

# The Unnecessity of Assuming Statistically Independent Tests in Bayesian Software Reliability Assessments

Kizito Salako, Xingyu Zhao

**Abstract**—When assessing a software-based system, the results of Bayesian statistical inference on operational testing data can provide strong support for software reliability claims. For inference, this data (i.e. software successes and failures) is often assumed to arise in an independent, identically distributed (i.i.d.) manner. In this paper we show how conservative Bayesian approaches make this assumption unnecessary, by incorporating one’s doubts about the assumption into the assessment. We derive conservative confidence bounds on a system’s probability of failure on demand (*pdf*), when operational testing reveals no failures. The generality and utility of the confidence bounds are illustrated in the assessment of a nuclear power-plant safety-protection system, under varying levels of skepticism about the i.i.d. assumption. The analysis suggests that the i.i.d. assumption can make Bayesian reliability assessments extremely optimistic – such assessments do not explicitly account for how software can be very likely to exhibit no failures during extensive operational testing despite the software’s *pdf* being undesirably large.

**Index Terms**—conservative Bayesian inference, CBI, dependability claims, independent software failures, operational testing, software reliability assessment, statistical testing

## 1 INTRODUCTION

CONSIDER a software-based on-demand system subjected to black-box operational testing. During testing, an assessor observes the system – in particular, the software – as it responds to each demand in a random sequence of demands. By noting those demands the software correctly responds to, and those it fails on, the assessor intends to gain enough confidence to support claims of the system being sufficiently reliable<sup>1</sup>. For example, confidence in the system’s unknown *probability of failure on demand (pdf)* –  $X$  say – being sufficiently small. A Bayesian approach to gaining such confidence typically requires 2 things of our assessor: **i)** *prior* to operational testing, the assessor must scrutinise all evidence related to the system’s operational readiness. Via such probing and the assessor’s domain expertise, the assessor forms beliefs about which ranges of *pdf* values are most likely, and which ranges are less so. Examples of evidence include formal analyses of a codebase, the performance of a system during operation [1], the historical performance of similar systems [2], [3], and improvements in software development approaches [4]; **ii)** the assessor must postulate a *statistical model* – in essence, a family of stochastic processes, any of which could characterise the occurrence of the system’s successes and failures during operation. These processes should exhibit statistical properties consistent with the assessor’s beliefs.

For the statistical model it is often assumed that software failures and successes are the outcomes of *independent, identically distributed* (i.i.d.) Bernoulli trials. This is mathematically

- K. Salako is with the Centre for Software Reliability, City, University of London, Northampton Square EC1V 0HB, U.K. (email: k.o.salako@city.ac.uk)
- X. Zhao is with the Department of Computer Science, University of Liverpool, Ashton Street L69 3BX, U.K. (email: xingyu.zhao@liverpool.ac.uk)

Manuscript received August 15, 2022; revised December 21, 2022.

1. This work focuses on software reliability; i.e. only software failures are considered to cause system failure.

convenient and guarantees the *strong law of large numbers* – i.e., operational testing statistics converge *almost surely* to dependability measures of interest (e.g. *pdf*). The i.i.d. assumption can also be reasonable on practical grounds: a typical justification is when demands are rare, and the states/properties of the software and its operational environment are effectively “reset” inbetween these rare demand occurrences. Nevertheless, we must point out that *any* statistical model used in reliability assessment – including one reliant on the i.i.d. assumption – is necessarily a postulate. One does not (cannot?) know with complete certainty that such model assumptions are valid in practice. The skeptical assessor allows for the possibility that the i.i.d. assumption *does not* hold, even if this is unlikely. By incorporating their skepticism into the assessment, our assessor can investigate whether their doubts have a significant impact on their confidence in the system. In the best case, their confidence is insensitive to significant departures from the i.i.d. assumption. But in the worst case, ignoring the slightest failure correlations could lead to seriously misplaced confidence and dangerously optimistic reliability claims.

To aid the skeptical assessor, this work presents a new application of *conservative Bayesian inference* (CBI) techniques for reliability assessments. The paper’s research contributions are:

- 1) incorporating a formal notion of “*doubting*” the i.i.d. assumption into reliability assessments (Section 5);
- 2) a novel CBI technique that accounts for correlated operational testing outcomes (Sections 4, 5 and Appendices B, C);
- 3) a theorem that gives conservative posterior confidence bounds on a system’s *pdf* (Section 5 and Appendices B, D).

The paper’s outline: Section 2 gives critical context, while Section 3 reviews CBI. Section 4 introduces a statistical model for (possibly) dependent system failures/successes. The conservative confidence bounds on *pdf* in Section 5 are applied in Section 6, and discussed in Section 7. Section 8 concludes the paper.

## 2 RELATED WORK

### 2.1 The i.i.d. Assumption in Reliability Assessments

Statistical models with the i.i.d. assumption have a long history of being used in software reliability assessment. Thayer *et al.*'s [5] model was one of the earliest, used in early works on random testing [6]. Regulatory bodies have recommended using the i.i.d. assumption in reliability assessments when appropriate [7].

But there are reasons to doubt the assumption. For instance, the possibility of “failure clustering”; where a system receives sequences of inputs from its operational environment that cause the system to fail. These inputs form trajectories through a system’s *failure region* – the subset of inputs that cause system failure. Failure regions can have topologically interesting properties that allow for failure clustering [8], [9]. These observations informed random testing approaches [10] and approaches for assessing systems with recovery block fault-tolerance [11], [12].

### 2.2 Weakening the i.i.d. Assumption: Statistical Models

Various statistical models that weaken the i.i.d. assumption include: Chen and Mill’s binary Markov chain [13], Goseva-Popstojanova and Trivedi’s Markov renewal process [14] (extended in [15] by Xie *et al.*), and Bondavalli *et al.*’s Markov model [16] with benign-failure states. Classical statistical inference produces “*point estimates*” for these models’ parameters, using only data from operational testing. Such estimates do not reflect an assessor’s uncertainty about whether the i.i.d. assumption holds or not. Indeed, fitted parameter values imply that the models either exhibit dependence, or they do not – there is no room for uncertainty here. In contrast, Bayesian inference – our preferred approach – allows for such uncertainty. In this paper, within a Bayesian framework, we model the system’s failure process as a Markov chain introduced by Klotz [17]. The Klotz model predates, is consistent with, and has (at most) as many states/parameters as, the models in [13], [14].

Another advantage of Bayesian approaches is the assessor’s beliefs (about the unknown *pdf*) are explicitly accounted for in the assessment; beliefs that are justified by various forms of reliability evidence. And, although models of dependent system failures/successes have been developed, none of the assessment approaches using these models provide demonstrably *conservative* statistical support for the skeptical assessor. Particularly when such support is justified by operational testing and other forms of reliability evidence. To the best of our knowledge, ours is the first approach to guarantee conservatism in the face of i.i.d. uncertainty.

### 2.3 (Conservative) Bayesian Methods for Assessments

Bayesian methods have been applied in various assessment scenarios, e.g. [2], [18] involve hierarchical models, while [19]–[22] all use families of Beta prior distributions. More recently, *conservative Bayesian inference* (CBI) methods have been developed to: **i**) address the usual challenge of eliciting a suitable prior distribution from an assessor – a prior that captures all, and only all, of an assessor’s beliefs/views about the *pdf*; **ii**) give support for the most pessimistic reliability claims allowed by the reliability evidence. Thereby, CBI prevents dangerously optimistic claims.

Bishop *et al.* [23] introduced CBI, illustrating its use in assessing safety-critical systems. Povyakalo *et al.* [24] use CBI to obtain the smallest probability of the system surviving future demands, and Salako [25] applies CBI to assessing binary classifiers. While

Flynn *et al.* [26], [27] apply CBI in the assessment of autonomous vehicle safety. Littlewood *et al.* [28] show how CBI supports dependability claims when evidence suggests a new system is an “improvement” over an older system it replaces. Salako *et al.* [29] extend this work, by considering more general “improvement arguments” for a wider range of assessment scenarios. Our application of CBI differs from these in 2 important ways: it allows for dependent testing outcomes during operational testing, and it incorporates skepticism of the i.i.d. assumption into assessments.

CBI methods are closely related to *robust Bayesian analysis* [30]–[33], which studies how the results of Bayesian inference are impacted by uncertainty about the inputs to inference – inputs such as the prior distribution and the statistical model. In particular, Lavine [32] outlines methods that reveal how uncertainty about the statistical model (specifically, about the so-called sampling distribution) impacts inference. This uncertainty is represented by a suitably general joint prior distribution over the family of stochastic processes defined by the statistical model. Subject to constraints on this prior, algorithms give the largest and smallest values for posterior measures of interest – e.g. posterior expectations. Our use of CBI parallels the statistical techniques of Lavine, but applied to reliability assessments. Also, Lavine considers likelihoods consisting of products of the same functional form, while we do not (in order to weaken the i.i.d. assumption).

Draper [34] also tackles the problem of statistical model uncertainty in Bayesian inference, but offers an alternative solution. If one is uncertain about a model’s assumptions – specifically, assumptions that constrain the structural/functional forms of the related family of stochastic processes – one could replace the model with an expanded model. This expanded model encompasses all of the stochastic processes defined under the original model, as well as other stochastic processes that violate the assumptions in question. A suitable prior distribution over this expanded model has to be defined. Draper argues for this Bayesian hierarchical model as a way of addressing model uncertainty. We concur, and further argue for the inference to be conservative; our results are guaranteed to be conservative, while Draper’s are not.

For uncertainty about finitely many alternative models, Pericchi *et al.* [35] use discrete prior distributions over these alternatives. In principle, this is a hierarchical model akin to Draper’s approach, but in more abstract terms within the robust Bayes framework. Our results lie at the intersection of Pericchi *et al.* and Draper’s ideas, within a reliability assessment context.

### 2.4 On-demand vs Continuously Operating Software

This work focuses on assessing on-demand systems: i.e. systems that do not operate continuously, taking action only when certain operating conditions (i.e. demands) arise [36], [37]. Reliability assessments for such systems can use “discrete-time” statistical models (e.g. *Bernoulli processes*) with appropriate reliability measures (e.g. *pdf*). Contrast this with continuously operating software<sup>2</sup> [38]; for assessing *these* systems, it is more appropriate to employ “continuous-time” statistical models (e.g. *non-homogeneous Poisson processes*) and consider reliability measures like failure-rates. *Software reliability growth models* (SRGMs) are an extensive family of stochastic processes used in predicting the future reliability of continuously operating, evolving software (e.g. software with bugs that are discovered and fixed overtime). If bug fixing is successful *without* introducing new software bugs then,

2. Such software *can* have downtime due to, say, maintenance or upgrades.

*ceteris paribus*, the software becomes more reliable with each fix; i.e. reliability “grows”. Singpurwalla and Wilson [39] give an overview of early SRGMs, while Miller [40] details a unified mathematical characterisation of large SRGM subfamilies. Also, see Bergman and Xie’s review of early Bayesian SRGMs [41].

### 3 REVIEW: BAYESIAN RELIABILITY ASSESSMENT

Statistical inference for reliability claims comes in different flavours. The classical “frequentist” confidence statement, e.g. 95% confidence in a *pdf* bound  $b$ , typically means that with a sufficiently large number of i.i.d. tests, there is no more than a 5% chance that the software succeeds on all the tests despite having a *pdf* worse than  $b$ . While the Bayesian approach, instead, treats *pdf* as a random variable, with a “prior” probability distribution representing an assessor’s evidence-based beliefs about the *pdf* before operational testing. The assessor updates their beliefs (via Bayes Theorem) upon seeing testing evidence. This yields a “posterior” distribution. Reliability claims can be made using this posterior; claims that reflect the assessor’s updated judgements. For example, after seeing a sufficiently large number of successful tests, the assessor’s Bayesian “confidence” in a *pdf* bound  $b$  – i.e. their conditional probability of the *pdf* being less than  $b$  – is 95%.

