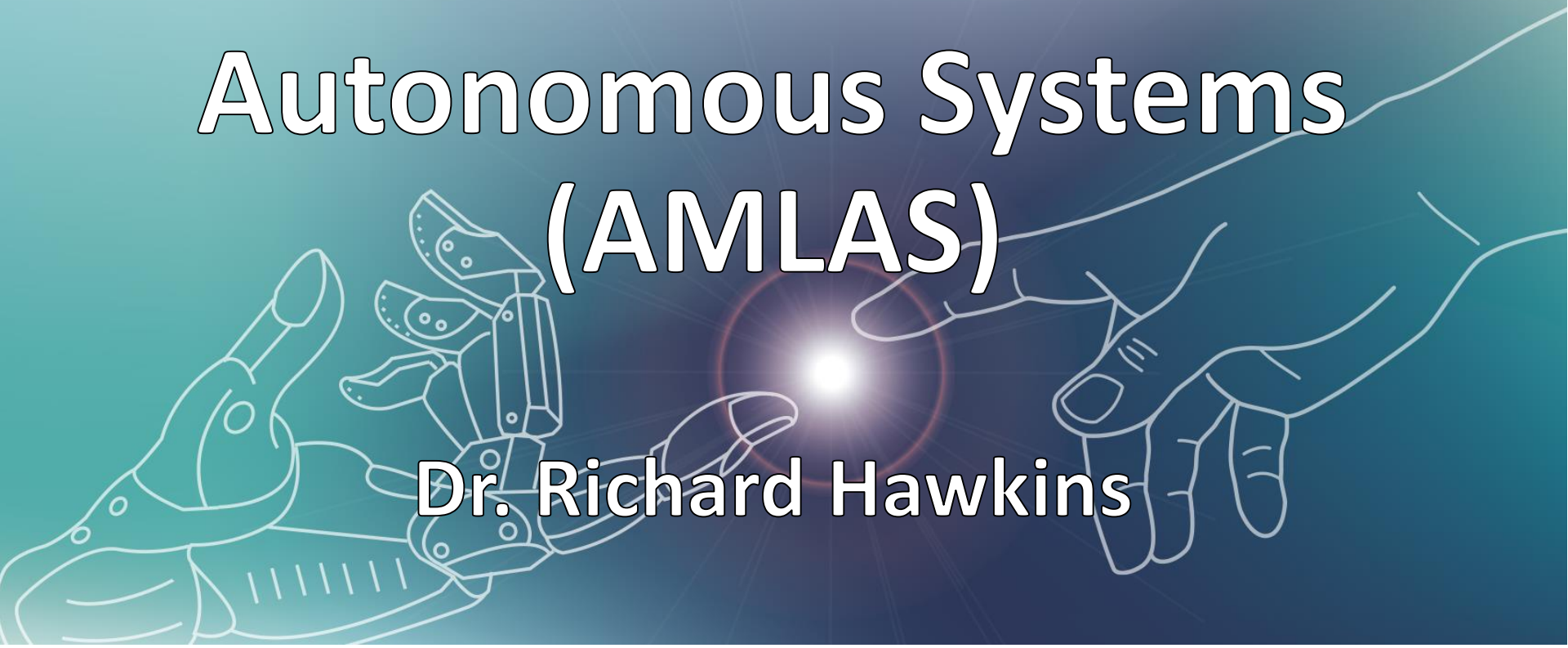


Assurance of ML for Autonomous Systems (AMLAS)

Dr. Richard Hawkins



What is AMLAS?

- AMLAS provides
 - Defined **process**
 - Set of **safety case patterns**
- AMLAS enables
 1. Integration of safety assurance into development of ML components
 2. Generation of evidence base for justifying acceptable safety
- Resulting in structured safety case for ML component

ASSURING AUTONOMY INTERNATIONAL PROGRAMME

Guidance on the Assurance of Machine Learning in Autonomous Systems
(AMLAS)

Richard Hawkins, Colin Paterson, Chiara Picardi, Yan Jia, Radu Calinescu and Ibrahim Habli.

