# ASSURANCE 2.0: A MANIFESTO

THE DEVELOPMENT AND APPLICATION OF ASSURANCE 2.0

Prof Robin E Bloomfield FREng
Adelard LLP and City, University of London
reb@adelard.com

Joint paper with John Rushby, SRI

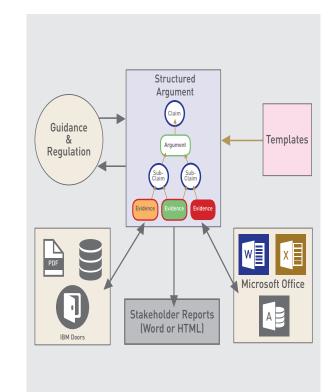Presentation to SSS'21.
Feb 10th 2021

PT/908/180001/9

# ADELARD

- Adelard is a specialized, influential product and services company working on safety, security and resilience

- Wide-ranging experience of assessing computer-based systems and components

- Work across different industrial sectors, including nuclear, transport, defence, financial, medical
  - Policy, methodology, technology
  - Product for managing safety and assurance cases (ASCE)

- Consultants PhD level, international team
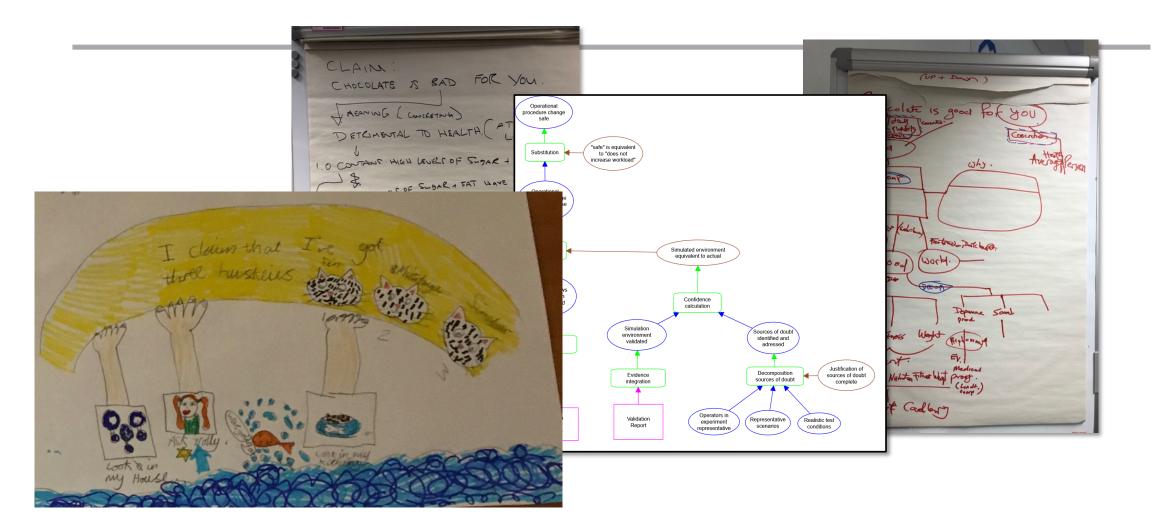
ASCE - in the wider environment

# OUTLINE

- Motivation
  - Briefly, why is Assurance 2.0 needed

- Summary of Assurance 2.0
  - Joint work with John Rushby, SRI

- Some application experience
  - Templates and guidance for Autonomous systems
  - Tool support
  - Industry courses

- Conclusions – from manifesto to methodology

# WHAT DOES GOOD LOOK LIKE?

# DRIVERS FOR CHANGE

- Trustworthy systems expensive and often slow to produce
  - And still have failures

- Assurance is essential – gaining confidence in the system
  - Essential for legal, reputational, market, ethical, commercial reasons
  - Can be slow to produce, slow to change

- Innovation challenges
  - New lifecycles, new technology
  - Higher tempo, varied supply chains. increased threats

- Address existing and emerging requirements for safety and assurance arguments
  - ISO26262, PAS11281, UL4600, EU Pegasus project, Safety First For Automated Driving, UK Regulation for the Fourth Industrial Revolution White Paper

# DRIVERS FOR NEW APPROACH

- Challenge from broadening approach to security and engineering justifications
  - The "non safety case" world using the approach
  - Long term study CAE adoption and CAE role in supporting innovation

- Commoditisation of risk assessment, loss of mindset
  - UK NCSC withdrawal of risk assessment guidance IS1 and IS2
  - https://www.ncsc.gov.uk/guidance/critical-appraisal-risk-methods-and-frameworks

- Challenge of
  - autonomous systems and those using AI/ML
  - automated certification

