



City Research Online

City, University of London Institutional Repository

Citation: Bloomfield, R. E. & Rushby, J. (2020). Assurance 2.0. City, University of London.

This is the accepted version of the paper.

This version of the publication may differ from the final published version.

Permanent repository link: <https://openaccess.city.ac.uk/id/eprint/24093/>

Link to published version:

Copyright: City Research Online aims to make research outputs of City, University of London available to a wider audience. Copyright and Moral Rights remain with the author(s) and/or copyright holders. URLs from City Research Online may be freely distributed and linked to.

Reuse: Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Assurance 2.0: A Manifesto

Robin Bloomfield

Adelard LLP and City, University of London
London UK

John Rushby

Computer Science Laboratory
SRI International, Menlo Park CA USA

Abstract—System assurance is confronted by significant challenges. Some of these are new, for example, autonomous systems with major functions driven by machine learning and AI, and ultra-rapid system development, while others are the familiar, persistent issues of the need for efficient, effective and timely assurance. Traditional assurance is seen as a brake on innovation and often costly and time consuming. We therefore propose a modernized framework, “Assurance 2.0,” as an enabler that supports innovation and continuous incremental assurance. Perhaps unexpectedly, it does so by making assurance more rigorous, with increased focus on the reasoning and evidence employed, and explicit identification of defeaters and counterevidence.

1. Introduction

Assurance is often seen as a drag on innovation and as an activity that is additional to (and generally comes after) the “real work” of design and implementation. We instead propose that assurance can be an enabler for innovation and a constructive element in a holistic design process. However, if assurance is employed from the early stages of design, it will necessarily be incomplete at those stages, so we need some measures to indicate if we are headed in the “right direction” and to help prioritize issues and solutions. Counterintuitively, perhaps, we propose that the way to address these and other concerns that we will introduce later, is by making assurance more rigorous, in a framework that we call “Assurance 2.0.”

This framework aims to support reasoning and communication about the behavior and trustworthiness of engineered systems and, ultimately, their certification. It builds on the notion of an “Assurance Case,” where claims about the system are justified by an argument based on evidence. In particular, it maintains a representation of the *structure* of the argument as a tree of claims linked by argument steps and supported by evidence (e.g., Figure 7) as in ASCAD CAE [1] and GSN [2]¹, but strengthens it with increased focus on the evidence and the reasoning (both logical and probabilistic) employed, and on exploration and assessment

of doubts and “defeaters.” We introduce the ideas in this section, and give details (and references) in subsequent ones.

In current practice, steps in an assurance argument are often *inductive*,² meaning the subclaims strongly support the parent claim, but do not ensure it, as a *deductive* step would. In Assurance 2.0 we advocate that argument steps should be deductive, and this can require additional evidence. For example, argument steps often iterate over some enumeration (e.g., over components, or over hazards) and for this to be deductive we need evidence that the enumeration is complete and that the claim distributes over its elements. In cases where it seems impossible to provide a deductive step, the “gap” must be acknowledged and given special attention. To support these recommendations, we advocate use of pre-analyzed argument templates such as *CAE Blocks* [3], which provide mechanisms for separating inductive and deductive lines of reasoning, and for managing the side conditions necessary to justify deductive steps and excuse inductive ones. This insistence that reasoning steps should be “as deductive as possible and inductive only as strictly necessary” is one of the ways in which Assurance 2.0 strengthens traditional assurance; deductive reasoning steps ensure that doubts have nowhere to hide and thereby help identify weak spots and focus attention in productive directions.

Arguments are grounded on evidence and we advocate explicit assessment of the “weight” of evidence offered in support of a claim. It is not enough for evidence to support a claim; it must also discriminate between a claim and its negation or *counterclaim*. We recommend interpreting “weight” using ideas and measures from *Confirmation Theory*, which do exactly this. Again, this aspect of Assurance 2.0 is more demanding than traditional estimates for the strength of evidence and requires explicit consideration of counterclaims. Claims supported by sufficient weight of evidence may be used as premises in a logical interpretation of the overall assurance argument and when, in addition, all the reasoning steps are deductive, we have a deductive thread from facts, established by evidence, to the top level claim and thereby satisfy a benchmark for informal reasoning known as *Natural Language Deductivism* (NLD).

1. We use a variant on CAE terminology: we say *claim* where GSN says *goal*, we say *argument step* where CAE says simply *argument* (and GSN says *strategy*) and we use *argument* for the whole tree of claims and argument steps. Our diagrams use the CAE style.

2. This is an unfortunate choice of words as the same term is used with several other meanings in mathematics and logic.

Although it is primarily motivated by practical considerations and by experience with current methods, the Assurance 2.0 framework aligns with modern developments in epistemology (notably, confirmation theory and NLD), and we strengthen this alignment through use of “indefeasibility” as the criterion for justified belief (e.g., that an assurance case establishes its claim). For a belief to be justified, the *Indefeasibility Criterion* requires that we must be so sure that all doubts and objections have been attended to that there is no (or, more realistically, we cannot imagine any) new information that would cause us to change our evaluation.