Specifically, we recall the standard Bayesian approach to assessing an on-demand system [7], [42], [43], in which the i.i.d. assumption is adopted. An i.i.d. Bernoulli process represents the stochastic failure behaviour of the system’s software. We denote by  $X$  the system’s unknown *pdf* due to software failures. According to the operational profile [44],  $n$  demands are randomly submitted to the software and no failures are observed (this is the usual requirement when assessing a safety-critical system using operational acceptance testing). Let  $b$  be the required upper bound on *pdf*. Bayesian inference then gives an assessor’s posterior confidence in  $b$  after observing  $n$  tests without failure<sup>3</sup>:

$$P(X \leq b \mid n \text{ demands without failure}) = \frac{P(X \leq b, n \text{ demands without failure})}{P(n \text{ demands without failure})} \quad (1)$$

In practice, Bayesian reliability assessments require that one specifies a prior distribution representing one’s beliefs about the possible values of *pdf*. CBI relaxes this by requiring only a *partial specification* of the prior distribution, when such specifications can be justified by evidence obtained prior to operational testing. Such partial specifications – so-called *prior knowledge* (PK) – take various forms. Most notably, the form of confidence bounds; e.g. being 90% confident that the *pdf* is no greater than  $10^{-3}$ , partly because IEC 61508 Safety Integrity Level 3 [45] was a strict requirement in the system’s development. When an assessor articulates their beliefs as PKs, there is an *infinitely* large set  $\mathcal{D}$  of all prior distributions that satisfy these PKs. CBI then determines the worst support the priors in  $\mathcal{D}$  can give for a reliability claim; e.g. the smallest posterior confidence (1) an assessor can have:

$$\inf_{\mathcal{D}} P(X \leq b \mid n \text{ demands without failure}) \quad (2)$$

The solution of (2) identifies a prior with posterior confidence that is the infimum value; *no other prior that satisfies the PKs can give a smaller value for posterior confidence* (1). In this sense, CBI results are conservative. This “worst case” prior encodes within it the most conservative assessor beliefs consistent with the PKs.

3. ... as well as  $P(\text{failure-free operation} \mid n \text{ demands without failure})$ .

## 4 A STATISTICAL MODEL OF TESTING OUTCOMES

The Klotz model [17] is a stationary random process consisting of possibly dependent Bernoulli trials. As a model of a system’s failure process it generalises the i.i.d. Bernoulli process used in reliability assessments. During operational testing, a random sequence of demands is submitted to the system. On each demand, the system either successfully handles the demand or fails. So, we have a sequence of random variables  $T_1, \dots, T_n$ , each taking the values 0 or 1, corresponding to success or failure, respectively. In this paper we follow the usual convention of upper-case letters for random variables and lower-case for their realisations.

The Klotz model is characterised by a “frequency” parameter  $x$  and a “dependence” parameter  $\lambda$ . Here,  $x$  is the system’s *pdf* while  $\lambda$  is the probability that a failure is followed by another failure. So,  $P(T_1 = 1) = P(T_i = 1) = x$  and  $P(T_i = 1 \mid T_{i-1} = 1) = \lambda$  for  $i = 2, \dots, n$ . By requiring the process be *1st-order stationary*, we have  $P(T_i = 1 \mid T_{i-1} = 0) = \frac{(1-\lambda)x}{1-x}$  for  $i = 2, \dots, n$  (see appendix A), which yields Figure 1. The transitions in Figure 1

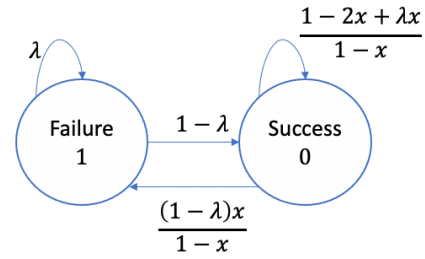


Fig. 1: The Klotz model with dependent Bernoulli trials [17].

need to lie between zero and one, so

$$0 \leq x < 1, \quad \max\{0, (2x-1)/x\} \leq \lambda \leq 1 \quad (3)$$

Inequalities (3) define a subset  $\mathcal{R}$  of the unit square (Figure 2a).

Suppose the system succeeds on all  $n$  demands during operational testing. For  $(x, \lambda) \in \mathcal{R}$ , the Klotz likelihood function gives the probability of observing this sequence of successes<sup>4</sup>:

$$L(x, \lambda; n) := (1-x) \left(1 - \frac{(1-\lambda)x}{1-x}\right)^{n-1} \quad (4)$$

Different values of  $(x, \lambda)$  can alter the dependence structure and functional form of the Klotz likelihood: **i)**  $x = \lambda$  is the special case of independent testing outcomes, with  $L(x, x; n) = (1-x)^n$ ; **ii)** when  $\lambda > x$ , successes (and failures) tend to cluster during testing – i.e. positive dependence. In the extreme,  $\lambda = 1$  and we have  $L(x, 1; n) = (1-x)$ ; **iii)** lastly, when  $\lambda < x$ , successes and failures tend to alternate more often – i.e. negative dependence. In the extreme,  $x = \frac{1}{2-\lambda}$  and we have  $L(\frac{1}{2-\lambda}, \lambda; n) = 0$ .

The assessor’s uncertainty about  $x$  and  $\lambda$  is captured by a joint prior distribution of  $(X, \Lambda)$  over  $\mathcal{R}$ . The assessor’s posterior confidence (1) is:

$$P(X \leq b \mid n \text{ demands without failure}) = \frac{\mathbb{E}[L(X, \Lambda; n) \mathbf{1}_{X \in [0, b]}]}{\mathbb{E}[L(X, \Lambda; n)]} \quad (5)$$

where the indicator function  $\mathbf{1}_{x \in S}$  equals 1 if  $x \in S$ , and is 0 otherwise.

4. We define  $L(1, 1; n) := \lim_{(x, \lambda) \rightarrow (1, 1)} L(x, \lambda; n) = 0$ , where the limiting process involves only  $(x, \lambda)$  values in  $\mathcal{R}$ .

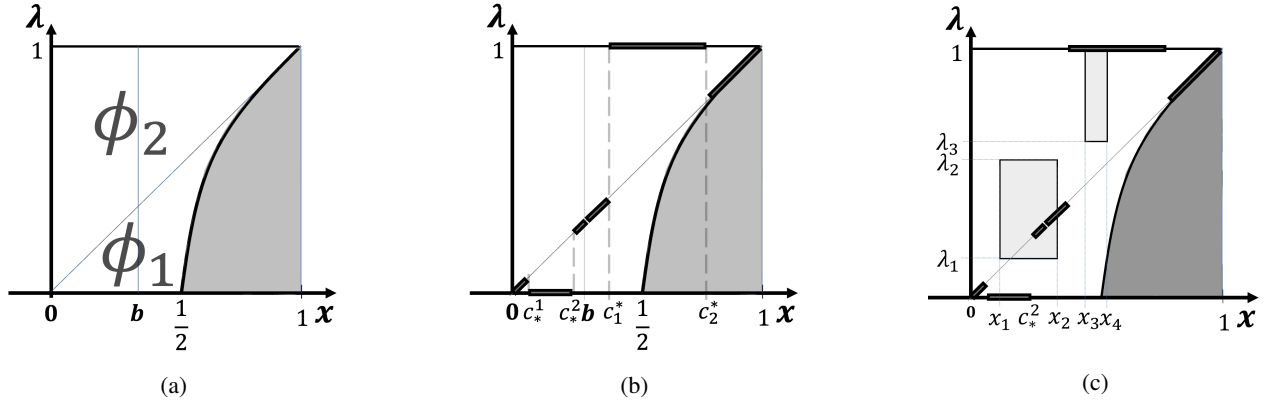


Fig. 2: **a)** A depiction of the region  $\mathcal{R}$  defined by inequalities (3), and the subsets of  $\mathcal{R}$  defined by PKs 2 and 3. All prior distributions of  $(X, \Lambda)$  have domain  $\mathcal{R}$ ; **b)** A prior distribution that gives the theorem’s infimum – depicted “from above”, looking down on its domain  $\mathcal{R}$ . The *pdfs*  $c_*^1, c_*^2, c_1^*, c_2^*$  satisfy the theorem’s constraints; **c)** Probabilities of the shaded rectangular regions are given by integrals of  $f(x)$  from PK1.

In Section 5, via CBI (i.e. constrained nonlinear optimisation of (5) in the vein of (2)), we derive conservative posterior confidence bounds on *pdf*.

## 5 CONSERVATIVE CONFIDENCE (PFD) BOUNDS

We now present conservative posterior confidence bounds on *pdf* for the skeptical assessor. We begin by formalising the assessor’s beliefs as a collection of 4 constraints called “prior knowledge”. These PKs only weakly specify the joint prior distribution of  $(X, \Lambda)$ . Firstly, the assessor’s prior distribution of  $X$  is continuous.

**Prior Knowledge 1** (continuous prior distribution of  $X$ ). A density function  $f(x)$  gives the assessor’s prior confidence in *pdf* bound  $u$  – i.e.  $P(X \leq u) = \int_0^u f(x) dx$  for all  $u \in [0, 1]$ .

Secondly, the assessor does not rule out negatively or positively correlated testing outcomes (see Figure 2a).

**Prior Knowledge 2** (confidence in negative correlation).  $\phi_1 \times 100\%$  confident that the outcomes of successive tests are negatively correlated, i.e.  $P(\Lambda < X) = \phi_1$ .

**Prior Knowledge 3** (confidence in positive correlation).  $\phi_2 \times 100\%$  confident that the outcomes of successive tests are positively correlated, i.e.  $P(\Lambda > X) = \phi_2$ .

Thirdly, the assessor is relatively confident that the i.i.d. assumption holds, and bound  $b$  is satisfied, but not so confident as to make testing unnecessary. A straightforward extension of the theorem presented here accounts for less confidence.

**Prior Knowledge 4** (confidence in bound and independence). For a target system *pdf*  $b$ ,  $\phi_1 \leq P(X \leq b) \leq 1 - \phi_2$

An assessor with these 4 beliefs has conservative posterior confidence bounds given by this theorem (see appendix B):

**Theorem.** Let  $\mathcal{D}$  be the set of all prior distributions over  $\mathcal{R}$  and assume  $0 < b < 1/2$ . Using (5), the optimisation problem

$$\inf_{\mathcal{D}} P(X \leq b \mid n \text{ demands without failure})$$

s.t. PK1, PK2, PK3, PK4

is solved by the prior in Figure 2b, since the infimum equals the value of  $P(X \leq b \mid n \text{ demands without failure})$  for this prior. The infimum takes the form  $\frac{1}{1+Q}$ , where  $Q$  is

$$\frac{\int_b^1 \left( (1-x)^n \mathbf{1}_{x \in (b, c_1^*) \cup (c_2^*, 1)} + (1-x) \mathbf{1}_{x \in (c_1^*, c_2^*)} \right) f(x) dx}{\int_0^b \left( (1-x)^n \mathbf{1}_{x \in (0, c_1^*) \cup (c_2^*, b)} + \frac{(1-2x)^{n-1}}{(1-x)^{n-2}} \mathbf{1}_{x \in (c_1^*, c_2^*)} \right) f(x) dx} \quad (6)$$

and the *pdfs*  $c_*^1, c_*^2, c_1^*, c_2^*$  are the unique values of  $r, s, v, w$ , respectively, that solve

$$\begin{aligned} & \arg \min_{0 \leq r < s \leq b} |g_l(r) - g_l(s)|, \quad \arg \min_{b \leq v < w \leq 1} |g_u(v) - g_u(w)| \\ \text{s.t.} \quad & g_l(0) \leq g_l(r), \quad g_l(b) \leq g_l(s), \\ & g_u(b) \leq g_u(v), \quad g_u(1) \leq g_u(w), \\ & \int_r^s f(x) dx = \phi_1, \quad \int_v^w f(x) dx = \phi_2, \\ & 0 \leq r < s \leq b \leq v < w \leq 1 \end{aligned}$$

for  $g_l : [0, 1/2] \rightarrow [0, 1]$  and  $g_u : [0, 1] \rightarrow [0, 1]$  defined as

$$\begin{aligned} g_l(x) &= (L(x, x; n) - L(x, 0; n)) \mathbf{1}_{x \in [0, \frac{1}{2}]} \\ g_u(x) &= (L(x, 1; n) - L(x, x; n)) \mathbf{1}_{x \in [0, 1]} \end{aligned} \quad (7)$$

Numerical estimates for  $c_*^1, c_*^2, c_1^*, c_2^*$  may be computed using “root-finding” algorithms such as that in Appendix D.