- Evolution of research on argumentation and assurance

- Overall need for
  - understanding, explanation, challenge, and learning

# ASSURANCE 2.0

- Our idea is to make assurance an enabler for innovation, not a brake

- Paradoxically, we think we can achieve this by making it more rigorous
  - Keep structure of traditional assurance cases
  - Strengthen focus on evidence and reasoning
  - Bring assurance thinking  forward within life-cycle
    – makes it clear what must be done and makes you do it earlier
  - Support assurance with known best practices
    – reduce the bewildering choice of free form cases with "pre-validated" blocks or templates
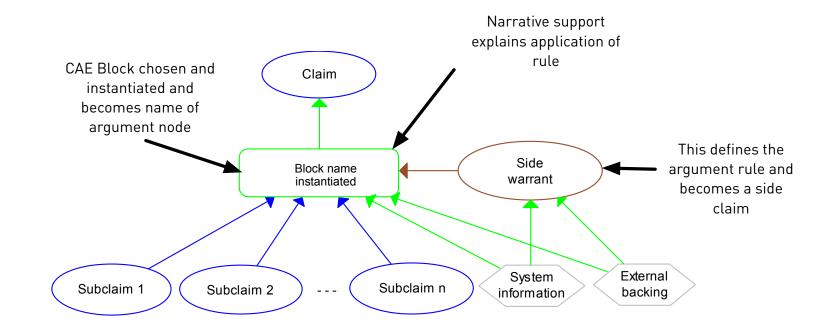
# ASSURANCE 2.0 - MANIFESTO

- Making explicit inference rules and the separation of inductive and deductive reasoning.
  - empirically based CAE Blocks provides a mechanism for separating inductive and deductive aspects of the reasoning. *Natural language deductivism*. **(NLD)**

- Explicit use of doubts and defeaters
  - both undercutting and rebuttal, that confidence an integral part of the justification
  - indefeasibility criterion

- Focus on evidence integration, addressing both the relevance and provenance of evidence.
  - evidential threshold, in which a claim can be reasoned about deductively might be used when considering the role of automated reasoning

- Confirmation theory to evaluate the strength of evidence and arguments.

- Explicit approach to reduce bias by the use of counter-cases and confirmation theory.

- Recognition of importance of both mindset and methodology

# CAE BUILDING BLOCKS - NLD

- Well defined argument fragments, empirically based, but rigorously defined, supporting reasoning both deductive and inductive

- Fragment that support a combined graphical and narrative approach

# DEDUCTIVE AND INDUCTIVE ARGUMENTS

- For valid deductive arguments the premises *logically entail* the conclusion, where the entailment means that the truth of the premises provides a *guarantee* of the truth of the conclusion

- An inductive logic is a system of evidential support that extends deductive logic to less-than-certain inferences

- In a good inductive argument the premises should provide some *degree of support* for the conclusion, where such support means that the truth of the premises indicates with some *degree of strength* that the conclusion is true.
  - acceptability, relevance and sufficiency
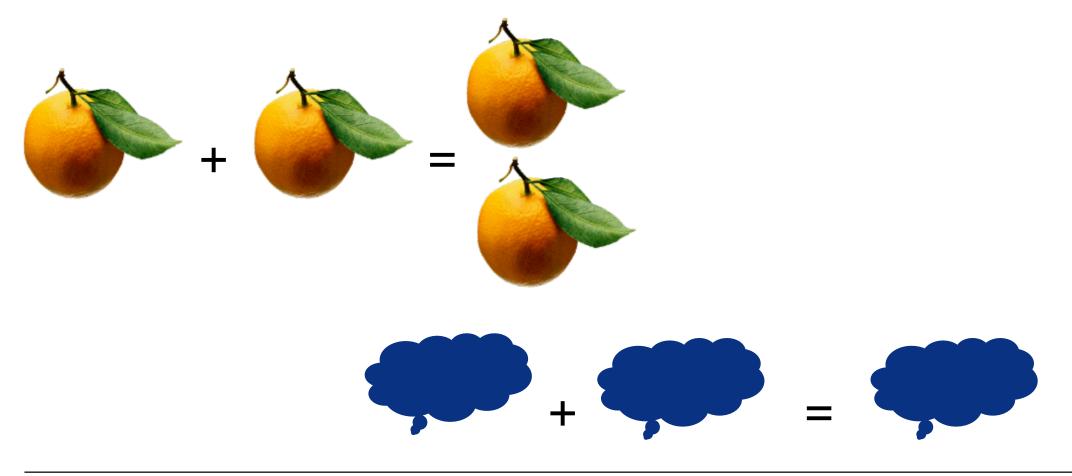
Adapted from *https://plato.stanford.edu/index.html*

# DEDUCTIVE AND INDUCTIVE ARGUMENTS

- For valid deductive arguments the premises *logically entail* the conclusion, where the entailment means that the truth of the premises provides a *guarantee* of the truth of the conclusion

- An inductive logic is a system of evidential support that extends deductive logic to less-than-certain inferences

- In a good inductive argument the premises should provide some *degree of support* for the conclusion, where such support means that the truth of the premises indicates with some *degree of strength* that the conclusion is true.
  - acceptability, relevance and sufficiency


Adapted from *https://plato.stanford.edu/index.html*

# EXAMPLE

# DEDUCTIVE AND INDUCTIVE ARGUMENTS –WHY SEPARATE OUT?

## Science of security – importance of deductive/inductive split

"We now detail security research failures to adopt accepted lessons from the history and philosophy of science.