Doubts and objections are exemplified as *defeaters* so the indefeasibility criterion applied to assurance cases requires a comprehensive search for defeaters to the argument. Once a potential defeater has been identified, it must itself be defeated, meaning that more detailed analysis shows that it is not, in fact, a defeater, or that the system and/or its assurance case are adjusted to negate it. In Assurance 2.0, we advocate that the search for defeaters, and their own defeat, should be systematized and documented as essential parts of the case (just as hazard analysis and the hazard log are essential parts of safety engineering). One systematic approach is through construction and dialectical consideration of *counterclaims* and *countercases*. Counterclaims arise naturally in confirmation measures and are discussed in Section 2.1, while a countercase is an assurance case for the negation of the top claim and is discussed in Section 3.2.2.

Confirmation bias—the tendency to interpret information in a way that confirms or strengthens our prior beliefs—is a natural hazard in assurance cases—after all, we are engaged in building a case to support the system. Competent and diligent external reviewers are good defenses against confirmation bias, but are typically involved only periodically and mostly toward the end of the development of a case. Several of the innovations in Assurance 2.0 are intended to provide systematic mitigations against confirmation bias at every step in the development of a case without the excessive conservatism that leads to prolix cases with unnecessary evidence presented “just in case,” or even to the rejection of good systems.

The paper is organized as follows. Section 2 describes the basic structure of an assurance case argument and the criteria for evaluating its soundness. Section 3 discusses confidence in the case and Section 4 provides brief conclusions.

2. Arguments, Step by Step

A key innovation in the development of modern assurance cases was the idea of a “structured safety case,” introduced in the 1970s, that required an *argument* to explain how the design of the system and the checks and tests performed during its development combine to ensure safety. Subsequent refinements in the 1990s led to the idea that the argument itself should be structured, that is, organized around goals or *claims*, and grounded on *evidence* about the system. Methods and notations such as Goal Structuring Notation (GSN) [2] and Claims, Argument, Evidence (CAE)

[1] emerged at this time and support a body of expertise and practice that thrives to this day.

The general structure of an assurance case argument is illustrated in later diagrams, such as Figure 7. An argument is organized as a tree of two kinds of basic *steps*: evidential (at the leaves) and reasoning (interior), which are described in the following subsections. Mixed forms are also possible.

2.1. Elementary Evidential Steps

Let us begin with the most basic kind of argument step: one where some item of evidence directly supports a claim. To make things concrete, we will suppose our examples are taken from a case in which random tests are used to support a claim of reliability (certain nuclear cases are like this [4]). One step in the argument for this case will concern soundness of the test oracle: that is, soundness of the means by which we judge the correctness of test outcomes. Figure 1 portrays this step: at the top is the (sub)claim that the oracle is sound (which will be backed by a description of what it means for an oracle to be sound); at the bottom is a description of the evidence for its soundness (which will be backed by reference to files containing the actual evidence), and in between is an argument that the evidence does indeed guarantee the claim; we say that this argument is one for *evidence incorporation* and we refer to the whole argument step (i.e., claim, argument, and evidence) as an *evidential step*.

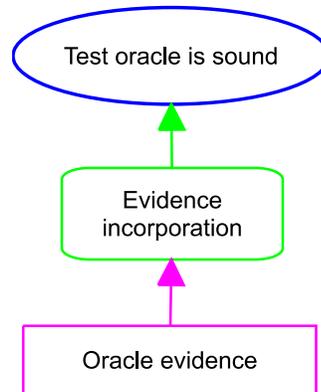


Figure 1. Elementary Evidential Step

Implicit in the previous sentence is the idea that claims and subclaims are logical propositions: that is, statements about the “world” (by which we mean the system of interest and its environment) that may be true or false. Evidence, on the other hand, is a description or pointer to some observation or experiment on the world. An argument for evidence incorporation documents a human assessment that the evidence persuasively attests the truth of the claim.

This assessment may be informal, or it may employ some systematic process. In the latter case, it is usual to talk of *weighing* the evidence and of accepting the claim when the weight of its supporting evidence crosses some

threshold. This raises the question of how weighing is performed and what units are employed. A standard treatment uses probabilities: if e is some evidence then $P(e)$ is the probability of seeing this evidence. Although it is possible to construct frequentist interpretations for this quantity, it is generally interpreted as a *subjective probability*, that is, a human judgment of likelihood expressed numerically from 0 (impossible) to 1 (certain). Similarly $P(c)$ is the subjective probability that the claim c is true. We might consider this the “background” or *prior* probability, which is then “boosted” by the evidence e to the *posterior* probability $P(c|e)$. Thus, $P(c|e) > t$ for some t might be considered a suitable criterion for accepting c on the basis of e .

Let us suppose that the evidence for soundness of our oracle is that it was extensively validated against the previous version of the system. This seems like fairly strong evidence so we might make the qualitative assessment that $P(c|e)$ is “high.” However, a critic might say that if the evidence is about a previous version of the system, how relevant can it be to soundness of the oracle for this version? A sharp and general version of this question asks whether the evidence can discriminate between a claim and its negation, or counterclaim. This suggests the weight of evidence should not be based on $P(c|e)$ alone, but should also consider the *difference* between this value and $P(\neg c|e)$. Difference can be measured as a ratio, or as arithmetic difference.