If the assessor has no doubts about the i.i.d. assumption, this CBI result is given by traditional Bayesian inference. That is,

**Corollary.** Let  $Q$  be given by (6) in the theorem, when  $\phi_1 = \phi_2 = 0$ . Then the infimum is given by traditional Bayesian inference, using the prior density  $f(x)$ . That is,

$$\inf_{\mathcal{D}} P(X \leq b \mid n \text{ demands without failure}) = \frac{\int_0^b (1-x)^n f(x) dx}{\int_0^1 (1-x)^n f(x) dx} \quad (8)$$

However, if the assessor begins operational testing with beliefs captured by PKs 1, 2, 3 and 4 – i.e. beliefs that only partially specify a joint prior distribution of  $(X, \Lambda)$  – then the CBI theorem identifies a joint prior that is consistent with these beliefs, is unique up to zero probability events, and that gives the smallest posterior confidence in the *pdf* upper-bound  $b$  (see Figure 2b).

This CBI prior assigns probability density only to the indicated thick black line segments in  $\mathcal{R}$ ; the rest of  $\mathcal{R}$  has zero probability<sup>5</sup> (Appendix B). For example, prior to testing, the conservative

5. Expressing this formally involves projections from  $\mathcal{R}$  to the interval  $[0, 1]$ .

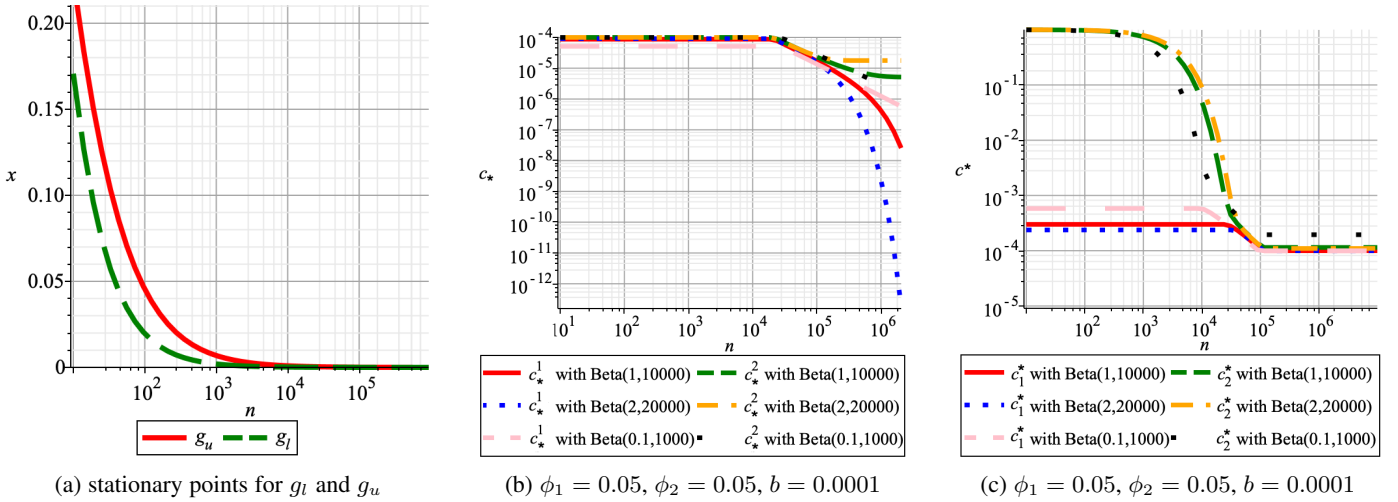


Fig. 3: The stationary points of  $g_l$ ,  $g_u$  and the *pfd*s  $c_*^1$ ,  $c_*^2$ ,  $c_1^*$ ,  $c_2^*$  are all monotonically decreasing functions of  $n$ .

assessor expects negatively, or positively, dependent test outcomes with probabilities  $\phi_1 = P(c_*^1 \leq X \leq c_*^2, \Lambda = 0) = \int_{c_*^1}^{c_*^2} f(x) dx$  and  $\phi_2 = P(c_1^* \leq X \leq c_2^*, \Lambda = 1) = \int_{c_1^*}^{c_2^*} f(x) dx$ , respectively. Or consider the probabilities of the 2 shaded rectangular regions in Figure 2c. The non-zero contributions to these probabilities come from the thick black line segments that intersect these regions; i.e.  $P(x_1 \leq X \leq x_2, \lambda_1 \leq \Lambda \leq \lambda_2) = \int_{c_1^*}^{c_2^*} f(x) d(x)$  and  $P(x_3 \leq X \leq x_4, \lambda_3 \leq \Lambda \leq 1) = \int_{c_1^*}^{c_2^*} f(x) d(x)$ .

## 5.1 Conservative Beliefs for Failure-free Testing

The CBI prior Figure 2b encodes conservative beliefs. In particular, the *pfd* values  $c_*^1$ ,  $c_*^2$ ,  $c_1^*$ ,  $c_2^*$  indicate where (in  $\mathcal{R}$ ) doubts in the i.i.d. assumption should be placed for conservative posterior confidence. The  $\mathcal{R}$  locations between  $(c_*^1, 0)$  and  $(c_*^2, 0)$  along the  $\lambda = 0$  edge represent statistical dependence that is *unlikely* to produce failure-free testing *if* the unknown *pfd* actually satisfies the bound. And if the *pfd* does not satisfy the bound, the locations between  $(c_1^*, 1)$  and  $(c_2^*, 1)$  along  $\lambda = 1$  represent dependence *likely* to produce failure-free testing. Note that the assessor has been unable to rule out these extreme beliefs prior to testing, since these beliefs are consistent with the PKs.

These beliefs depend on the number  $n$  of required failure-free tests; the “ $c$ ”s become smaller if  $n$  becomes bigger. Because, the “ $c$ ”s are defined with respect to the stationary points of the theorem’s “ $g$ ” functions – themselves dependent on  $n$ . Figure 3a plots how the  $x$  values for these stationary points tend to zero as  $n$  increases. Figure 3b shows the consequence of this –  $c_*^1$  decreases towards zero and  $c_*^2$  decreases towards a non-zero value dependent on  $f(x)$ . Figure 3c tells a similar story –  $c_1^*$  decreases towards the bound  $10^{-4}$ , and  $c_2^*$  towards a non-zero value dependent on  $f(x)$ . This asymptotic behaviour is reasonable: the greater the amount of failure-free operation required during testing, the smaller the *pfd* is expected to be for a system that performs this well (even when one is doubtful of the tests being i.i.d.).

Informally, beliefs in i.i.d. tests “lie on the boundary” of the set of conservative beliefs that express doubts about i.i.d. tests, and confidence from beliefs in i.i.d. tests (8) is well-approximated by conservative confidence from the theorem. Indeed, by the *dominated convergence theorem* [46], (6) tends to  $\frac{\int_0^1 (1-x)^n f(x) dx}{\int_0^1 (1-x)^n f(x) dx}$  as  $\phi_1$  and  $\phi_2$  tend to 0 (since  $c_*^1$ ,  $c_1^*$  tend to  $c_*^2$ ,  $c_2^*$  respectively).

## 6 RESULTS: ASSESSMENT USING THE THEOREM

### 6.1 Practical Context and Guidance

The theorem gives conservative confidence in a *pfd* bound  $b$ , when an on-demand system is subjected to black-box operational testing. A bound such as  $b = 10^{-4}$ ; a target *pfd* used in the assessment of the Sizewell-B nuclear power plant safety protection system in the United Kingdom<sup>6</sup> [47], [48]. To gain 99% confidence in this bound – using the i.i.d. assumption under a classical statistical inference approach – requires between  $10^4$  and  $10^5$  test demands [48], [49]. In this section, we will use similar orders of magnitude of test demands to illustrate the theorem’s use.

In particular, the theorem can be used during acceptance testing to check the robustness of confidence (8) to doubts about the i.i.d. assumption; we illustrate how to do this in the rest of this section. When applying statistical techniques like the theorem, one may follow the guidance from Littlewood and Strigini [50] and Lyu *et al.* [38] on performing statistical testing. See Parnas *et al.* [49] for additional discussion on evaluating safety-critical software, including random test-case (e.g. demand) selection. For nuclear safety applications in particular, [7], [51] give guidance on reliability assessment using statistical techniques.

More generally, the theorem can be applied in any black-box testing phase where failure-free operation can be used to gain confidence (i.e. (8)) in the software. Such testing would typically involve subjecting (some part of) the software to a large number of randomly generated demands in a simulated environment. For example, in accordance with integrity level 4 (see IEEE 1012:2016 [52]), one may apply the theorem during component, or integration, testing phases for safety-critical software. For nuclear safety-critical software testing phases, see also IEC 60880 [53].

### 6.2 Examples: Prior Beliefs and Confidence in a Bound

The analyses in the rest of this section show how: **i)** confidence based on the i.i.d. assumption can be very optimistic as failure-free tests accumulate; **ii)** some forms of doubt about the i.i.d. assumption significantly impact confidence, while other forms do not; **iii)** surprisingly, failure-free testing can eventually undermine confidence in the system satisfying the bound.

6. This was the target *pfd* for a hardwired secondary safety subsystem; a software-based primary safety subsystem had a more modest  $10^{-3}$  target *pfd*.



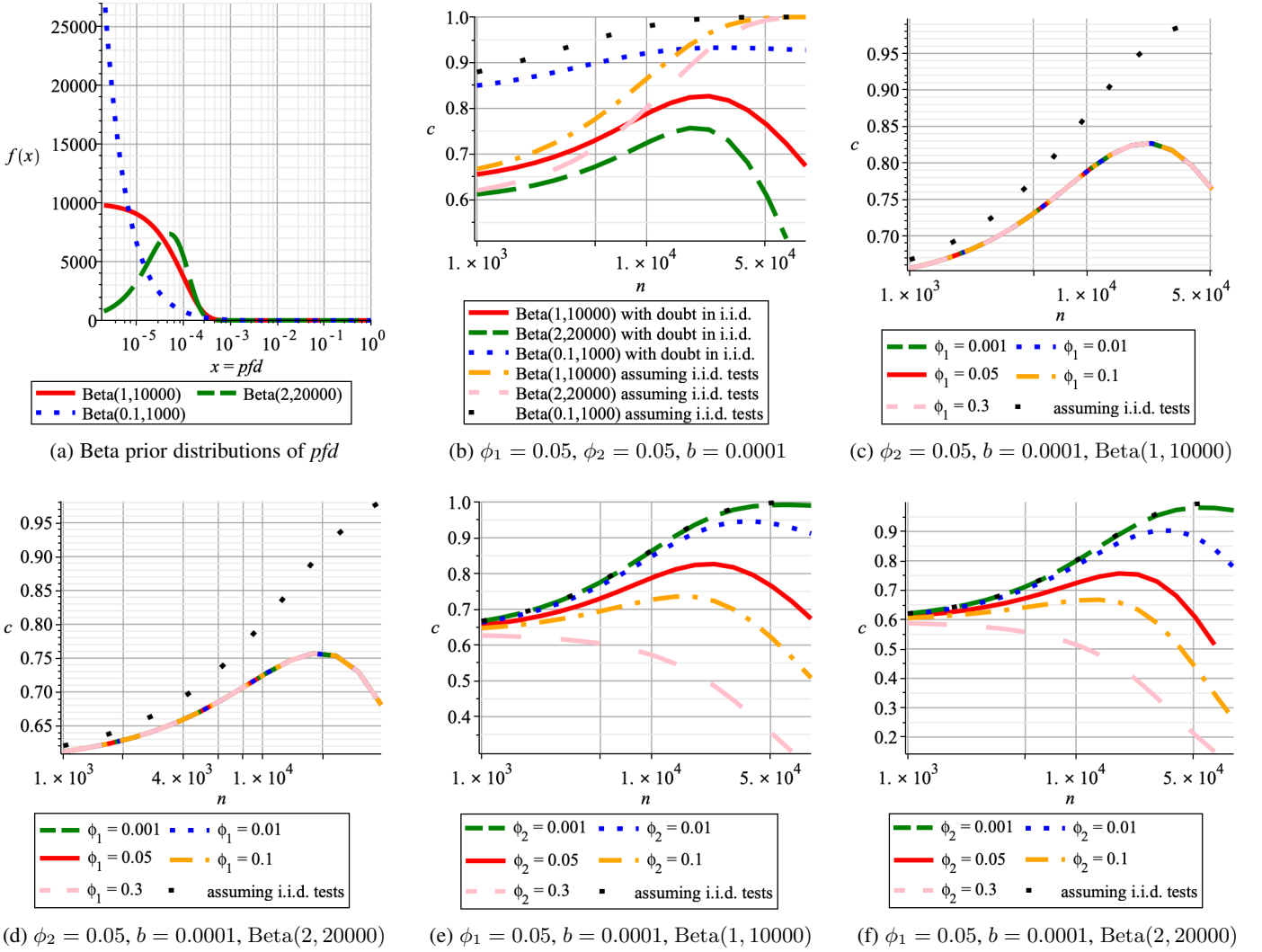


Fig. 4: Sensitivity analyses showing which forms of PK have the biggest impact on (conservative) posterior confidence  $c$ .