*A. Failure to observe inductive-deductive split*

Despite broad consensus in the scientific community, in Security there is repeated failure to respect the separation of inductive and deductive statements "

## SoK: Science, Security, and the Elusive Goal of Security as a Scientific Pursuit

Cormac Herley
Microsoft Research, Redmond, WA, USA
cormac@microsoft.com

P.C. van Oorschot
Carleton University, Ottawa, ON, Canada
paulv@scs.carleton.ca

# DEDUCTIVE AND INDUCTIVE ARGUMENTS – WHY SEPARATE OUT?

- Side claim provides a mechanism for factoring
  - Inductive argument-A = Deductive argument + Inductive argument-B
  - Where deductive gives some leverage e.g. analysis, tool support
  - Inductive argument-B is easier to show than Inductive argument-A (then we have made progress!

- Examples
  - Application of deductive models
    - Infer properties
      - Testing evidence -> reliability
      - Abstract interpretation -> run time errors
    - Architecture
      - Property distributes over components (e.g. confidentiality)
    - System properties
      - Fire, flood, earthquakes
  - Each time need to address validity of model and proper application via a side claim

# FIVE CAE BUILDING BLOCKS

- Well defined argument fragments
  - Empirically based, but rigorously defined
  - Supporting both deductive and inductive reasoning

- Fragments support a combined graphical and narrative approach

**Decomposition**
Partition some aspect of the claim
Divide and conquer

**Substitution**
Refine a claim about an object into claim about an equivalent object

**Evidence incorporation**
Evidence supports the claim
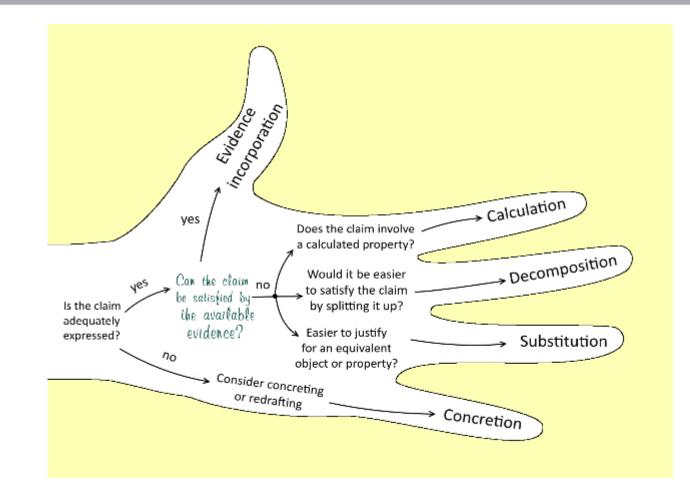Emphasis on direct support

**Concretion**
Some aspect of the claim is given a more precise definition

**Calculation or proof**
Some value of the claim can be computed or proved

# 'HELPING HAND' – GUIDANCE ON SELECTING BLOCKS

# DEFEATERS – EXPLICITLY DEALING WITH SOURCES OF DOUBT

- One concept used to address stopping rules and over-confidence is "defeaters". The concept of defeaters is used to articulate reasons why a claim might **not** be supported.

- Two kinds of defeaters:
  - Rebutting defeaters, which are reasons for believing the negation of the conclusion, and
  - Undercutting defeaters, which provide a reason for doubting that claim.

- Identification and mitigation of defeaters are foundational to assurance
  - Think of as hazard analysis applied to arguments

- In CAE
  - Rebutting defeaters can be addressed with negated subclaims
  - Undercutting defeaters can be addressed by explicitly showing them in the CAE structure
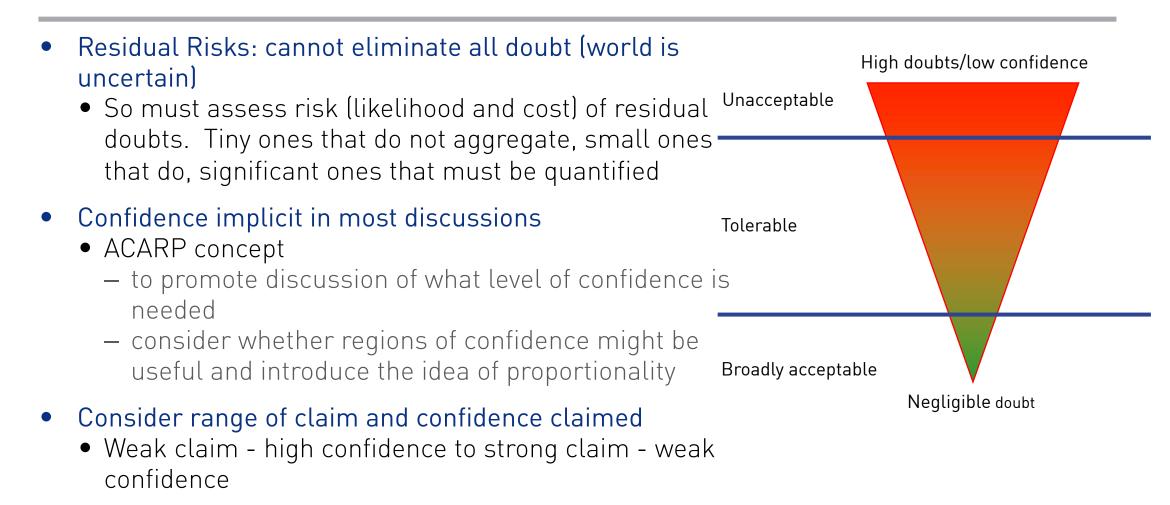
# CONFIDENCE