An attractive variant turns these conditional probabilities around: instead of the posterior probability of the claim $P(c|e)$, we consider the *likelihood* of the evidence given the claim, $P(e|c)$, and compare this to its likelihood given the counterclaim, $P(e|\neg c)$. Likelihood and posterior probability are related by Bayes’ rule and so the choice of one over the other might seem moot. However, it is often easier to estimate the likelihood of concrete observations, given a claim about the world, than vice-versa (i.e., it is easier to estimate a likelihood than a posterior). Furthermore, the likelihood $P(e|c)$ has a more “causal” flavor—we think of (the property underlying) the claim causing the evidence rather than vice-versa.

These ideas, and the general topics of evaluating and measuring “weight of evidence,” date back to the World War II codebreaking work of Turing and Good [5], where Good’s original measure for weight of evidence was $\log \frac{P(e|c)}{P(e|\neg c)}$. Today, these topics are studied in Bayesian Confirmation Theory (a subfield of Bayesian Epistemology [6]) and many *confirmation* (i.e., weight) *measures* have been proposed [7]. Among these, that of Kemeny and Oppenheim is popular:

$$\frac{P(e|c) - P(e|\neg c)}{P(e|c) + P(e|\neg c)}.$$

This measure is positive for strong evidence, near zero for weak evidence, and negative for counterevidence.

Returning to our example, we need to estimate the likelihood of the evidence about the oracle (i.e., it exhibited good performance against a previous version of the system), given a) the claim that the oracle is sound, and b) the counterclaim that it is not. An oracle evaluates tests and their outcomes against requirements, so we need to ask whether

the requirements have changed between the previous and current versions of the system. Let us suppose the answer is “yes, a little.” It’s good that we asked, for the proffered evidence tells us nothing about the performance of the oracle against those requirements that have changed since the previous system—unless we know more about the oracle structure and the modularity of the requirements. Without further evidence about the nature of the requirements and the oracle, the Kemeny-Oppenheim measure is zero and we conclude that the proffered evidence is of no value.

Although in most cases we do not advocate assessment of numerical valuations for confirmation measures, nor their constituent probabilities, we believe that informal consideration as was done here (and “qualitative” assessments such as *small*, *medium*, and *large*) can provide significant benefits in the evaluation of evidence.

What are these benefits? There are just a couple of ways in which an assurance case can be flawed or, as we say, *defeated* [8]. One is that the evidence supporting a claim is inadequate to justify the confidence required; philosophers call this *undercutting defeat*. It could be that the evidence is merely insufficient (e.g., we did testing, but not enough of it) or it could be that there is a gap or flaw (e.g. the case just considered of an oracle evaluated against a previous version of the system). Confirmation measures, even when assessed informally, provide rational quantification for the weight of evidence and thereby guard against undercutting defeat.

The other kind of defeat is when there is evidence that contradicts a claim; this is called a *rebutting defeater*. Confirmation measures require consideration of the extent to which evidence supports counterclaims, and thereby force a search for rebutting defeaters within evidential steps.

Defeaters for an assurance case are rather like hazards for a critical system, and just as the search for hazards is an essential element in the engineering of critical systems, so the search for defeaters is an essential element in the evaluation of assurance cases. Confirmation measures are an attractive tool in this search as they identify both kinds of defeat in evidential steps and thereby provide a valuable and necessary antidote to *confirmation bias*, which some consider an endemic vulnerability in assurance cases [9].

This section considered only elementary evidential steps; a less elementary step may incorporate several items of evidence in support of a single claim. The overall confidence measure may then involve conditional probabilities and likelihoods for evidential items that are not independent of each other. Tools for Bayesian Belief Nets (BBNs) can assist in construction and evaluation of numeric models for these circumstances. Although we do not advocate numerical assessments for the probabilities involved, “what if” experiments with a range of possibilities can prove very enlightening. An example is given in [10].

2.2. Elementary Reasoning Steps

We have considered an elementary evidential argument step—one where we assess the extent to which evidence supports a claim—and now turn to a similarly elementary

reasoning step—one where several (sub)claims combine to support a parent claim. Figure 2 illustrates such a step. Here we suppose we have three subclaims concerning a test procedure, each supported by evidence or an entire subargument (these are not shown): one asserts that the test oracle is sound (as in the previous section), another that the test procedure is sound, and the third that the tested software is the actual software. The step asserts that if these three subclaims are true, then we may conclude that the overall test process is sound. Each of these claims will be backed by a description of what it means and bound together by an argument that the subclaims “lead to” the parent claim.

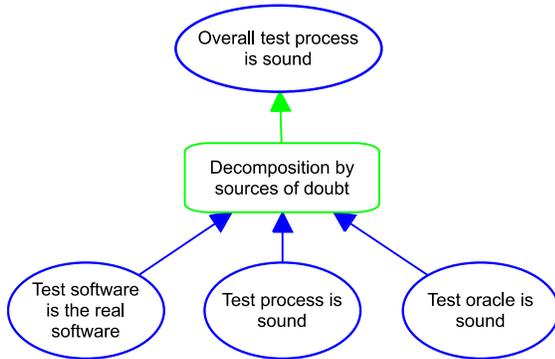


Figure 2. Elementary Reasoning Step

We say the subclaims “lead to” the parent because we have not yet established the relationship that is intended. In some early interpretations for an assurance case, the intended relationship was structural rather than logical: it simply indicated that the case for the parent claim decomposed into subcases for each of the three subclaims. In modern interpretations, the intended relationship is logical but it may be deductive (i.e., the subclaims imply or entail the parent claim) or inductive (i.e., the subclaims “suggest” the parent claim). In text presentations, an annotation on the central argument box can indicate which of these is intended.