TABLE 1: A summary of 3 Beta prior distributions of  $pdf$

$\alpha$	$\beta$	$\mathbb{E}[X] = \frac{\alpha}{\alpha+\beta}$	$P(X \leq 10^{-4})$
2	20000	0.0001	0.6
1	10000	0.0001	0.63
0.1	1000	0.0001	0.83

The PK1 density  $f(x)$  can be from *any family* of continuous distributions over the interval  $[0, 1]$ . Beta densities are often used in practice [7], [19], [21]. Consider 3 alternative  $f(x)$ ; the Beta distributions in Figure 4a with parameters/properties in Table. 1. Let these represent prior beliefs of 3 assessors that differ in how confident they are that the system satisfies the bound.

Suppose the assessors are a little skeptical of the tests being i.i.d. (e.g.  $\phi_1 = \phi_2 = 0.05$ ). Figure 4b shows how posterior confidence  $c$  evolves as operational testing evidence mounts. For each Beta prior, the posterior confidence (8) under an i.i.d. assumption is plotted against the CBI posterior confidence  $\frac{1}{1+Q}$  for our skeptic. In all cases, confidence from assuming independence is initially comparable to conservative confidence – the relevant pair of curves for each Beta prior almost overlap initially.

However, as more failure-free testing is observed – i.e. as  $n$

grows – “i.i.d.”-based confidence grows and tends towards certainty. While conservative confidence grows more slowly, reaches a maximum, and then tends towards zero. Appendix C proves this zero limiting behaviour will occur whenever an assessor allows for the possibility of positively correlated tests (i.e.  $\phi_2 > 0$ ). In the limit of large  $n$ , “i.i.d.”-based confidence can be the *most optimistic* confidence *can* be, while still remaining consistent with an assessor’s informed views/beliefs about the unknown  $pdf$ .

If the evidence available to an assessor before testing justifies being very confident in the bound, then the initial closeness between the “i.i.d.”-based posterior confidence and the confidence given by CBI can continue for longer before ultimately diverging. In Figure 4b the pair of curves for the Beta(0.1,1000) prior – i.e. for the very confident assessor – stay closer together for longer, compared to the curves for the least confident assessor with prior Beta(2,20000). The greater the prior confidence in the bound, the greater the posterior confidence when testing begins (i.e. small  $n$ ).

### 6.3 Sensitivity Analysis: How Failure-free Tests and Skepticism about i.i.d. Tests Impact Confidence

The nature of an assessor’s skepticism determines whether their conservative posterior confidence ultimately grows or shrinks during operational testing. At one extreme, some forms of doubt have

no noticeable impact on confidence when no failures are observed during testing. The possibility (however likely) of negatively correlated tests has no apparent effect on conservative confidence. For example, an assessor might intentionally seek to “stress” the software during testing, by randomly including a disproportionate number of test demands that are thought will likely cause the software’s failure. If stressful demands are adequately interspersed with significantly less stressful ones, one might expect the testing outcomes (i.e. the software’s successes/failures) to exhibit some negative dependence (i.e. non-zero  $\phi_1$ ) – so a failure is quickly followed by successes, then another failure relatively soon afterwards, etc. However, when no failures occur, Figures 4c, 4d show that the “rise and fall” of CBI confidence in Figure 4b is unaffected by varying one’s prior confidence in negative dependence. Intuitively, the more successes occur, the less likely these are from a system undergoing negatively dependent tests.

At the other extreme are positively correlated tests. Figures 4e, 4f both show that the smaller  $\phi_2$  is, the closer conservative confidence gets to the confidence under the i.i.d. assumption. When  $\phi_2 = 0$ , conservative confidence grows to certainty as the number of successes grows (Appendix C). While the larger  $\phi_2$  is, the more conservative confidence becomes. Here, confidence in positive correlations (i.e. large  $\phi_2$ ) may be due to pessimistic reasons for the failure-free tests – i.e. “success clustering” can occur even if the software is unreliable. The tests could be unrepresentatively “easy” for the software to correctly respond to, or the test oracle is incorrect so failures go undetected [54], [55].

So, failure-free testing can undermine one’s confidence in a system’s *pdf*. Such conservatism is not unique to CBI – even with classical inference, confidence bounds can be quite conservative initially, becoming optimistic (compared to the CBI bounds) after many tests. Indeed, using the Klotz likelihood (4) when  $\lambda = 1$ , the probability of succeeding on all  $n$  tests (despite the *pdf* being worse than  $10^{-4}$ ) is at most  $(1 - 10^{-4})^n = 0.9999^n$ . That is, a system with a *pdf* worse than  $10^{-4}$  may be almost certain to succeed on all tests *if* the tests are strongly positively correlated. This bleak result holds *for all n*; so, increasing the number of failure-free tests does not increase one’s confidence in the system.

## 7 DISCUSSION

### 7.1 Skepticism about Model Assumptions

Software reliability assessments should be conservative: to wit, only when test results stand up to the most critical scrutiny can confidence be justifiably placed in the system satisfying a *pdf* bound. Conservative assessments require a skeptical assessor. In Bayesian terms, our assessor holds conservative beliefs about what the evidence implies for a system’s reliability, and about the validity of statistical modelling assumptions.

This paper illustrates a general, incremental approach to dealing with doubts about *any* statistical model assumptions: we offer a demonstrably conservative form of Draper’s ideas [34]. For a model property one is doubtful of, one can check the sensitivity of claims based on the model by using a slightly more general model (that has the original model as a special case and weakens the property in question) for inference. “Slight” model generalisations keep models from becoming unnecessarily complex, ensure generalisations cover all scenarios covered by the original models, and minimize eliciting increasingly complex prior distributions.

This incremental approach is a “win-win”. If the i.i.d. assumption is *not* too optimistic, “i.i.d.”-based confidence doesn’t depart

significantly from the CBI theorem’s conservative confidence based on the Klotz model. If the i.i.d. assumption is too optimistic, sensitivity analysis using the theorem can reveal this – in such circumstances, caution is warranted when relying on “i.i.d.”-based reliability claims. And if, in turn, one has doubts about the Klotz model, then a generalisation of the Klotz model can be used to check if “Klotz model”-based confidence is sensitive to doubts.

### 7.2 Limitations and Future Work

Let us highlight some Klotz model shortcomings. It does not distinguish between different success/failure types. Future work might consider using the models of [16], [56], [57] with CBI, to check the robustness of “Klotz model”-based claims. These models account for the cumulative impact of benign failures.

The Klotz model uses the relative sizes of  $x$  and  $\lambda$  to characterise *all* pairs of successive Bernoulli trials as being identically positively, negatively, or zero correlated. Consequently, the Klotz model is unable to express non-stationary dependence, such as may be due to software updates that remove, or inadvertently add, faults to the software. The model cannot capture dependence across time either; such as periodic correlations over relatively short, or relatively long, runs of demands. Such periodicity can arise if demands that cause software failure are more likely at certain times when (or certain locations where) operating conditions tend to be more stressful (e.g. for software ensuring aircraft safety, unfavourable weather along a flight path may be more likely at certain times of the year, or more likely along certain flight paths). CBI models with time-dependent correlations are worth exploring.

On the use of *pdfs* we make the following comment. When the failure process is stationary, *pdfs* make sense. The probability of the system failing on the  $n$ -th demand is the same for all  $n$ . But for time-dependent failure probabilities, there are more suitable dependability measures – such as the probability of future failure-free operation. Strigini [57] makes a related point in a classical inference context. Even in the present context of posterior confidence bounds on *pdf*, it’s worth investigating whether other dependability measures are more/less robust to i.i.d. doubts.

We have formalised (via PK1) and illustrated how *any* Bayesian reliability assessment using a continuous prior,  $f(x)$ , can conservatively incorporate doubts about “i.i.d.”. However, when eliciting  $f(x)$  proves too challenging, future work could extend the theorem to work with partially specified  $f(x)$ . For example, prior to testing, one might justifiably have some confidence in the software containing no faults [58]. In effect,  $f(x)$  becomes discontinuous, with a non-zero probability of the *pdf* being zero. For such scenarios, preliminary results suggest significantly better agreement between “i.i.d.”-based and “Klotz model”-based posterior confidence, even with a very large number of successful tests.

In general, the “elicitation challenge” remains an open problem for Bayesian approaches. In light of this, best practice approaches should be followed when eliciting PKs [59], [60], while sensitivity analyses (as illustrated in Section 6) is crucial and practical for checking the robustness of confidence to PK changes.

This work has not considered model selection or validation. One might envisage applying CBI to conservatively gain confidence in the i.i.d. assumption, or using Bayes factors and CBI to conservatively determine which modelling assumptions lead to more trustworthy predictions about future system reliability.

The theorem could be extended to account for failures during testing, e.g. when assessing machine learning applications, or

extended to support conservative claims for software modules in a fault-tolerant configuration (e.g. extending Singh *et al.* [20]).

## 8 CONCLUSIONS

When assessing software using operational testing it is natural to ask, “is it appropriate to assume the software’s failures and successes arise in an i.i.d. manner?”. For many practical scenarios there are well-known reasons to doubt the i.i.d. assumption. A few statistical models which weaken the assumption have been proposed for use in reliability assessments [13], [14], [16], [56]. However, none of these proposals allow an assessor to remain unsure about whether i.i.d. holds or not, nor allow the assessor to see the impact of their uncertainty on their confidence in a *pdf* bound. This, despite it often being the case that an assessor may have good reason to believe in i.i.d., but not enough reason to be certain that it holds. Furthermore, these proposals do not directly support *conservative* reliability claims.

Using conservative Bayesian techniques – in particular, CBI – we show how doubts about i.i.d. can be formally included in software reliability assessments (see Sections 5 and 6). In this way, we obviate the need imposed by the previous proposals – the need for assessors to either assume/conclude that the software test outcomes are i.i.d., or assume/conclude that they aren’t. Instead, our method allows a skeptical assessor’s confidence in a target *pdf*, and their confidence in (and doubts about) i.i.d., to grow or shrink in response to seeing the software operate without failure.

