}(\bullet,\bullet,\bullet)$. If assessors can specify their beliefs about upper bounds on the channels $pdfs$, P_{Amax} and P_{Bmax} , and system pdf , $f_{P_{AB}}(\bullet)$, these can be combined into the simplified prior, $f^*_{P_{AB},P_A,P_B}(\bullet,\bullet,\bullet)$, to obtain conservative prediction.

Table 3 illustrates a small part of the numerical experiments we carried out with different priors and assumed testing results. It presents the *typical* cases in which the observations are consistent with the priors: the number of channel failures are within the variation due to the random failures. In all cases the simplified prior produced more conservative predictions than the true prior. It must be noted, however, that for some extreme case of observations which are not consistent with the priors (their occurrence is virtually impossible with the assumed priors) the conservatism of the simplified prior is not guaranteed. General conditions under which the simplified prior is guaranteed to generate conservatism are yet to be identified.

Table 3

Percentiles		10%	50%	75%	90%	95%
Prior		0.00025	0.00225	0.0035	0.00435	0.00485
$r_1 = 0, r_2 = 0$	$f_{P_{AB}}(\bullet n, r_1, r_2, r_3)$	0.00278	0.003525	0.00406	0.00445	0.00473
$r_3 = 0$	$f^*_{P_{AB}}(\bullet n, r_1, r_2, r_3)$	0.0055	0.00635	0.0067	0.00705	0.00725
Black-box posterior		0	0	0.000125	0.00035	0.0005
$r_1 = 1, r_2 = 0$	$f_{P_{AB}}(\bullet n, r_1, r_2, r_3)$	0.0029	0.00372	0.0042	0.00455	0.0048
$r_3 = 0$	$f^*_{P_{AB}}(\bullet n, r_1, r_2, r_3)$	0.0055	0.00638	0.0067	0.00685	0.00735
Black-box posterior		0	0.00033	0.00058	0.00092	0.00115
$r_1 = 1, r_2 = 24$	$f_{P_{AB}}(\bullet n, r_1, r_2, r_3)$	0	0.0001	0.00035	0.00062	0.00076
$r_3 = 20$	$f^*_{P_{AB}}(\bullet n, r_1, r_2, r_3)$	0.00035	0.00123	0.00178	0.00235	0.00275
$r_1 = 0, r_2 = 20$	$f_{P_{AB}}(\bullet n, r_1, r_2, r_3)$	0	0	0.00022	0.00049	0.00067
$r_3 = 15$	$f^*_{P_{AB}}(\bullet n, r_1, r_2, r_3)$	0.00015	0.00135	0.00224	0.0029	0.00331

Table 3: The percentiles of the prior marginal distribution $f_{P_{AB}}(\bullet)$ and the following three posterior distributions: $f_{P_{AB}}(\bullet|n, r_1, r_2, r_3)$, $f^*_{P_{AB}}(\bullet|n, r_1, r_2, r_3)$ and black-box posterior.

The usefulness of the conservative prior $f^*_{P_{AB},P_A,P_B}(\bullet,\bullet,\bullet)$ seems *limited*. Indeed, for the important special case of testing which does not reveal any failure ($r_1 = 0, r_2 = 0, r_3 = 0$), the conservative result is *too conservative* and hence not very useful: the posterior will be more pessimistic than the prior, due to the phenomenon explained in section 4.2. This fact reiterates the main point of this paper: prior elicitation is difficult and there does not seem to exist easy ways out of this difficulty.

The last two cases presented in Table 3 with testing results ($r_1 = 1, r_2 = 24, r_3 = 20$) and ($r_1 = 0, r_2 = 20, r_3 = 15$), respectively, illustrate the interplay between the black-box and the white-box inferences. In the case with a single system failure the black-box posterior is more pessimistic than the full white-box posterior, while in the case with no system failure the black-box posterior gives more optimistic prediction about system reliability. In the case ($r_1 = 1, r_2 = 24, r_3 = 20$) we have evidence of negative correlation between failures of channels. The expected number of system failures under the assumption of independence is 1.4 in 4000 tests, while we only observed 1. In the case ($r_1 = 0, r_2 = 20, r_3 = 15$) even though no system failure is observed the evidence of negative correlation is weaker (lower number of individual failures is observed). As a result, the white-box prediction is worse than the black-box one.