When a deductive interpretation is indicated, the argument must make the case that the subclaims truly entail the parent claim. Sometimes a convincing case can be made with no additional information, but often an additional subclaim will be needed to substantiate the case. Logically, this additional subclaim is just like the others and conjoins with them to entail the parent claim; however, it is contextually somewhat different, so we call it a “side condition” or “side claim” and draw it in a different position and color (but same shape), as shown in Figure 3. In this case, we are claiming the three conditions considered in the original subclaims are the only threats to overall soundness of the testing process and the side condition, which asserts this, will need to be supported by evidence akin to hazard analysis to justify it.

Observe that some elements that may appear in claims (e.g., $\frac{x}{y}$) may not “make sense” unless another (assumption) claim (e.g., $y \neq 0$) is true. Since all the subclaims in an argument must be true if we are to conclude its top claim, we could allow each subclaim to be interpreted under the

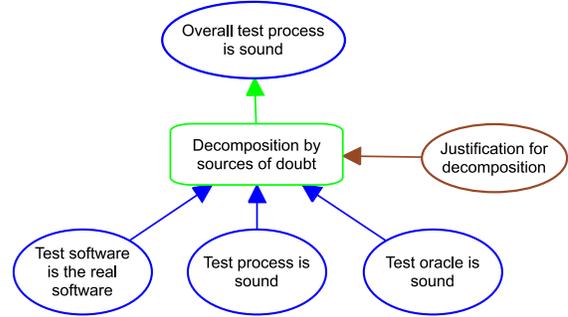


Figure 3. Reasoning Step with Side Condition

assumption that all other subclaims are true. However, it can require additional analysis to ensure there is no circularity in this reasoning, so a useful compromise is to impose some standard order of evaluation.

In Assurance 2.0, we advocate that all reasoning steps eventually should be deductive as this raises the bar on the quality of argumentation required, and is necessary to satisfy the indefeasibility criterion for justified belief in the overall argument. However, the precision and rigor we advocate needs to be reconciled with the need for concise communication and constructive progress during system development. Thus, we accept that argument steps may be inductive during the early stages of system development and assurance exploration. But it is desirable that tools should assist in keeping track of these transitional compromises. A workaround for tools based on deductivism would be to supply inductive steps with a nugatory “something missing here” side claim that is asserted to make the step deductive, but is unsupported by evidence. This allows progress, while the unsupported side claim acts as a constant reminder of the imperfection in this argument step.

The energetic search for defeaters is a rational guard against hubris and confirmation bias in the construction of assurance cases. For reasoning steps, a systematic organization of the search can be based on challenges to their deductiveness, which goes hand-in-hand with (re)formulation of their side conditions: strengthening a side condition, and its supporting evidence, is one way to defeat a successful defeater of this type.

This section considered only elementary reasoning steps: those where a claim is supported by subclaims. In less elementary reasoning steps, a claim may be supported by a combination of subclaims and evidence. The most useful construction of this kind is best interpreted not as a reasoning step, but as an evidential step with side claims that function as assumptions. Reference [11] provides more discussion of these topics.

3. Soundness and Confidence Assessment

In Assurance 2.0, the interpretation that we apply to an assurance case is a systematic instance of “Natural Language Deductivism” (NLD) [12], which regards its informal argument as an approximation to a deductively valid proof. NLD

differs from proof in formal mathematics and logic in that its premises are “reasonable or plausible” rather than certain, and hence its conclusions are likewise reasonable or plausible rather than certain. Our requirements that evidential steps cross some threshold for credibility (e.g., as assessed by a confirmation measure), and that a thorough search for defeaters persuades stakeholders that the case is indefeasible and all reasoning steps are deductive, systematizes what it means for the premises to be “reasonable or plausible” and thereby give us confidence that the overall argument is sound and the top claim is true. But then we might ask, how much confidence, and how much do we need?

Some assurance cases may be more persuasive than others, and not all (sub)systems need the highest levels of assurance: indeed, several standards speak of “Safety Integrity Levels” (SILs) from 1 (low) to 4 (high). Thus, we need ways to assess confidence in a case, and principled ways to organize cases so that the lower SILs are easier and cheaper to achieve. The confidence we need depends on the nature of the claim and the decision being made. In some cases (we call them “quantitative”), the claim may include a numeric estimate for some parameter (e.g., reliability) and our confidence then reflects epistemic uncertainty in this quantity. In others (we call them “qualitative”), the claim may be that the system has no faults, and our confidence in this claim (sometimes called “probability of perfection” [13]) can be used to estimate long run survival without critical failures [14].

3.1. Qualitative Cases

A natural measure for confidence in the claim of an evidential step is $P(c|e)$; as explained in Section 2.1, we do not use this as a measure for the *weight* of evidence because that must also account for the ability of the evidence to discriminate between the claim and a counterclaim, but once the evidence has been accepted on the basis of its weight, it is reasonable to use $P(c|e)$ as our confidence in its claim.