Moreover, CBI’s conservative confidence bounds continuously invite our assessor to be skeptical, and to question whether seemingly favourable reliability evidence from testing does, in fact, corroborate actual reliability. For example, while failure-free operation is generally an indicator of a desirable *pdf*, our results highlight why this might be (at best) a sanguine view in some situations. There may be undesirable reasons for why failure-free operation is occurring; reasons that ultimately undermine one’s confidence in the *pdf* being sufficiently small (see section 6). CBI, by weakening the i.i.d. assumption and producing conservative confidence bounds, can call into question the representativeness of failure-free operation (as indicating a reliable system). When this happens, it’s incumbent on the assessor to rule out potential problems during testing that “masquerade” as failure-free operation, and to incorporate these efforts into any further use of CBI.

## ACKNOWLEDGMENTS

We are grateful to the anonymous reviewers whose comments were very helpful in improving the presentation. We are also grateful to Bev Littlewood for giving very helpful feedback on an initial draft of this paper. This work was partly funded by the European Union’s Horizon 2020 Research and Innovation Programme under grant agreement No 956123, and by the UK EPSRC through the End-to-End Conceptual Guarding of Neural Architectures [EP/T026995/1]. Xingyu Zhao’s contribution is partially supported through Fellowships at the Assuring Autonomy International Programme.

## REFERENCES

- [1] Thomas E. Wierman, Scott T. Beck, Michael B. Calley, Steven A. Eide, Cindy D. Gentillon, and William E. Kohn, “Reliability study: combustion engineering reactor protection system, 1984–1998,” Idaho National Engineering and Environmental Laboratory, U.S. Nuclear Regulatory Commission Washington, DC, Tech. Rep. NUREG/CR-5500 Vol.10, 2001. [Online]. Available: <https://www.nrc.gov/docs/ML0220/ML022050271.pdf>
- [2] C. Bunea, T. Charitos, R. M. Cooke, and G. Becker, “Two-stage Bayesian models—application to ZEDB project,” *Reliability Engineering & System Safety*, vol. 90, no. 2, pp. 123–130, 2005.
- [3] K. Pörn, “The two-stage Bayesian method used for the T-Book application,” *Reliability Engineering & System Safety*, vol. 51, no. 2, pp. 169–179, 1996.
- [4] N. E. Fenton and M. Neil, “Software metrics: successes, failures and new directions,” *Journal of Systems and Software*, vol. 47, no. 2, pp. 149–157, 1999.
- [5] R. A. Thayer, M. Lipow, and E. C. Nelson, *Software Reliability*. North-Holland, 1978.
- [6] J. W. Duran and S. C. Ntafos, “An evaluation of random testing,” *IEEE Transactions on Software Engineering*, vol. SE-10, no. 4, pp. 438–444, 1984.
- [7] C. Atwood, J. LaChance, H. Martz, D. Anderson, M. Englehardt, D. Whitehead, and T. Wheeler, “Handbook of parameter estimation for probabilistic risk assessment,” U.S. Nuclear Regulatory Commission, Washington, DC, Report NUREG/CR-6823, 2003.
- [8] P. E. Ammann and J. C. Knight, “Data diversity: an approach to software fault tolerance,” *IEEE transactions on computers*, vol. 37, no. 4, pp. 418–425, 1988.
- [9] P. Bishop, “The variation of software survival time for different operational input profiles (or why you can wait a long time for a big bug to fail),” in *FTCS-23 The Twenty-Third International Symposium on Fault-Tolerant Computing*. IEEE, 1993, pp. 98–107.
- [10] R. Huang, W. Sun, Y. Xu, H. Chen, D. Towey, and X. Xia, “A survey on adaptive random testing,” *IEEE Transactions on Software Engineering*, vol. 47, no. 10, pp. 2052–2083, 2021.
- [11] A. Csenki, “Reliability analysis of recovery blocks with nested clusters of failure points,” *IEEE transactions on reliability*, vol. 42, no. 1, pp. 34–43, 1993.
- [12] L. A. Tomek, J. K. Muppala, and K. S. Trivedi, “Modeling correlation in software recovery blocks,” *IEEE Transactions on Software Engineering*, vol. 19, no. 11, pp. 1071–1086, 1993.
- [13] S. Chen and S. Mills, “A binary Markov process model for random testing,” *IEEE Transactions on Software Engineering*, vol. 22, no. 3, pp. 218–223, 1996.
- [14] K. Goseva-Popstojanova and K. S. Trivedi, “Failure correlation in software reliability models,” *IEEE Transactions on Reliability*, vol. 49, no. 1, pp. 37–48, 2000.
- [15] Y. Dai, M. Xie, and K. Poh, “Modeling and analysis of correlated software failures of multiple types,” *IEEE Transactions on Reliability*, vol. 54, no. 1, pp. 100–106, 2005.
- [16] A. Bondavalli, S. Chiaradonna, F. Di Giandomenico, and L. Strigini, “Dependability models for iterative software considering correlation between successive inputs,” in *Proc. of IEEE Int. Computer Performance and Dependability Symposium*. Erlangen, Germany: IEEE, 1995, pp. 13–21.
- [17] J. Klotz, “Statistical inference in Bernoulli trials with dependence,” *The Annals of Statistics*, vol. 1, no. 2, pp. 373–379, 1973.
- [18] K. Pörn, “The two-stage Bayesian method used for the T-Book application,” *Reliability Engineering & System Safety*, vol. 51, no. 2, pp. 169–179, 1996.
- [19] K. W. Miller, L. J. Morell, R. E. Noonan, S. K. Park, D. M. Nicol, B. W. Murrill, and M. Voas, “Estimating the probability of failure when testing reveals no failures,” *IEEE Transactions on Software Engineering*, vol. 18, no. 1, pp. 33–43, 1992.
- [20] H. Singh, V. Cortellessa, B. Cukic, E. Gunel, and V. Bharadwaj, “A Bayesian approach to reliability prediction and assessment of component based systems,” in *Proc. 12th Int. Symp. on Software Reliability Engineering*. IEEE Computer Society, 2001, pp. 12–21.
- [21] B. Littlewood, P. Popov, and L. Strigini, “Assessing the reliability of diverse fault-tolerant software-based systems,” *Safety Science*, vol. 40, no. 9, pp. 781–796, 2002.
- [22] P. Popov, “Bayesian reliability assessment of legacy safety-critical systems upgraded with fault-tolerant off-the-shelf software,” *Reliability engineering & system safety*, vol. 117, pp. 98–113, 2013.
- [23] P. Bishop, R. Bloomfield, B. Littlewood, A. Povyakalo, and D. Wright, “Toward a formalism for conservative claims about the dependability of software-based systems,” *IEEE Transactions on Software Engineering*, vol. 37, no. 5, pp. 708–717, 2011.
- [24] L. Strigini and A. Povyakalo, “Software fault-freeness and reliability predictions,” in *Computer Safety, Reliability, and Security*, ser. LNCS, vol. 8153. Springer Berlin Heidelberg, 2013, pp. 106–117.
- [25] K. Salako, “Loss-size and reliability trade-offs amongst diverse redundant binary classifiers,” in *Quantitative Evaluation of Systems*, M. Grib-



- audio, D. N. Jansen, and A. Remke, Eds. Springer International Publishing, 2020, pp. 96–114.
- [26] X. Zhao, V. Robu, D. Flynn, K. Salako, and L. Strigini, “Assessing the safety and reliability of autonomous vehicles from road testing,” in *the 30th Int. Symp. on Software Reliability Engineering*. Berlin, Germany: IEEE, 2019, pp. 13–23.
- [27] X. Zhao, K. Salako, L. Strigini, V. Robu, and D. Flynn, “Assessing safety-critical systems from operational testing: A study on autonomous vehicles,” *Information and Software Technology*, vol. 128, p. 106393, 2020.
- [28] B. Littlewood, K. Salako, L. Strigini, and X. Zhao, “On reliability assessment when a software-based system is replaced by a thought-to-be-better one,” *Reliability Engineering & System Safety*, vol. 197, p. 106752, 2020.
- [29] K. Salako, L. Strigini, and X. Zhao, “Conservative confidence bounds in safety, from generalised claims of improvement & statistical evidence,” in *51st Annual IEEE/IFIP Int. Conf. on Dependable Systems and Networks*, ser. DSN’21. Taipei Taiwan: IEEE/IFIP, 2021, pp. 451–462.
- [30] J. O. Berger, “An overview of robust Bayesian analysis,” *Test*, vol. 3, no. 1, pp. 5–124, 1994.
- [31] —, “Robust Bayesian analysis: Sensitivity to the prior,” *Journal of statistical planning and inference*, vol. 25, pp. 303–328, 1990.
- [32] M. Lavine, “Sensitivity in Bayesian statistics: the prior and the likelihood,” *Journal of the American Statistical Association*, vol. 86, no. 414, pp. 396–399, 1991.
- [33] J. Berger and E. Moreno, “Bayesian robustness in bidimensional models: Prior independence,” *Journal of statistical planning and inference*, vol. 40, no. 2, pp. 161–176, 1994.
- [34] D. Draper, “Assessment and propagation of model uncertainty,” *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 57, no. 1, pp. 45–97, 1995.
- [35] L. R. Pericchi and M. E. Perez, “Posterior robustness with more than one sampling model,” *Journal of Statistical Planning and Inference*, vol. 40, no. 2, pp. 279–294, 1994.
- [36] PD IEC TR 63161:2022, “Assignment of safety integrity requirements: basic rationale,” IEC, Geneva, Switzerland, Standard, 2022.
- [37] M. Rausand, *Reliability of safety-critical systems: theory and applications*. Hoboken, New Jersey: Wiley & sons, 2014.
- [38] M. R. Lyu, Ed., *Handbook of Software Reliability Engineering*. USA: McGraw-Hill, Inc., 1996.
- [39] N. D. Singpurwalla and S. P. Wilson, “Software reliability modeling,” *International Statistical Review / Revue Internationale de Statistique*, vol. 62, no. 3, pp. 289–317, 1994. [Online]. Available: <http://www.jstor.org/stable/1403763>
- [40] D. R. Miller, “Exponential order statistic models of software reliability growth,” *IEEE Transactions on Software Engineering*, vol. SE-12, no. 1, pp. 12–24, 1986.
- [41] B. Bergman and M. Xie, “On Bayesian software reliability modelling,” *Journal of Statistical Planning and Inference*, vol. 29, no. 1, pp. 33–41, 1991. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/037837589290119D>
- [42] B. Littlewood and L. Strigini, “Validation of ultra-high dependability for software-based systems,” *Comm. of the ACM*, vol. 36, pp. 69–80, 1993.
- [43] B. Littlewood and D. Wright, “Some conservative stopping rules for the operational testing of safety critical software,” *IEEE Transactions on Software Engineering*, vol. 23, no. 11, pp. 673–683, 1997.
- [44] J. D. Musa, “Operational profiles in software-reliability engineering,” *IEEE Software*, vol. 10, no. 2, pp. 14–32, 1993.
- [45] IEC 61508:2010, “Functional Safety of Electrical/ Electronic/Programmable Electronic Safety Related Systems,” IEC, Geneva, Switzerland, Standard, 2009.
- [46] R. L. Schilling, *Measures, Integrals and Martingales*. Cambridge University Press, 2005.
- [47] D. M. Hunns and N. Wainwright, “Software-based protection for Sizewell B: the regulator’s perspective,” in *Int. Conf. on Electrical and Control Aspects of the Sizewell B Power Station*, 1992, pp. 198–203.
- [48] NuSAC study group on the safety of operational computer systems, *The Use of Computers in Safety-critical Applications*. Health and safety Commission, London, UK, 1998.
- [49] D. L. Parnas, A. J. van Schouwen, and S. P. Kwan, “Evaluation of safety-critical software,” *Commun. ACM*, vol. 33, no. 6, p. 636–648, 1990.
- [50] L. Strigini and B. Littlewood, “Guidelines for statistical testing.” City University London, Project Report PASCON/WO6-CCN2/TN12, 1997.
- [51] The International Regulator Task Force on Safety Critical Software, “Licensing of safety critical software for nuclear reactors. common position of international nuclear regulators and authorised technical support organisations,” 2022. [Online]. Available: <https://www.onr.org.uk/software.pdf>
- [52] IEEE 1012:2016, “IEEE standard for system, software, and hardware verification and validation,” IEEE, New York, U.S.A, Standard, 2016.
- [53] IEC 60880, “Nuclear power plants. Instrumentation and control systems important to safety. Software aspects for computer-based systems performing category A functions.” IEC, Geneva, Switzerland, Standard, 2006.
- [54] B. Littlewood and D. Wright, “The use of multilegged arguments to increase confidence in safety claims for software-based systems: A study based on a BBN analysis of an idealized example,” *IEEE Transactions on Software Engineering*, vol. 33, no. 5, pp. 347–365, 2007.
- [55] E. T. Barr, M. Harman, P. McMinn, M. Shahbaz, and S. Yoo, “The oracle problem in software testing: A survey,” *IEEE Transactions on Software Engineering*, vol. 41, no. 5, pp. 507–525, 2015.
- [56] A. Bondavalli, S. Chiaradonna, F. Di Giandomenico, and S. La Torre, “Modelling the effects of input correlation in iterative software,” *Reliability Engineering & System Safety*, vol. 57, no. 3, pp. 189–202, 1997.
- [57] L. Strigini, “On testing process control software for reliability assessment: the effects of correlation between successive failures,” *Software Testing, Verification and Reliability*, vol. 6, no. 1, pp. 33–48, 1996.
- [58] B. Littlewood and J. Rushby, “Reasoning about the reliability of diverse two-channel systems in which one channel is ‘possibly perfect’,” *IEEE Tran. on Software Engineering*, vol. 38, no. 5, pp. 1178–1194, 2012.
- [59] PRA Working Group, “A review of NRC staff uses of probabilistic risk assessment,” U. S. Nuclear Regulatory Commission, Tech. Rep. NUREG-1489, 1994. [Online]. Available: <https://www.nrc.gov/docs/ML0635/ML063540593.pdf>
- [60] A. O’Hagan, C. Buck, A. Daneshkhah, J. Eiser, P. Garthwaite, D. Jenkinson, J. Oakley, and T. Rakow, *Uncertain Judgements: Eliciting Experts’ Probabilities*, ser. Statistics in Practice. Wiley, 2006.
- [61] W. Rudin, *Principles of Mathematical Analysis*, 3rd ed., ser. International series in pure and applied mathematics. McGraw-Hill, 1976.
- [62] V. Bryant, *Metric Spaces: Iteration and Application*. Cambridge University Press, 1985.
- [63] C. D. Aliprantis and K. C. Border, *Infinite Dimensional Analysis: a Hitchhiker’s Guide*, 2nd ed. Berlin: Springer, 1999.