- The purpose of an assurance case is to assist in making, justifying, and communicating the *decision* to deploy a system or service in a given context

- Top level requirement is that the justification should be indefeasible.
    - Meaning it is so well supported and all credible doubts & objections have been so thoroughly considered & countered
    - That no credible doubts remain that could change the decision

- Confidence is strength of our belief that case is indefeasible

- We do not think is can be reduced to some single assessment of the case

- Instead, we identify three perspectives, and assessments and measures within those
    - Assessment of confidence based on all three perspectives

# THREE PERSPECTIVES ON CONFIDENCE

- Positive: extent to which case makes positive case to justify belief in its claims
  - Soundness: logical criterion using Natural Language Deductivism (NLD)
    - Based on weight of evidence, deductive reasoning
  - Probabilistic valuation: probabilistic criterion using Bayesian framework (CBI, BBN)
    - This is what many others mean by confidence: usually flawed (Graydon &Holloway)
    - We require case to be sound, only 5 argument blocks: avoids flaws

- Negative: extent to which doubts have been investigated and addressed
  - Doubts are vague, become defeaters when sharpened, recorded in the case
    - Together with justification for their own defeat (eliminative argumentation)
    - Use systematic methods to find credible defeaters (cf. hazard analysis)
    - May also be possible to invert positive perspective on counterclaims

- Residual Risks: cannot eliminate all doubt (world is uncertain)
  - So must assess risk (likelihood and cost) posed by residual doubts. Tiny ones that do not aggregate, small ones that do, Significant ones that must be quantified

# ACARP - ANALOGY WITH ALARP

- Residual Risks: cannot eliminate all doubt (world is uncertain)
  - So must assess risk (likelihood and cost) of residual doubts. Tiny ones that do not aggregate, small ones that do, significant ones that must be quantified

- Confidence implicit in most discussions
  - ACARP concept
    - to promote discussion of what level of confidence is needed
    - consider whether regions of confidence might be useful and introduce the idea of proportionality

- Consider range of claim and confidence claimed
  - Weak claim - high confidence to strong claim - weak confidence

High doubts/low confidence

Unacceptable

Tolerable

Broadly acceptable

Negligible doubt

# WEIGHT OF EVIDENCE – STRENGTH OF CLAIM

- It's not enough for evidence to support a claim

- It must distinguish a claim from its negation

- Confirmation measures do this: e.g., Kemeny-Oppenheim
  - Goes back to work of Good and Turing in WW2 codebreaking

- These force you to look at counterclaims
  - These are potential defeaters

- Can do this informally/qualitatively, don't need numerical probabilities

$$confirmation\_ratio(Evidence, Claim)$$

$$= \frac{\Pr(Evidence \mid Claim\_true) - \Pr(Evidence \mid Claim\_false)}{\Pr(Evidence \mid Claim\_true) + \Pr(Evidence \mid Claim\_false)}$$

Probability that you see the evidence if the claim is true

Probability that you see the evidence if the claim is false

# CONFIRMATION – ROLE OF DIFFERENT EVIDENCE

Probability see evidence if claim true

| | | very unlikely | perhaps | quite probable | very likely |
|---|---|---|---|---|---|
| | | 0.05 | 0.1 | 0.6 | 0.95 |
| very unlikely | 0.05 | 0.00 | 0.33 | 0.85 | 0.90 |
| perhaps | 0.1 | -0.33 | 0.00 | 0.71 | 0.81 |
| quite probable | 0.6 | -0.85 | -0.71 | 0.00 | 0.23 |
| very likely | 0.95 | -0.90 | -0.81 | -0.23 | 0.00 |

Probability see evidence if claim false

$$confirmation\_ratio(Evidence, Claim)$$

$$= \frac{\Pr(Evidence|Claim\_true) - \Pr(Evidence|Claim\_false)}{\Pr(Evidence|Claim\_true) + \Pr(Evidence|Claim\_false)}$$

# CREATING COUNTER CASES

| Group #1 | Group #2 |
|---|---|

- Chocolate is good for you

- Chocolate is bad for you

# SUMMARY – ASSURANCE 2.0 MANIFESTO

- Assurance 2.0 – key components

- Basic Concepts CAE

- CAE Blocks
  - Empirically based
  - Potential for deductive/inductive split

- Defeaters and confidence
  - Indefeasibility and residual rikss

- Evidence
  - Relevance and provenance
  - Confirmation theory and strength of arguments and evidence