Next, we need a method to “combine” the confidence measures from the evidentially supported subclaims of a reasoning step to yield a confidence measure for its parent claim, and so on to the top of the tree where we obtain a confidence measure for the top claim. Probability and logic build on completely different foundations and their combination is difficult. Graydon and Holloway [15] examined 12 proposals for using probabilistic methods to quantify confidence in assurance case arguments: 5 based on Bayesian Belief Networks (BBNs), 5 based on Dempster-Shafer or similar forms of evidential reasoning, and 2 using other methods. By perturbing the original authors’ own examples, they showed that all the proposed methods can deliver implausible results.

However, in Assurance 2.0 we have a very simple special case. Ideally, all our reasoning steps are deductive conjunctive implications (i.e., definite clauses), so confidence in a parent claim is given by the product of confidence in the subclaims (provided they are independent). Iterating this

over the whole argument tree, confidence in the top claim is the product of confidence in all the evidentially supported claims. If we have reasoning steps that are not deductive, then it is sound (though often highly conservative) to calculate doubt (i.e., $1 - \text{confidence}$) in the parent claim as no worse than the sum of doubts of the subclaims [16].

Confidence in individual claims may itself be expressed qualitatively (e.g., “high,” “medium,” “low”) and so it will be necessary to develop plausible rules for the “product” of such estimates (e.g., the product of 15 to 25 “highs” yields “medium”). Adjusting a case, or a case template, for different SILs can be accomplished by weakening claims, and by reducing the quantity or quality of evidence demanded; this may in turn allow some subclaims and their supporting argument to be eliminated: e.g., if we replace static analysis by human inspection, we no longer need a subcase for soundness of the static analyzer (but we will need a subcase for reviewer efficacy). Subclaims should not otherwise be removed, for that necessarily makes the case inductive, but we could reduce the threshold at which minor caveats and defeaters are considered mitigated.

In the early stages of system development, the assurance case may be very incomplete yet we would still like to get guidance on areas where attention should be focused. One possibility is to assign exaggeratedly precise assessments for projected confidence in various subclaims (e.g., 73% for this one and 91% for that, and 3% for a nugatory side claim) and then “run the numbers” and do “what if” exercises to learn where the largest impacts reside. The tools supporting these calculations could also take challenges and defeaters into account: a subclaim that has not been challenged would have its confidence reduced, and undefeated defeaters would do the same.

3.2. Quantitative Cases

Next, we consider an example where confidence is an explicit part of the top claim. In addition to confidence, this example also illustrates a more complex development, where defeaters and counterclaims play an important part.

The example is a case based on statistical testing, as used in certain nuclear applications. The idea is that random tests (that follow the “operational profile”) can justify a reliability claim, such as probability of failure on demand (*pdf*) [4].

The evidence offered is a report of the tests performed. The analyst reviews this and integrates it into an assurance case justifying a claim for a certain *pdf* x that is held with confidence c_1 , based on what the analyst considers to be a well known theory of statistical testing. This initial assurance case is shown in Figure 4, where the top claim is a predicate that could be used in a larger case, such as one that combines reasoning about reliability with evidence for correctness.

However, on reflection, the analyst decides that this initial case is subject to significant doubts, because the argument for evidence incorporation does not reach the threshold for indefeasibility: for example, can we be sure the theory of statistical testing was applied correctly? In a fully

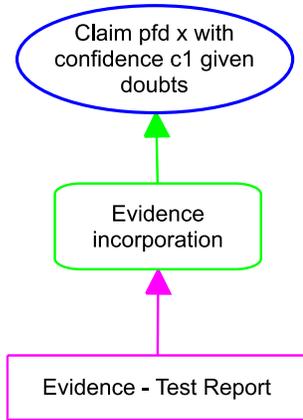


Figure 4. Initial Case for Statistical Testing, with Doubt Annotation

tool-supported environment, there would be ways to indicate this potential defeater but, for the text-based description used here, the analyst simply marks the claim as one with doubts.

The root problem is that the evidence incorporation step combines both the extraction of facts from the test report, and their analysis and interpretation with respect to a model of statistical testing.

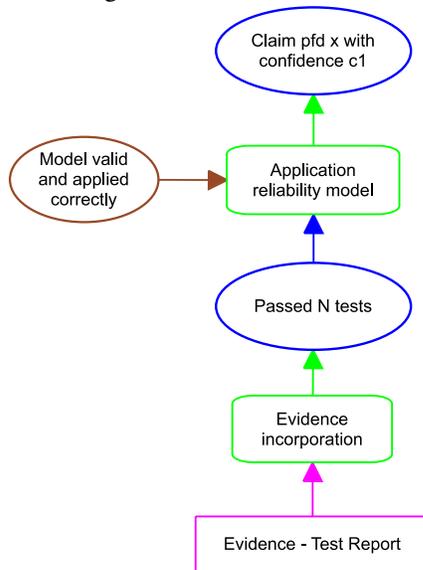


Figure 5. Separation of Facts from Test Report and Inference of Reliability

Consequently, in Figure 5 these two aspects are separated: evidence incorporation extracts the purported facts, and a reasoning step provides the argument that these justify the top claim, with a side claim (that will eventually need to be justified by evidence) to support the validity and correct application of the statistical testing and reliability model.