## APPENDIX A TRANSITION PROBABILITIES IN THE KLOTZ MODEL

*1st-order stationarity* requires that the probability of being in a given state after  $n$  trials is the same for all  $n$ . In particular, the probability of being in a successful state after two trials is the same as the probability after one trial, i.e.  $1 - x$ . So, upon writing the shorthand  $p = P(T_2 = 0 \mid T_1 = 0)$ , we have  $1 - x = x(1 - \lambda) + (1 - x)p$ . Solving for  $p$  gives  $P(T_2 = 0 \mid T_1 = 0) = p = 1 - \frac{(1-\lambda)x}{1-x}$  and  $P(T_2 = 1 \mid T_1 = 0) = 1 - p = \frac{(1-\lambda)x}{1-x}$ , for  $0 \leq x < 1$ .

## APPENDIX B PROOF OF THE THEOREM

*Proof.* Choose any  $F \in \mathcal{D}$  that satisfies the constraints of the optimisation and denote the Klotz likelihood (4) as  $L$ . The objective function (5) in the theorem, computed using  $F$ , is  $\frac{\int_{\{x \leq b\} \cap \mathcal{R}} L dF}{\int_{\mathcal{R}} L dF} = \left(1 + \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF}\right)^{-1}$ . Consequently, we focus on the equivalent optimisation (subject to the same constraints)

$$\sup_{\mathcal{D}} \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF} \quad (9)$$

From  $F$ , one can construct a sequence of priors  $\{F_k^*\}$  (for  $k = 1, 2, \dots$ ) that: **i)** all give larger values than  $F$  for the objective function in (9); and **ii)** give objective function values that converge to the objective function value given by some  $F^* \in \mathcal{D}$ . The construction is as follows. Consider the sequence  $\{\mathcal{P}_k\}$  of partitions of the interval  $[0, 1]$ , defined by  $\mathcal{P}_k = \{[0, 1/2^k), [1/2^k, 2/2^k), \dots, [1 - 1/2^k, 1]\}$ . Each partition induces a partition of  $\mathcal{R}$  into vertical strips, as illustrated in Figure 5a. Within the  $i$ th strip, denote the region above the diagonal as  $r_{ia}$ , the region below the diagonal as  $r_{ib}$ , and the diagonal segment within the strip as  $r_{id}$ . Let  $i^*$  denote the unique index for the strip containing the vertical line  $x = b$ . Then, for each  $F$ , partition  $\mathcal{P}_k$  allows the objective function in (9) to be rewritten,

$$\frac{\sum_{i^* < i \leq 2^k} \int_{r_{ia} \cup r_{ib} \cup r_{id}} L dF + \int_{\{x \in (b, i^*/2^k)\} \cap \mathcal{R}} L dF}{\sum_{1 \leq i < i^*} \int_{r_{ia} \cup r_{ib} \cup r_{id}} L dF + \int_{\{x \in [(i^* - 1)/2^k, b]\} \cap \mathcal{R}} L dF} \quad (10)$$

$L$  is continuous and bounded over  $\mathcal{R}$ . So we may bound (10) from above by reallocating the probability mass that  $F$  assigns within each region/diagonal segment in each strip. All of the mass is reassigned to a point in the relevant region/segment, within  $\frac{1}{2^k}$  distance from where  $L$  is largest (when  $x > b$ ) or smallest (when  $x \leq b$ ). These locations at which  $L$  takes its largest and smallest values are *limit points*<sup>7</sup> of the respective regions/diagonal segment within each strip, as illustrated in Figure 5b. The reallocations define a prior  $F_k^*$  with a discrete marginal distribution of  $pdf$ . For each  $k$ ,  $F_k^*$  satisfies  $\frac{\int_{\{x > b\} \cap \mathcal{R}} L dF_k^*}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF_k^*} > \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF}$ .

7. Definition: for the ‘‘open balls’’ topology associated with the 2D Euclidean plane, a *limit point* of a subset of the plane is a point that is arbitrarily well-approximated by sequences of points within the subset [61], [62].

By construction, the objective function values from the  $F_k^*$  converge to the objective function value for some prior  $F^*$  with continuous marginal density  $f(x)$ . So, for each  $F$ ,

$$\frac{\int_{\{x > b\} \cap \mathcal{R}} L dF^*}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF^*} \geq \inf_k \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF_k^*}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF_k^*} \geq \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF} \quad (11)$$

Since this holds for any feasible prior  $F \in \mathcal{D}$ , we have

$$\sup_{\mathcal{D}^*} \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF^*}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF^*} \geq \sup_{\mathcal{D}} \inf_k \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF_k^*}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF_k^*} \geq \sup_{\mathcal{D}} \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF} \quad (12)$$

where  $\mathcal{D}^*$  contains all of the  $F^*$  priors. Because the objective function values for priors in  $\mathcal{D}^*$  are the limits of objective function values for feasible priors in  $\mathcal{D}$ , we also have

$$\sup_{\mathcal{D}^*} \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF^*}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF^*} \leq \sup_{\mathcal{D}} \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF} \quad (13)$$

Thus (12) and (13) imply three equivalent forms of optimisation,

$$\sup_{\mathcal{D}^*} \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF^*}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF^*} = \sup_{\mathcal{D}} \inf_k \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF_k^*}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF_k^*} = \sup_{\mathcal{D}} \frac{\int_{\{x > b\} \cap \mathcal{R}} L dF}{\int_{\{x \leq b\} \cap \mathcal{R}} L dF} \quad (14)$$

So, we can restrict the optimisation to sequences of priors  $\{F_k^*\}$ .

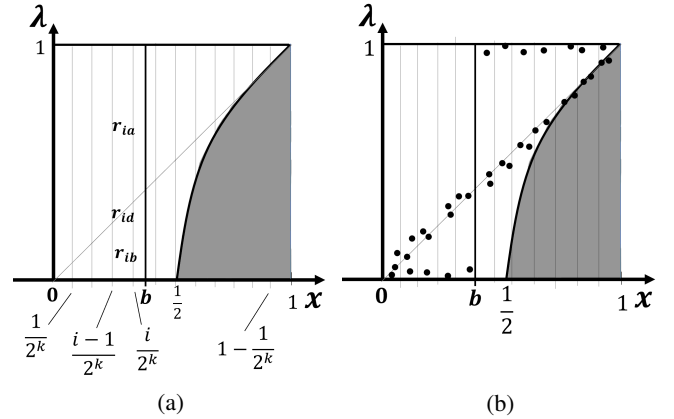


Fig. 5: (a)  $\mathcal{P}_k$  partitions  $\mathcal{R}$  into vertical strips; (b) example support of  $F_k^*$

For all sufficiently large  $k$ , the width of the strips can be made as small as we please. Consequently, by considering sufficiently large  $k$ , we may treat the location of masses within each strip as lying on the same vertical line, with masses on the diagonal segment or on the  $\lambda = 0, 1$  borders of  $\mathcal{R}$ . Consider then an arbitrary prior  $F_k^*$  (with discrete marginal) for sufficiently large  $k$ . The probability masses in a pair of strips can be reallocated within each strip to construct a new prior that gives a larger objective function value. One does this as follows.

Let the functions  $g_l(x)$  and  $g_u(x)$  be as defined in (7). Denote the unique  $x$  values at which  $g_l(x)$  and  $g_u(x)$  attain their maxima as  $x_l$  and  $x^u$ , respectively. There are 4 possibilities for reallocating probability masses, based on the relative sizes of  $x_l$ ,  $x^u$  and  $b$ .

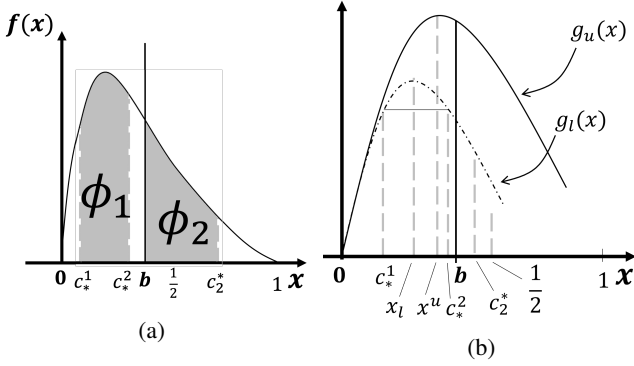


Fig. 6: (a) An example density  $f(x)$ ; (b) example maxima for  $g_l, g_u$

**Case 1)  $x_l < b$  and  $x^u < b$**

Let  $F_k^*$  be as depicted in Figure 5b. Consider two vertical strips, as shown in Figure 8a. The strips lie to the left of the vertical line  $x = x_l$ . For  $\Delta = 0$ , let the probabilities  $M_1 - \Delta, M_2 + \Delta, M_3 + \Delta$  and  $M_4 - \Delta$  be initially assigned to the 4 depicted locations (2 in each strip). These “ $M$ ”s are constant and consistent with the PKs, and  $\Delta$  is a sufficiently small probability mass. The derivative of the objective function with respect to  $\Delta$  exists, because the objective function is a rational function of  $\Delta$ . The sign of this derivative is determined by the function  $g_l(x)$  in Figure 6b. That is, the expression for the derivative is negative *iff*  $g_l(x_2) - g_l(x_1) > 0$  (where, for  $x_1 < x_2 < x_l < b$ ,  $x_1$  is in the leftmost strip and  $x_2$  is in the other strip). But this is true because  $g_l(x)$  is unimodal.