- Explicit approach to bias
  - Counter-cases and confirmation theory

# DEVELOPMENT AND APPLICATION – WILL IT WORK?

- Security applications

- Impact on regulation of systems incorporating AI/machine learning

- Developed autonomous system "templates and guidance"

- Tool support
  - building on Adelard ASCE tool within a program on automated certification

- Teaching concepts to professional engineers
  - many disciplines

Theory into practice

DSTL sponsored research

# SAFETY CASE TEMPLATES FOR AUTONOMOUS SYSTEMS

http://arxiv.org/abs/2102.02625

# DEVELOPMENT OF TEMPATES FOR AV

# GENERIC MONITOR GUARD ARCHITECTURE



F3269-17 Standard Practice for Methods to
Safely Bound Flight Behavior of Unmanned
Aircraft Systems Containing Complex Functions,
ASTM International

# DEFEATER WORKSHOP – MONITOR/GUARD ARCHITECTURE

- Colour coded issues and organisations

- Identified issues on-line with international team
  - Briefing
  - Silent brainstorm
  - Collaborative
  - Grouping, sentencing

- Work in progress
  - Still exploring how to capture and present defeaters

## DEFEATERS

- Summary tables – with supporting narrative

| Description | Part of monitor pattern | Possible mitigations |
|---|---|---|
| Operating out of permitted operational envelope not detectable/detected. | Guard/recovery action. | Well-defined operating requirements, testing. Operational restrictions. Make an explicit part of case to detect out of envelope (see Section 7.2.1.1). |
| AI/ML guard functional behaviour not fully verifiable. | Guard. | Restrict design to verifiable ML algorithms in guards. Use reliability rather correctness arguments. |
| AI/ML guard functional behaviour too complex in practice. | Guard. | Simplify guards and place restrictions on operation. |
| Not enough of diversity/independence in sensor and guard. Common cause issues, e.g. due to external common systems GPS or due to sensors finding similar situations difficult. | Architecture level. | Functional diversity – use different type of input data provides some defence. Architectural diversity – different computer system for guards. Justify a level of dependence and use a confidence evaluation that takes this into account. |
| Architecture sensitive to complex failures, e.g. dataflow between sensor | Architecture level. | Adopt appropriate explicit fault models, validate these and engineer |

# TECHNICAL GUIDANCE

- **Confidence measures for ML**
  - Conformal Prediction
  - Inductive Conformal Prediction
  - Attribution-based confidence
  - Learning confidence

- **Performance of ML based components**
  - Performance metrics for binary classifiers
  - Object detection
  - Experimental performance