In constructing a justification for the reasoning step, consideration of this side claim will force realization that the supporting claim “Passed N tests” needs to have a more precise interpretation: namely, that the tests demonstrated N failure-free demands in succession, and that no other failures were observed. Thus, this claim should be changed to “ N successive failure-free demands and no other failures.” If the

lower, evidential step can support this claim (as opposed to a weaker claim where some failures may have been observed) then we can retain Figure 5 as our assurance case, but with the claim “Passed N tests” replaced by the more precise form, as shown in Figure 6.

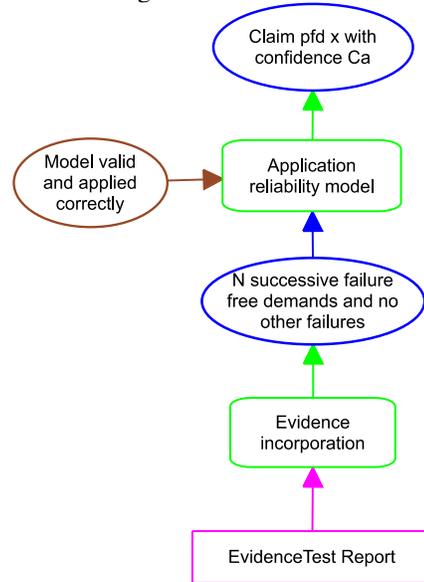


Figure 6. More Precise Claims

Consideration of the side claim also forces realization that confidence in the top claim is with respect to *aleatoric uncertainty*³ (based on the extent of testing) and this is reflected in the revised top claim where C_a replaces C_1 .⁴

Consideration of indefeasibility and side claims suggested improvements in the case; we now look at defeaters.

3.2.1. Defeaters. Further reflection, or challenging peer review, might ask how do we know that the tests were performed correctly, and that issues such as correctness of the test oracle were addressed appropriately? The analyst recognizes that these are legitimate defeaters and the case needs to be strengthened by including a subcase similar to that previously illustrated in Figure 3.

One approach would be to construct a new assurance case in which Figure 6 is a subcase dealing with reliability and confidence, and an elaborated version of Figure 3 is a subcase dealing with soundness of the overall test procedure.

A slight variant, which is appropriate because the top claim explicitly states the confidence associated with the *pdf* x , is to interpret Figure 6 as a subcase dealing with *aleatoric* uncertainty and an elaborated Figure 3 as a subcase dealing with *epistemic* uncertainties. This approach is shown in Figure 7. Note that this and subsequent examples are not complete cases: some claims lack supporting evidence.

3. Aleatoric (or aleatory) uncertainty is uncertainty *in* the world: if I toss a fair coin 100 times, the number of heads is subject to aleatoric uncertainty; epistemic uncertainty is uncertainty *about* the world: if I give you a coin and invite you to toss it 100 times, there is additional uncertainty about the number of heads because you do not know if the coin is fair or not.

4. C_1 and C_a are numbers, but they are annotated with descriptions of their interpretation and it is these that change.

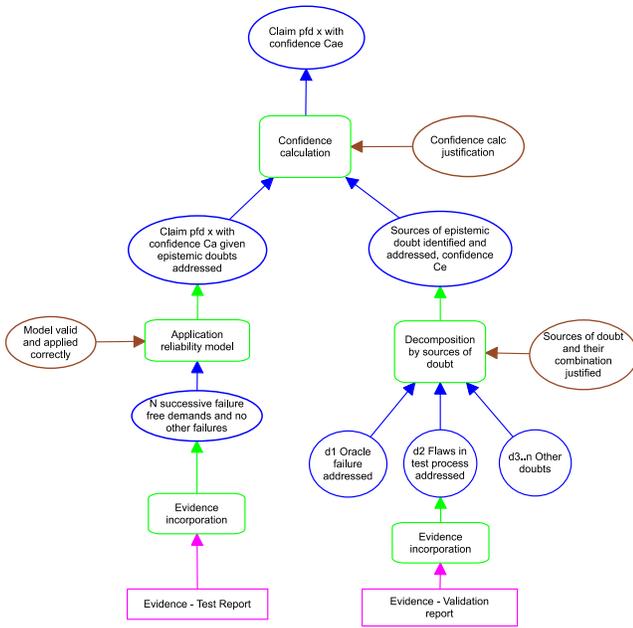


Figure 7. Showing Defeaters Have Been Incorporated

Figure 7 provides a pattern in which we separate reasoning about aleatoric doubts (left-hand leg) from that about epistemic doubts (right-hand leg). However, we may sometimes need to reason about aleatoric and epistemic aspects within the same framework, as when we wish to model their interactions and dependencies. If we were able to provide a quantified judgment of our confidence in the soundness of the oracle and the test process in the form of conditional probability distributions, then we could combine them in a BBN model that does this, as illustrated in Figure 8.