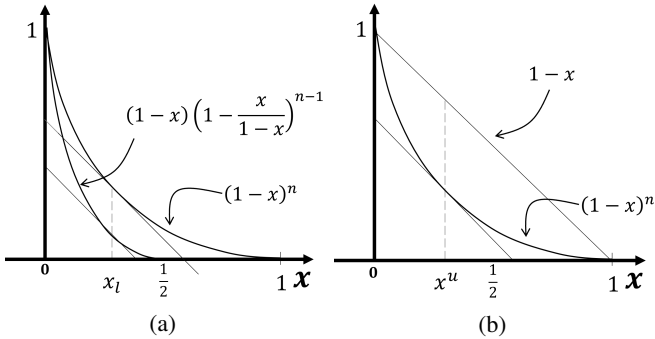


Fig. 7: Plots of the pair of functions that define (a)  $g_l$ , (b)  $g_u$

The unimodality of both  $g_l(x)$  and  $g_u(x)$  can be seen from arguments illustrated by Figures 7a and 7b. These figures depict the pair of functions that define each “ $g$ ” function. Each pair consists of two convex, monotonically decreasing functions that are equal at  $x = 0$ . In Figure 7a, over  $0 \leq x \leq \frac{1}{2}$ , the pair of functions are initially relatively convex (with relative derivative 1 at  $x = 0$ ) then relatively concave (with relative derivative 0 at  $x = \frac{1}{2}$ ). While in Figure 7b, the pair of functions are relatively concave over  $0 \leq x \leq \frac{1}{2}$  (with relative derivatives  $n$  and  $n(\frac{1}{2})^{n-1}$ , at  $x = 0$  and  $x = \frac{1}{2}$  respectively). Because of these, in each figure, the functions have the same tangent slope at a nontrivial  $x$  value in their shared domain. This is the  $x$  value at which the respective “ $g$ ” function attains its maximum – the values  $x_l$  and  $x^u$ . These values lie in the interval  $0 \leq x < \frac{1}{2}$ .

Since the objective function’s derivative with respect to  $\Delta$  is negative,  $\Delta$  should be made as small as possible, which makes the objective function as large as possible. Roughly speaking, mass

in the “ $x_1$ ” strip should be placed on the diagonal, while mass in the “ $x_2$ ” strip should be placed along the  $\lambda = 0$  line. Similar arguments justify the mass re-allocations illustrated in Figures 8b and 8c, using  $g_l(x)$  and  $g_u(x)$  respectively.

The general rule is, for a pair of strips containing  $x$  values less than  $b$ , the strip that is closest to containing  $x_l$  should have as much mass as possible below the diagonal, while the other strip should have as much mass as possible on or above the diagonal. Similarly, for two strips with  $x$  values greater than  $b$ , the strip closest to containing  $x^u$  should have as much mass as possible above the diagonal, while the other strip should have as much mass as possible on or below the diagonal.

So, by construction, a discrete prior  $F_k^*$  (e.g. Figure 5b) is replaced by a more extreme  $F_k^{**}$  (e.g. Figure 8d). Further reallocation is impossible when  $c_*^1, c_*^2, c_1^*$  and  $c_2^*$  have been found that solve

$$\begin{aligned} & \arg \min_{0 \leq r < s \leq b} |g_l(r) - g_l(s)|, \quad \arg \min_{b \leq v < w \leq 1} |g_u(v) - g_u(w)| \\ \text{s.t.} \quad & g_l(0) \leq g_l(r), \quad g_l(b) \leq g_l(s), \\ & g_u(b) \leq g_u(v), \quad g_u(1) \leq g_u(w), \\ & \int_{\{x \in [r, s]\} \cap \mathcal{R}} dF_k^{**} = \int_r^s f(x) dx = \phi_1, \\ & \int_{\{x \in [b, w]\} \cap \mathcal{R}} dF_k^{**} = \int_b^w f(x) dx = \phi_2, \\ & 0 < r < x_l < s \leq b, \quad 0 < x^u < b \leq v < w \leq 1 \end{aligned}$$

In particular, since  $x^u < b$  implies  $c_1^* = b$ , we can restrict the optimisation to these more extreme priors  $F_k^{**}$ . For such priors, the objective function (10) is comprised of sums that are integrals (with respect to  $F_k^{**}$ ) of simple functions. That is:

$$\begin{aligned} & \frac{\sum_{i^* < i \leq 2^k} \int_{\mathbf{r}_{ia} \cup \mathbf{r}_{ib} \cup \mathbf{r}_{id}} L dF_k^{**} + \int_{\{x \in (b, 2^{i^*}/2^k)\} \cap \mathcal{R}} L dF_k^{**}}{\sum_{1 \leq i < i^*} \int_{\mathbf{r}_{ia} \cup \mathbf{r}_{ib} \cup \mathbf{r}_{id}} L dF_k^{**} + \int_{\{x \in [2^{i^*-1}/2^k, b]\} \cap \mathcal{R}} L dF_k^{**}} \\ &= \frac{\int_{\{x \in (b, c_2^*)\} \cap \mathcal{R}} L dF_k^{**} + \int_{\{x \in (c_2^*, 1)\} \cap \mathcal{R}} L dF_k^{**}}{\int_{\{x \in (c_*^1, c_*^2)\} \cap \mathcal{R}} L dF_k^{**} + \int_{\{x \in (0, c_*^1) \cup (c_*^2, b)\} \cap \mathcal{R}} L dF_k^{**}} \\ &= \frac{\sum_{i^* < i \leq 2^k} (L_1 \mathbf{1}_{x_i \in (b, c_2^*)} + L_3 \mathbf{1}_{x_i \in (c_2^*, 1)}) \int_{\text{ith strip}} dF_k^{**}}{\sum_{1 \leq i \leq i^*} (L_2 \mathbf{1}_{x_i \in (c_*^1, c_*^2)} + L_3 \mathbf{1}_{x_i \in (0, c_*^1) \cup (c_*^2, b)}) \int_{\text{ith strip}} dF_k^{**}} \\ &= \frac{\sum_{i^* < i \leq 2^k} (L_1 \mathbf{1}_{x_i \in (b, c_2^*)} + L_3 \mathbf{1}_{x_i \in (c_2^*, 1)}) \int_{\frac{i-1}{2^k}}^{\frac{i}{2^k}} f(x) dx}{\sum_{1 \leq i \leq i^*} (L_2 \mathbf{1}_{x_i \in (c_*^1, c_*^2)} + L_3 \mathbf{1}_{x_i \in (0, c_*^1) \cup (c_*^2, b)}) \int_{\frac{i-1}{2^k}}^{\frac{i}{2^k}} f(x) dx} \end{aligned}$$

where  $x_i \in [\frac{i-1}{2^k}, \frac{i}{2^k}]$ ,  $L_1 := (1 - x_i)$ ,  $L_2 := \frac{(1-2x_i)^{n-1}}{(1-x_i)^{n-2}}$  and  $L_3 := (1 - x_i)^n$

By the *dominated convergence theorem* (see [46]), the continuity of  $L$  over  $\mathcal{R}$  implies that the sums converge to integrals with respect to the density  $f(x)$  as  $k \rightarrow \infty$ , so

$$\frac{\int_b^1 \left( (1-x) \mathbf{1}_{x \in (b, c_2^*)} + (1-x)^n \mathbf{1}_{x \in (c_2^*, 1)} \right) f(x) dx}{\int_0^b \left( \frac{(1-2x)^{n-1}}{(1-x)^{n-2}} \mathbf{1}_{x \in (c_*^1, c_*^2)} + (1-x)^n \mathbf{1}_{x \in (0, c_*^1) \cup (c_*^2, b)} \right) f(x) dx}$$

**Case 2)  $x_l < b$  and  $x^u > b$**

An analogous argument to case 1 gives the solution

$$\frac{\int_b^1 \left( (1-x) \mathbf{1}_{x \in (c_1^*, c_2^*)} + (1-x)^n \mathbf{1}_{x \in (b, c_1^*) \cup (c_2^*, 1)} \right) f(x) dx}{\int_0^b \left( \frac{(1-2x)^{n-1}}{(1-x)^{n-2}} \mathbf{1}_{x \in (c_*^1, c_*^2)} + (1-x)^n \mathbf{1}_{x \in (0, c_*^1) \cup (c_*^2, b)} \right) f(x) dx}$$

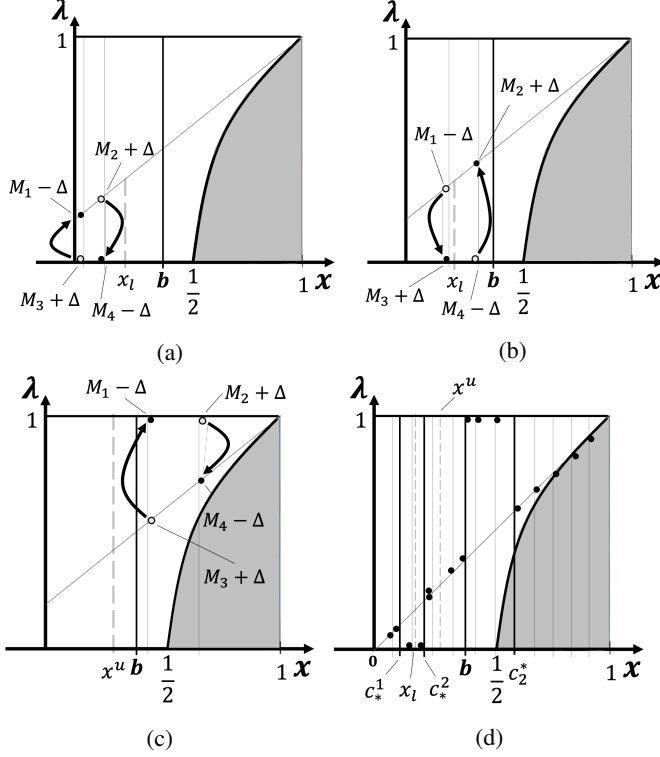


Fig. 8:  $F_k^{**}$  is constructed from  $F_k^*$  by reallocating probability mass

where  $c_1^*$ ,  $c_2^*$ ,  $c_1^*$  and  $c_2^*$  have been identified that solve

$$\begin{aligned} & \arg \min_{0 \leq r < s \leq b} |g_l(r) - g_l(s)|, \quad \arg \min_{b \leq v < w \leq 1} |g_u(v) - g_u(w)| \\ \text{s.t.} \quad & g_l(0) \leq g_l(r), \quad g_l(b) \leq g_l(s), \\ & g_u(b) \leq g_u(v), \quad g_u(1) \leq g_u(w), \\ & \int_{\{x \in [r, s]\} \cap \mathcal{R}} dF_k^{**} = \int_{[r, s]} f(x) dx = \phi_1, \\ & \int_{\{x \in [v, w]\} \cap \mathcal{R}} dF_k^{**} = \int_{[v, w]} f(x) dx = \phi_2, \\ & 0 < r < x_l < s \leq b, \quad 0 < b \leq v < x^u < w \leq 1 \end{aligned}$$

**Case 3)**  $x_l > b$  and  $x^u < b$

An analogous argument to case 1 gives the solution

$$\frac{\int_b^1 \left( (1-x) \mathbf{1}_{x \in (b, c_2^*)} + (1-x)^n \mathbf{1}_{x \in (c_2^*, 1)} \right) f(x) dx}{\int_0^b \left( \frac{(1-2x)^{n-1}}{(1-x)^{n-2}} \mathbf{1}_{x \in (c_1^*, b)} + (1-x)^n \mathbf{1}_{x \in (0, c_1^*)} \right) f(x) dx}$$

where  $c_1^*$ ,  $c_2^*$ ,  $c_1^*$  and  $c_2^*$  have been identified that solve

$$\begin{aligned} & \arg \min_{0 \leq r < s \leq b} |g_l(r) - g_l(s)|, \quad \arg \min_{b \leq v < w \leq 1} |g_u(v) - g_u(w)| \\ \text{s.t.} \quad & g_l(0) \leq g_l(r), \quad g_l(b) \leq g_l(s), \\ & g_u(b) \leq g_u(v), \quad g_u(1) \leq g_u(w), \\ & \int_{\{x \in [r, b]\} \cap \mathcal{R}} dF_k^{**} = \int_r^b f(x) dx = \phi_1, \\ & \int_{\{x \in [b, w]\} \cap \mathcal{R}} dF_k^{**} = \int_b^w f(x) dx = \phi_2, \\ & 0 < r < s \leq b < x_l, \quad 0 < x^u < b \leq v < w \leq 1 \end{aligned}$$

In particular, because  $x_l > b$ ,  $x^u < b$ , we have  $c_2^* = b$ ,  $c_1^* = b$ .