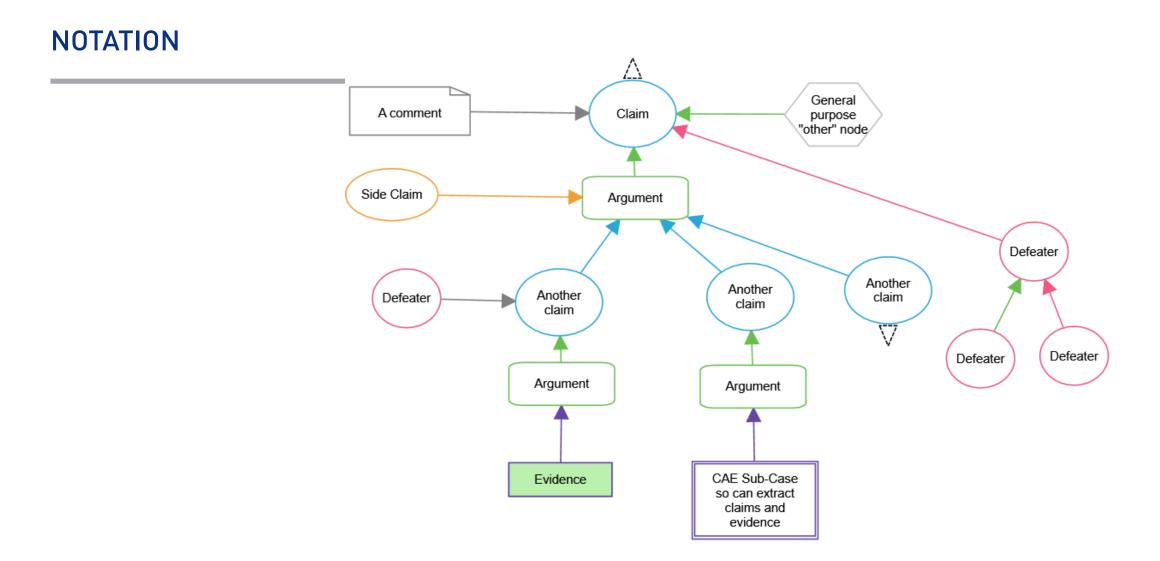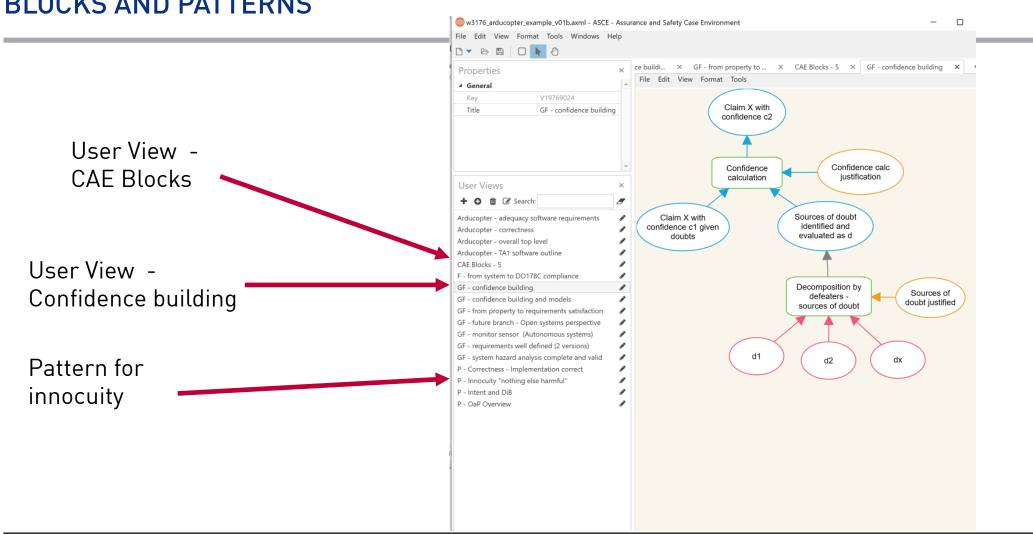| Evidence | Example | Role in case | Example claim |
|---|---|---|---|
| Temporal redundancy | The "Person of Interest" tracker tracked 41% of pedestrians and lost 19% of pedestrians over 20 consecutive frames.<br><br>The traffic light detection system detected all red lights in the test data within 1.6 seconds at a distance of at least 80 metres. | If the sensor output is processed further to produce a model of the world, then the frequency with which each vehicle/pedestrian is detected can support claims about the accuracy of the model.<br><br>Evidence regarding temporal redundancy is particularly relevant in detecting static objects such as traffic lights or a stop sign, which need not be detected every frame, but must be detected within a suitably short timeframe.<br><br>The sensor must also be resilient against single event upsets (if not detected or if falsely detected) to ensure the stability of its outputs. | The pedestrian tracking system identifies 80% of pedestrians which are visible for at least one second[1].<br><br>All red traffic lights are detected from a distance greater than the stopping distance of the vehicle. |
| Additional information (e.g. GPS) | The traffic light detection system correctly identified all traffic lights in the test using predictions from YOLOv3, GPS data and a map of traffic light locations.<br><br>Keeping maps up-to-date used for navigation and locations of static objects of interest (traffic lights, stop signs, junctions) needs to be made in the system is safe in the future branch. | Information such as GPS location can be combined with object detection algorithms to provide better performance for a sensor. A performance claim can be made for this combined system.<br><br>Additional information such as GPS location could also be used as a guard by, e.g. setting a maximum speed if a traffic light is not detected when expected, or geofencing the area in which the AV can operate autonomously. | The addition of a GPS guard reduces false positive traffic light detections by 80%.<br><br>The traffic light detection system correctly identifies 95% of traffic lights in Vitoria with confidence 60%[2].<br><br>The AV only operates autonomously within the city of Vitoria. |