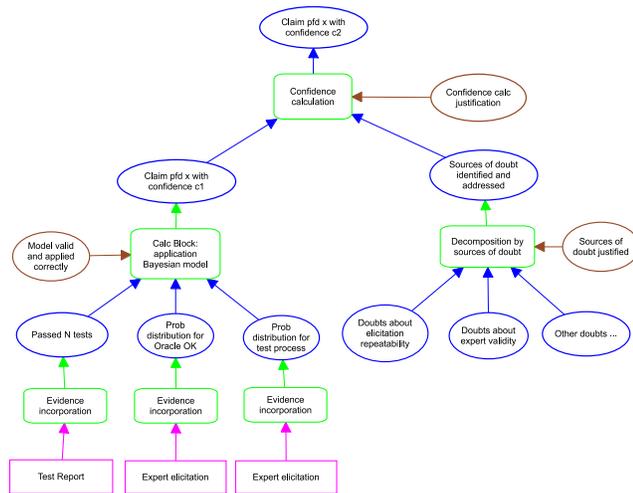


Figure 8. Incorporating BBN Modeling

Here, the left-hand leg uses Bayesian reasoning to provide a probability distribution for the property of interest, and from that derives a confidence figure in the claimed *pdf* x . There is a new side claim that requires justification for the application and validity of the BBN model. The right

hand leg deals with defeaters to the BBN approach; it has identified two (represented as negated claims): validity of elicitation of the probability distributions from experts, and its repeatability. Furthermore, we have a side claim asserting these are the *only* sources of doubt. We are not sure they are, so this argument step is inductive; we choose to represent this by adding a negatory third “other doubts” claim.

These defeaters are formidable: it is seldom credible that we can derive full conditional distributions as needed here. If we can, then the benefit is that the confidence calculated by the left hand leg may be much greater than can be supported by weaker assumptions and conservative calculations as in Section 3.1. An intermediate position in the tradeoff between confidence in the claim and doubts about assumptions is to reduce criticality of the claim and increase confidence in its reduced form: if we are 90% confident that a subcase establishes SIL3 then, with some additional modeling and assumptions, we might become 99% confident that it is better than SIL2, and this could be sufficient to argue that it meets the evidential threshold. If challenged to deal with the remaining doubt, we could use a chain of confidence [17] that combines a firm judgment about the 99% with a conservative judgment about the other 1%.

3.2.2. Countercases. Another way of identifying defeaters or sources of doubt is to develop an explicit countercase that aims to refute the claim under consideration. This task could be assigned to a different team, which, given its different viewpoint, might generate challenging and unexpected defeaters for the base case. There is some tension here: a totally independent countercase might have an argument structure completely different to the base case, and thereby generate irrelevant defeaters.

However, there seems to be a useful transformation from case to a parallel countercase (and vice-versa): mitigated defeaters become claims and claims become a source of defeaters. Thus, the assurance case pattern in Figure 7 is transformed into the countercase shown in Figure 9. (The left and right hand legs are reversed because we draw claims on the left and doubts on the right.)

4. Conclusion

We have described and illustrated Assurance 2.0, whose purpose is to support the assurance challenges posed by recent developments in system design and deployment, and to provide a framework in which assurance can become an enabler of innovation. Assurance 2.0 retains the structure of Assurance Cases and can build on much recent and current research and tooling. Where it differs is in stressing rigor in assessment of the evidence and reasoning employed, and a focus on challenges to confirmation bias through use of confirmation measures, counterclaims and countercases.

Assurance cases have served traditional safety-critical systems well [18], but we have observed them floundering when confronted by radically new challenges such as autonomous systems driven by machine learning and AI, by applications with a security focus, and with new stakeholders such as the AI community. Assurance 2.0 renews the original focus of assurance by asking for a natural language

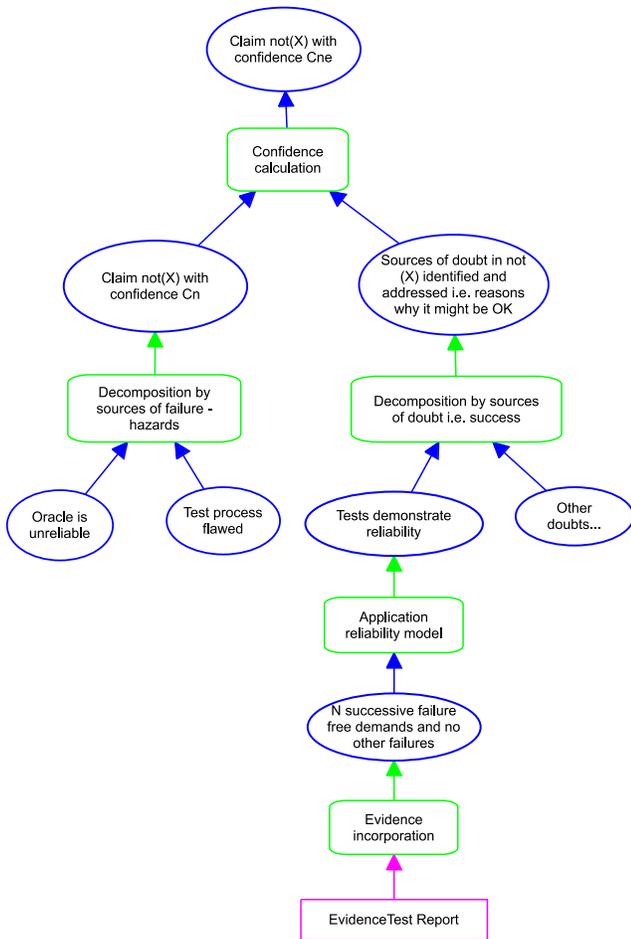


Figure 9. Systematically Derived Countercase

explanation why the proposed system satisfies the properties claimed for it, while balancing this with systematic methods for identifying defeaters. A completed Assurance 2.0 Case attests to the relevance and strength of its evidence and the deductive validity of its reasoning, and also records the defeaters to which it has responded, thereby establishing not merely its plausibility but its indefeasibility.