**Case 4)**  $x_l > b$  and  $x^u > b$

An analogous argument to case 1 gives the solution

$$\frac{\int_b^1 \left( (1-x) \mathbf{1}_{x \in (c_1^*, c_2^*)} + (1-x)^n \mathbf{1}_{x \in (b, c_1^*) \cup (c_2^*, 1)} \right) f(x) dx}{\int_0^b \left( \frac{(1-2x)^{n-1}}{(1-x)^{n-2}} \mathbf{1}_{x \in (c_1^*, b)} + (1-x)^n \mathbf{1}_{x \in (0, c_1^*)} \right) f(x) dx}$$

where  $c_1^*$ ,  $c_2^*$ ,  $c_1^*$  and  $c_2^*$  have been identified that solve

$$\begin{aligned} & \arg \min_{0 \leq r < s \leq b} |g_l(r) - g_l(s)|, \quad \arg \min_{b \leq v < w \leq 1} |g_u(v) - g_u(w)| \\ \text{s.t.} \quad & g_l(0) \leq g_l(r), \quad g_l(b) \leq g_l(s), \\ & g_u(b) \leq g_u(v), \quad g_u(1) \leq g_u(w), \\ & \int_{\{x \in [r, b]\} \cap \mathcal{R}} dF_k^{**} = \int_r^b f(x) dx = \phi_1, \\ & \int_{\{x \in [v, w]\} \cap \mathcal{R}} dF_k^{**} = \int_v^w f(x) dx = \phi_2, \\ & 0 < r < s \leq b < x_l, \quad 0 < b \leq v < x^u < w \leq 1 \end{aligned}$$

In particular, because  $x_l > b$ , we must have  $c_2^* = b$ .  $\square$

## APPENDIX C ASYMPTOTICS OF POSTERIOR CONFIDENCE BASED ON FAILURE-FREE OPERATION

*Claim.* In the theorem,  $\lim_{n \rightarrow \infty} Q = \begin{cases} 0, & \text{if } \phi_2 = 0 \\ \infty, & \text{if } \phi_2 > 0 \end{cases}$ . Since the assessor's conservative posterior confidence in the bound  $b$  is  $\frac{1}{1+Q}$ , the assessor either becomes certain that  $b$  has been satisfied, or they become certain that it has not.

*Proof.* We will show that

$$\frac{\int_b^1 \left( (1-x)^n \mathbf{1}_{x \in (b, c_1^*) \cup (c_2^*, 1)} + (1-x) \mathbf{1}_{x \in (c_1^*, c_2^*)} \right) f(x) dx}{\int_0^b \left( (1-x)^n \mathbf{1}_{x \in (0, c_1^*) \cup (c_2^*, b)} + \frac{(1-2x)^{n-1}}{(1-x)^{n-2}} \mathbf{1}_{x \in (c_1^*, c_2^*)} \right) f(x) dx} \xrightarrow{n \rightarrow \infty} \begin{cases} 0, & \text{if } \phi_2 = 0 \\ \infty, & \text{if } \phi_2 > 0 \end{cases} \quad (15)$$

Since  $\phi_2 = 0$  implies  $c_1^* = c_2^*$ , the lhs of (15), i.e.  $Q$ , becomes

$$\frac{\int_b^1 (1-x)^n f(x) dx}{\int_0^b \left( (1-x)^n \mathbf{1}_{x \in (0, c_1^*) \cup (c_2^*, b)} + \frac{(1-2x)^{n-1}}{(1-x)^{n-2}} \mathbf{1}_{x \in (c_1^*, c_2^*)} \right) f(x) dx} \quad (16)$$

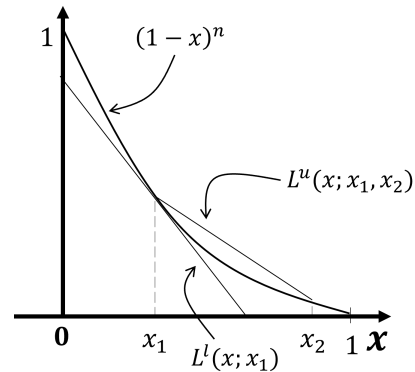


Fig. 9: Geometric illustration of Jensen's inequalities

The integrals of  $(1-x)^n$  in (16) can be bounded by a suitable choice of straight lines. We construct these as follows.

For constants  $x_1$  and  $x_2$  such that  $0 \leq x_1 < x_2 \leq 1$ , define the straight lines  $L^u(x; x_1, x_2)$  and  $L^l(x; x_1)$  (see Figure 9):

$$L^u(x; x_1, x_2) = (1 - x_1)^n \left( \frac{x_2 - x}{x_2 - x_1} \right) + (1 - x_2)^n \left( \frac{x - x_1}{x_2 - x_1} \right)$$

$$L^l(x; x_1) = \left( n(x_1 - x) + 1 - x_1 \right) (1 - x_1)^{n-1}$$

The curve  $(1 - x)^n$  is convex, so  $L^u$  lies above the curve when  $x_1 < x < x_2$ . While  $L^l$  is tangent at  $x = x_1$ , so lies below the curve. These are *Jensen's inequalities* [63]. Therefore:

$$L^l(\mathbb{E}_*^1 X; \mathbb{E}_*^1 X) \leq \frac{\int_0^{c_1^*} (1 - x)^n f(x) dx}{\int_0^{c_1^*} f(x) dx} \leq L^u(\mathbb{E}_*^1 X; 0, c_1^*) \quad (17)$$

$$L^l(\mathbb{E}_*^2 X; \mathbb{E}_*^2 X) \leq \frac{\int_{c_2^*}^b (1 - x)^n f(x) dx}{\int_{c_2^*}^b f(x) dx} \leq L^u(\mathbb{E}_*^2 X; c_2^*, b) \quad (18)$$

$$L^l(\mathbb{E} X; \mathbb{E} X) \leq \frac{\int_b^1 (1 - x)^n f(x) dx}{\int_b^1 f(x) dx} \leq L^u(\mathbb{E} X; b, 1) \quad (19)$$

where  $\mathbb{E}_*^1 X = \frac{\int_0^{c_1^*} x f(x) dx}{\int_0^{c_1^*} f(x) dx}$ ,  $\mathbb{E}_*^2 X = \frac{\int_{c_2^*}^b x f(x) dx}{\int_{c_2^*}^b f(x) dx}$ ,  $\mathbb{E} X = \frac{\int_b^1 x f(x) dx}{\int_b^1 f(x) dx}$ .

Using the bounds in (17)–(19), we can bound  $Q$  (i.e. (16)):

$$0 \leq Q \leq \frac{L^u(\mathbb{E} X; b, 1) \int_b^1 f(x) dx}{\frac{L^l(\mathbb{E}_*^1 X; \mathbb{E}_*^1 X) \int_0^{c_1^*} f(x) dx + L^l(\mathbb{E}_*^2 X; \mathbb{E}_*^2 X) \int_{c_2^*}^b f(x) dx}{\left( \frac{1 - \mathbb{E} X}{1 - b} \right) \int_b^1 f(x) dx + \left( \frac{1 - \mathbb{E}_*^2 X}{1 - b} \right) \int_{c_2^*}^b f(x) dx}} \quad (20)$$

We used  $\int_0^b \frac{(1-2x)^{n-1}}{(1-x)^{n-2}} \mathbf{1}_{x \in (c_1^*, c_2^*)} f(x) dx > 0$  to bound  $Q$  from above – by removing this term from  $Q$ 's denominator. Since  $0 < \mathbb{E}_*^2 X < b < 1$ , we have  $\left( \frac{1 - \mathbb{E}_*^2 X}{1 - b} \right) > 1$ . So, as  $n \rightarrow \infty$  in (20),  $c_2^*$  tends to a non-zero value less than  $b$ ,  $\int_{c_2^*}^b f(x) dx$  tends to a non-zero value less than 1,  $\mathbb{E}_*^2 X$  tends to a non-zero value less than  $b$ ,  $\left( \frac{1 - \mathbb{E}_*^2 X}{1 - b} \right)^n$  tends to  $\infty$ , and  $\lim_{n \rightarrow \infty} Q = 0$ .

If instead,  $\phi_2 > 0$ , then<sup>8</sup>  $c_1^* < c_2^*$  and  $Q$  is the quotient on the lhs of (15). As  $n \rightarrow \infty$ , integrals of  $(1 - x)^n$  in  $Q$  all tend to 0 by the *monotone convergence theorem* (m.c.t.) [46]. The m.c.t. also implies  $\lim_{n \rightarrow \infty} \int_0^b \frac{(1-2x)^{n-1}}{(1-x)^{n-2}} \mathbf{1}_{x \in (c_1^*, c_2^*)} f(x) dx = 0$ . Therefore,  $\lim_{n \rightarrow \infty} Q = \infty$ .  $\square$

8. In particular,  $\int_{c_1^*}^{c_2^*} (1 - x) f(x) dx > 0$

## APPENDIX D

### ALGORITHM FOR NUMERICAL ESTIMATES OF $c_*^1, c_*^2$

For brevity, we omit the analogous algorithm for  $c_*^1, c_*^2$ .

---

#### Algorithm Bisection Method based Algorithm for $c_*^1, c_*^2$

---

**Input:** The *pdf* density  $f(x)$ , an intermediate function  $g_l(x)$ , the target *pdf* bound  $b$ , a tolerance  $\epsilon$  and the doubts  $\phi_1, \phi_2$ .

**Output:**  $c_*^1, c_*^2$

```

1: if  $\int_0^b f(u) du > \phi_1$  then
2:    $x_l = \arg \max g_l(x)$ 
3:   if  $x_l > b$  then
4:      $c_*^1 = \text{solve}(\int_x^b f(u) du = \phi_1, x \in [0, b])$ 
5:      $c_*^2 = b$ ; return  $c_*^1, c_*^2$ 
6:   else
7:      $c = \text{solve}(g_l(x) = g_l(b), x \in [0, x_l])$ 
8:     if  $\int_c^b f(u) du < \phi_1$  then
9:        $c_*^1 = \text{solve}(\int_x^b f(u) du = \phi_1, x \in [0, b])$ 
10:       $c_*^2 = b$ ; return  $c_*^1, c_*^2$ 
11:     else  $\triangleright$  Start of the bisection method
12:        $c_*^2 = b$ 
13:        $tmp_{lb} = x_l$ 
14:        $tmp_{ub} = b$ 
15:        $tmp_{\phi_1} = \int_c^b f(u) du$ 
16:       while  $|tmp_{\phi_1} - \phi_1| > \epsilon$  do
17:         if  $tmp_{\phi_1} > \phi_1$  then
18:            $tmp_{ub} = c_*^2$ 
19:            $c_*^2 = \frac{c_*^2 + tmp_{lb}}{2}$ 
20:         else
21:            $tmp_{lb} = c_*^2$ 
22:            $c_*^2 = \frac{c_*^2 + tmp_{ub}}{2}$ 
23:         end if
24:          $c_*^1 = \text{solve}(g_l(x) = g_l(c_*^2), x \in [0, x_l])$ 
25:          $tmp_{\phi_1} = \int_{c_*^1}^{c_*^2} f(u) du$ 
26:       end while
27:     end if  $\triangleright$  End of the bisection method
28:     return  $c_*^1, c_*^2$ 
29:   end if
30: else  $\triangleright$  This is the case when  $\int_0^b f(u) du < \phi_1$ 
31:   print("PK4 violated!")
32: end if

```

---