# TOOL SUPPORT

# NOTATION

# BLOCKS AND PATTERNS

User View  -
CAE Blocks

User View  -
Confidence building

Pattern for
innocuity
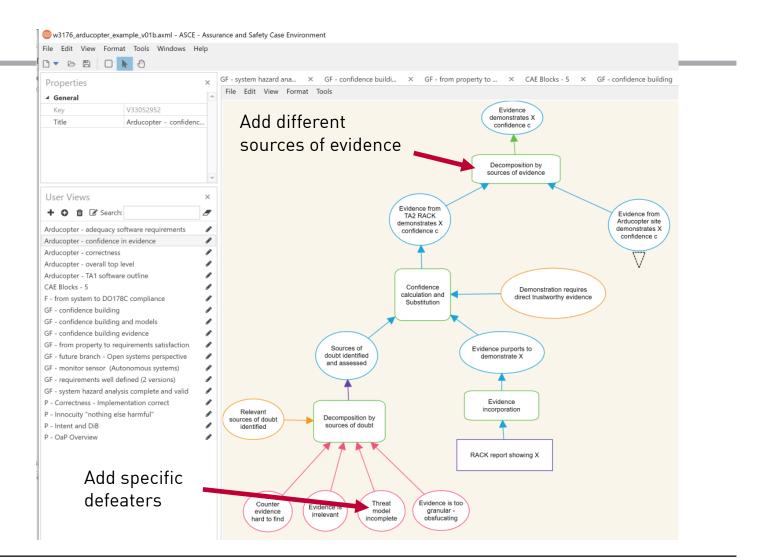
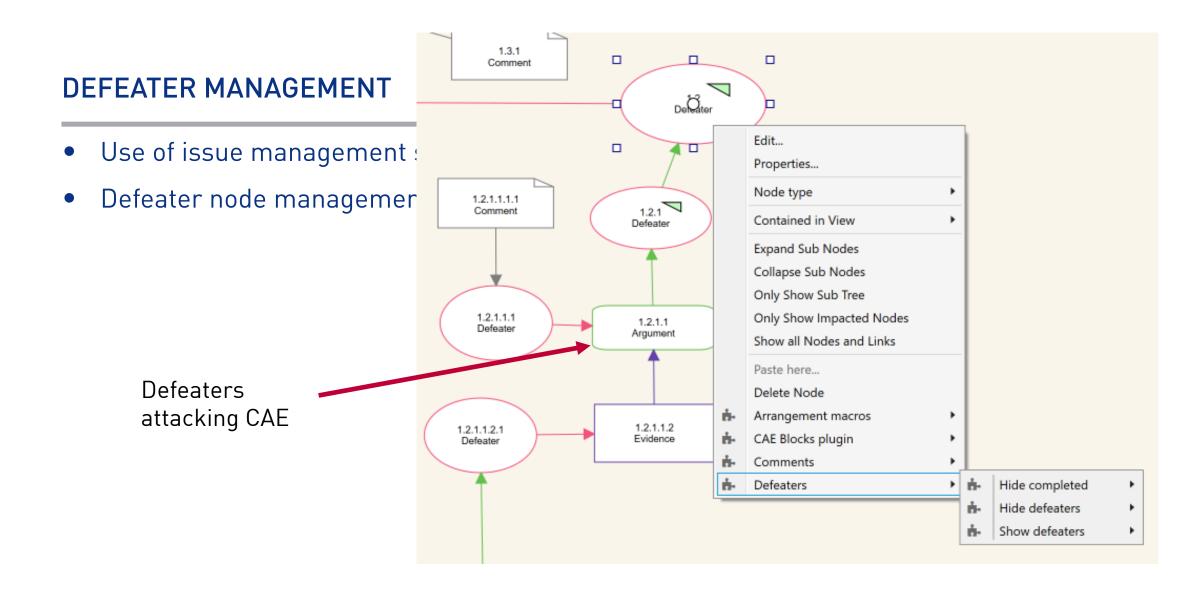# SYNTHESIS

- Evidence Integration + Confidence pattern

- Different sources of evidence
  - Added Decomposition

- Added specific defeaters

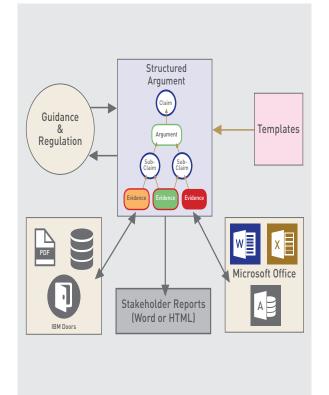# DEFEATER MANAGEMENT

- Use of issue management s
- Defeater node managemer

Defeaters
attacking CAE

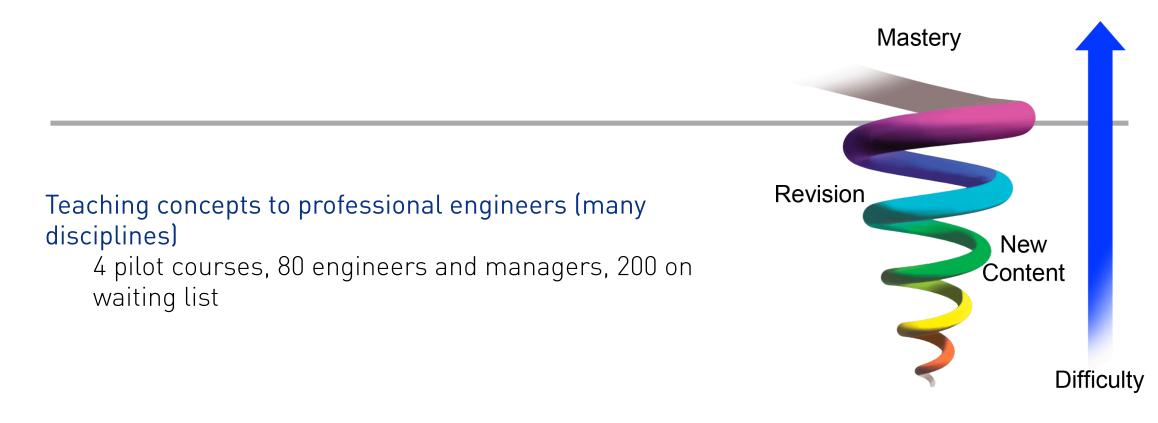# EMBEDDED DEFEATERS

# NEXT STEPS

- Assurance 2.0 support in Adelard ASCE tool
    - Available in new release, March 2021
    - If interested in beta versions please get in touch