Ideas underlying Assurance 2.0 have been used with some success in training several groups of engineers and managers and applied in research projects with regulators and industry. A significant element in successful deployment, only briefly mentioned here, is use of a small library of “pre-validated” argument steps (called “blocks”) [3] that reduce the bewildering choice in free-form arguments.

For the future, we hope to see application of these ideas to significant modern systems, supported by training across a wide range of disciplines and the development of constructive tool support. The formal nature of the reasoning and evidential analysis that underlies Assurance 2.0 should enable productive interaction with tools for logical and probabilistic reasoning and formal argumentation, together with novel automation in the search for defeaters, the construction of cases and countercases, and the management and representation of dialectical examination. We plan to prototype

and evaluate the approach in industrial applications and research projects including the DARPA ARCOS program.

Acknowledgments. We thank colleagues at Adelard, City, and SRI for many stimulating discussions on these topics.

This work was funded by Adelard and by SRI.

References

- [1] *ASCAD: Adelard Safety Case Development Manual*, Adelard LLP, London, UK, 1998.
- [2] S. P. Wilson, J. A. McDermid, C. Pygott, and D. J. Tombs, “Assessing complex computer based systems using the goal structuring notation,” in *2nd IEEE International Conference on the Engineering of Complex Computer Systems (ICECCS)*. Montreal, Canada: IEEE Computer Society, Oct. 1996, pp. 498–505.
- [3] R. Bloomfield and K. Netkachova, “Building blocks for assurance cases,” in *ASSURE: 2nd Intl. Workshop on Assurance Cases for Software-Intensive Systems*. Naples, Italy: IEEE Intl. Symp. on Software Reliability Engineering, Nov. 2014, pp. 186–191.
- [4] *Dependability Assessment of Software for Safety Instrumentation and Control Systems at Nuclear Power Plants*, International Atomic Energy Agency, 2018, nuclear Energy Series, NP-T-3.27.
- [5] I. J. Good, “Weight of evidence: A brief survey,” in *Bayesian Statistics 2: Proceedings of the Second Valencia International Meeting*, J. Bernardo *et al.*, Eds., Valencia, Spain, Sep. 1983, pp. 249–270.
- [6] L. Bovens and S. Hartmann, *Bayesian Epistemology*. Oxford University Press, 2003.
- [7] K. Tentori, V. Crupi, N. Bonini, and D. Osherson, “Comparison of confirmation measures,” *Cognition*, vol. 103, pp. 107–119, 2007.
- [8] J. B. Goodenough, C. B. Weinstock, and A. Z. Klein, “Eliminative induction: A basis for arguing system confidence,” in *Proceedings International Conference on Software Engineering, New Ideas and Emerging Results*. San Francisco, CA: IEEE Computer Society, May 2013, pp. 1161–1164.
- [9] N. Leveson, “The use of safety cases in certification and regulation,” *Journal of System Safety*, vol. 47, no. 6, pp. 1–5, 2011.
- [10] J. Rushby, “On the interpretation of assurance case arguments,” in *New Frontiers in AI, Revised Selected Papers*, Springer LNAI, vol. 10091. Kanagawa, Japan: Nov. 2015, pp. 331–347.
- [11] —, “The indefeasibility criterion for assurance cases,” in *Shonan Workshop on Implicit and Explicit Semantics Integration in Proof Based Developments of Discrete Systems*, Kanagawa, Japan, Nov. 2016, postproceedings to be published by Springer LNCS in 2020.
- [12] L. Groarke, “Deductivism within pragma-dialectics,” *Argumentation*, vol. 13, no. 1, pp. 1–16, 1999.
- [13] B. Littlewood and J. Rushby, “Reasoning about the reliability of diverse two-channel systems in which one channel is “possibly perfect”,” *IEEE Transactions on Software Engineering*, vol. 38, no. 5, pp. 1178–1194, Sep./Oct. 2012.
- [14] L. Strigini and A. Povyakalo, “Software fault-freeness and reliability predictions,” in *SAFECOMP 2013: Proceedings 32nd International Conference on Computer Safety, Reliability, and Security*, Springer LNCS, vol. 8153. Toulouse, France: Sep. 2013, pp. 106–117.
- [15] P. J. Graydon and C. M. Holloway, “An investigation of proposed techniques for quantifying confidence in assurance arguments,” *Safety Science*, vol. 92, pp. 53–65, Feb. 2017.
- [16] E. W. Adams, *A Primer of Probability Logic*. Center for the Study of Language and Information (CSLI), Stanford University, 1998.
- [17] 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.
- [18] D. Rinehart, J. Knight, and J. Rowanhill, “Current practices in constructing and evaluating assurance cases with applications to aviation,” NASA Langley Research Center, NASA Contractor Report NASA/CR-2015-218678, Jan. 2015,