In summary, using the black-box inference for predicting system reliability may lead to overestimating or underestimating the system reliability.

5. Conclusions

We have studied how to use the knowledge that a system is internally a fault-tolerant system, of which we can observe the individual channels, to improve the confidence in its reliability that we can derive from observing its

behaviour under realistic testing. I.e., we have studied what we call 'white box' inference, in which failures of the channels, masked by fault tolerance, are taken into account, as opposed to 'black box' inference in which they are ignored.

Bayesian inference is the correct method for 'white box' inference about a fault-tolerant systems: we do not see a better way for consistently performing inference about multiple parameters (the *pdfs* of the channels and of the whole systems) linked by reciprocal constraints.

We have described the proper application of this approach to infer the *pdf* of a 1-out-of-2, on-demand system.

Recognising that this method, albeit correct, is practically very cumbersome to apply, we have then looked for ways of simplifying its practical use. The most standard method - using prior distributions from a conjugate family of distributions, which is often a somewhat arbitrary but useful approximation - turns out to be useless for this particular problem as it is equivalent to ignoring the fault-tolerant structure of the system. We have then explored more ad hoc methods for simplifying the correct inference method. In one simplification, the reliabilities of the channels are taken as known with certainty. It turns out that this approximation, plausible in some situations, produces the artefact of counterintuitive (and useless) conclusions in the important case of no observed failures. Last, we showed that even if this assumption is not justified it may be used in some cases (see 4.3) as it seems to allow a conservative approximation that dispenses with the need to specify complete prior distributions: this may be useful in practice, but not universally so.

In conclusion, it is for the time being unavoidable to adapt the application of the inference procedure to the specific case at hand, selecting those specific approximations that work best for the conditions observed. The practical difficulties could be alleviated by better, specialised software tools, relieving the burden of the multiple sensitivity analyses and 'what if' analyses that may be necessary. The main requirement is that such tools must guarantee the necessary numerical precision, to avoid the risk of decisions being driven from mere artefacts of numerical error.

The immediate conclusion is that it is important to be aware of how 'white box' inference should be performed, but in many cases its difficulties will make it unattractive. We expect that in some cases its outcome will appear immediately useful (e.g. as a way of trusting that a certain claimed *pdf* is conservative), and hope that these will lead to improving mathematical techniques and tools and thus reduce the mechanical difficulties of applying the approach. The more basic difficulty - the dependence on prior distributions - is actually at the same time the basic advantage of Bayesian inference: it requires one to make explicit the assumptions underlying the inference activity and it clearly measures how much *added* confidence can really be derived from observing the system's behaviour.

Acknowledgement

This work was supported partially by British Energy Generation (UK) Ltd. under the 'Diverse Software PrOject' (DISPO) and by EPSRC under the 'Diversity In Safety Critical Software' (DISCS) project. The authors wish to thank Martin Newby for helpful discussions.

References

- [Butler & Finelli 1991] R. W. Butler and G. B. Finelli, "The Infeasibility of Experimental Quantification of Life-Critical Software Reliability", in *ACM SIGSOFT '91 Conference on Software for Critical Systems*, in *ACM SIGSOFT Software Eng. Notes*, Vol. 16 (5), New Orleans, Louisiana, pp.66-76, 1991.
- [Knight & Leveson 1986] J. C. Knight and N. G. Leveson, "An Experimental Evaluation of the Assumption of Independence in Multi-Version Programming", *IEEE Transactions on Software Engineering*, SE-12 (1), pp.96-109, 1986.
- [Littlewood *et al.* 2000a] B. Littlewood, P. Popov and L. Strigini, "Assessment of the Reliability of Fault-Tolerant Software: a Bayesian Approach", in *19th International Conference on Computer Safety, Reliability and Security (SAFECOMP 2000)*, Rotterdam, Netherlands, Lecture Notes in Computer Science, Springer-Verlag, 2000a. Also available at http://www.csr.city.ac.uk/people/peter.popov/papers/SAFECOMP2000_Copyright.pdf
- [Littlewood *et al.* 2000b] B. Littlewood, P. Popov and L. Strigini, "Modelling software design diversity - a review", *ACM Computing Surveys*, (to appear) 2000.
- [Littlewood & Strigini 1993] B. Littlewood and L. Strigini, "Validation of Ultra-High Dependability for Software-based Systems", *Communications of the ACM*, 36 (11), pp.69-80, 1993.