- Safety Case Templates for Autonomous Systems
    - Example templates for autonomous systems will be available too based on work for DSTL. Report is

    - http://arxiv.org/abs/2102.02625

**ASCE - in the wider environment**

Teaching concepts to professional engineers (many disciplines)

4 pilot courses, 80 engineers and managers, 200 on waiting list

# APPLICATION – MAJOR HAZARDS SITE

# OUTLINE – ONLINE COURSE

- **Session 1: CAE concepts**
  - Claims, Arguments, Evidence (CAE): concepts and background
  - Inductive and deductive reasoning
  - Application of CAE concepts
  - Introduction to defeaters
  - Short exercise

- **Session 2: Theory into practice**
  - Short exercise
  - The CAE blocks and guidance
  - Discussion of Operations Room example
  - Workshop exercise and discussion

- **Session 3: Learning by doing, workshop exercises and discussion**

- **Session 4: Challenge, review and deployment**
  - Build confidence into the justification
  - Review and challenge
  - Summary

- **Session 5: Wrap up and discussion**
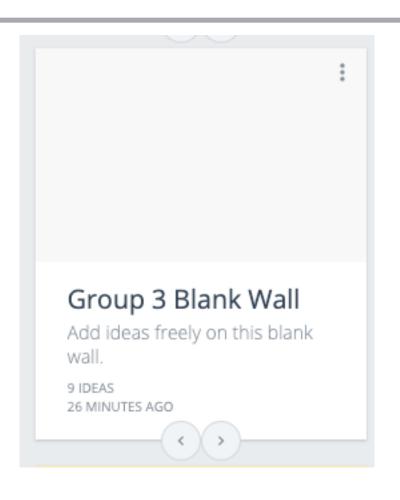  - Putting it all together  and next steps, work projects

# EXERCISES

- Objective is to practice using the CAE Blocks

- Work in groups with a canvas per group

- Stages
  - Decomposition Block example
  - An example of putting the Blocks together
  - Examples of all 5 Blocks

- Add questions and comments to us as you go

- Review

**Group 3 Blank Wall**

Add ideas freely on this blank wall.

9 IDEAS
26 MINUTES AGO

# EXERCISE - DOUBTS AND SIMULATION VALIDATION

- Objective is to express defeaters
  - What might defeat the reasoning that the simulator is valid i.e. sufficiently realistic?
  - "Simulated environment equivalent to actual"

- Work individually

- Add questions and comments to us as you go



Trial Defeater - validation of models/simulators

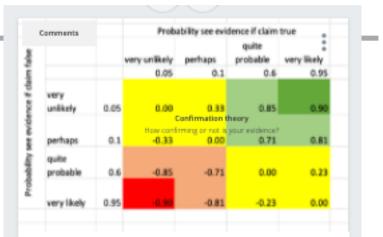To identify and group defeaters so we can improve assurance

2 IDEAS
16 HOURS AGO

## EXERCISE

- In groups discuss examples of claims and evidence asking
  - How likely I am to see the evidence if the claim is true?
  - How likely I am to see the evidence if the claim is false?

- and put on the grid along with any comments



Confirmation theory trial

6 IDEAS
2 DAYS AGO

# APPLICATION IN MAJOR HAZARDOUS SITE – CONCLUSIONS TO DATE

- **Can get ideas across with a day course**
  - Teaching concepts to professional engineers (many disciplines)
  - Often those *without* safety case background find it easier
  - Wide range of responses – struggle, OK, great

- **Follow up application on real projects required**
  - Over several months
  - Surgeries and support

- **Experience and feedback**
  - In progress
  - So far 4 pilot courses, 80 engineers and managers, 200 on waiting list
    - CAE Blocks , defeaters, counter cases ☺
  - Will review and publish experience after ~100 through course

# FROM MANIFESTO TO MATURE METHODOLOGY

- Empirically based CAE Blocks separate inductive and deductive aspects

- Explicit use of doubts and defeaters

- Increased focus on evidence integration, addressing both relevance and provenance

- Confirmation theory to evaluate the strength of evidence and arguments.

- Explicit approach to bias by the use of counter-cases and confirmation theory.

- Recognition of both mindset and methodology

- Publish and apply
  - Different maturity

- Real applications
  - Engineering justifications, safety and security

- Teaching and learning - evaluation
  - >100 industry by April

- Further development of methodology
  - Defeater identification and management
  - Synthesis approaches
  - Confidence and defeaters

- Assurance 2.0 and templates + tools
  - Evaluation and further development

Prof Robin E Bloomfield FREng

Adelard LLP and City, University of London

reb@adelard.com

r.e.bloomfield@city.ac.uk

Joint paper with John Rushby, SRI

Presentation to SSS'21, Feb 10th 2021