Assuring Autonomy International Programme (AAIP)
University of York

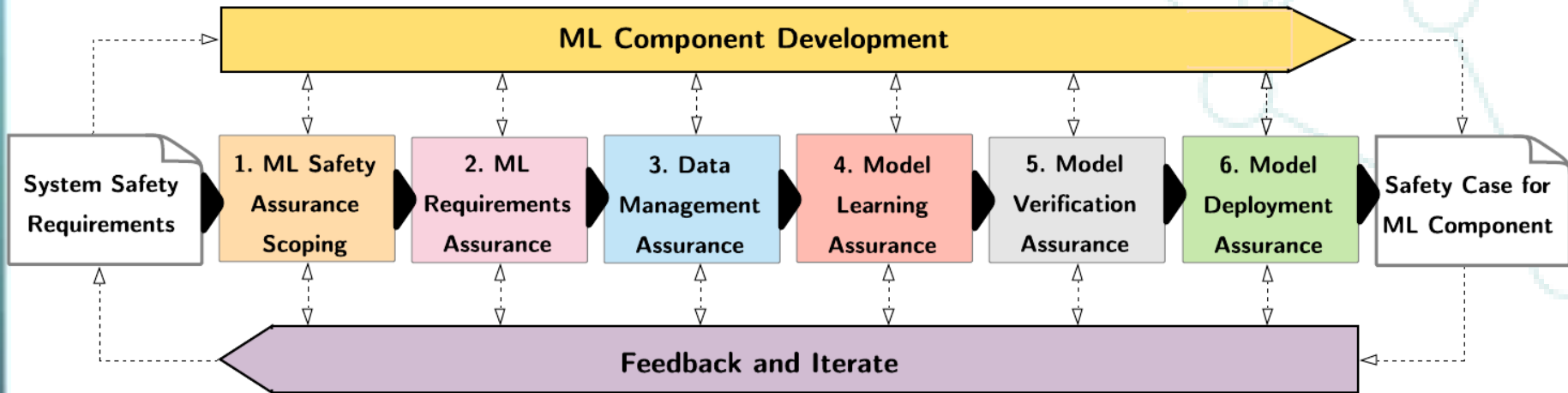
Version 1, February 2021

The material in this document is provided as guidance only. No responsibility for loss occasioned to any person acting or refraining from action as a result of this material or any comments made can be accepted by the authors or The University of York.

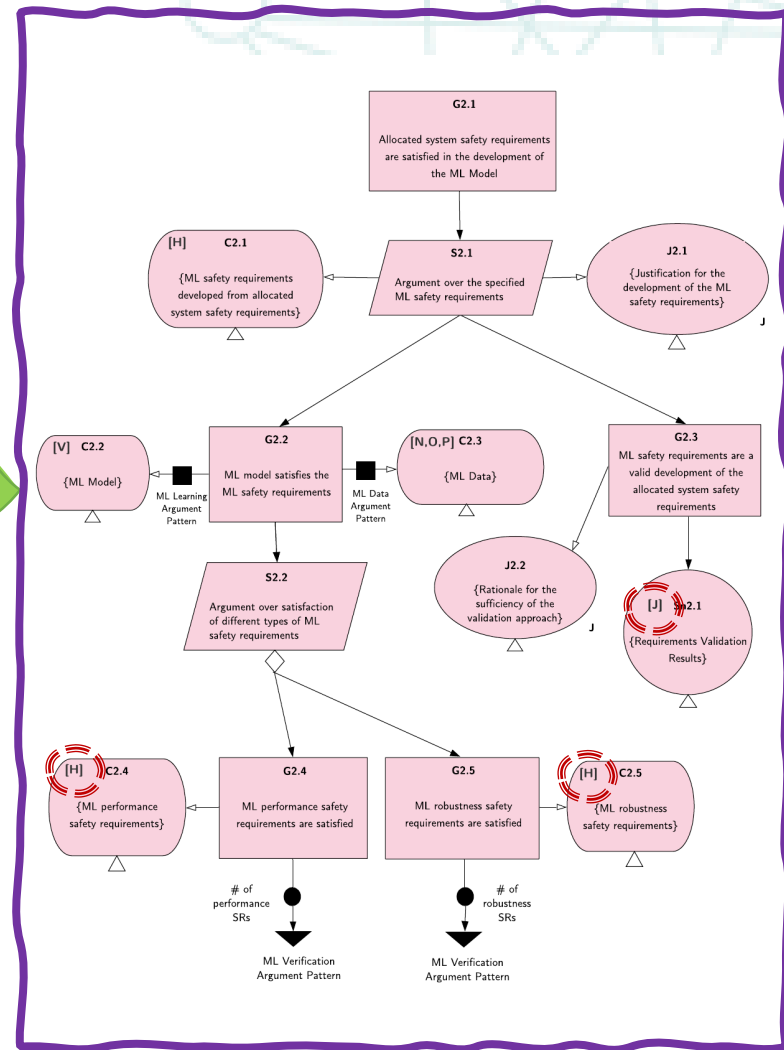
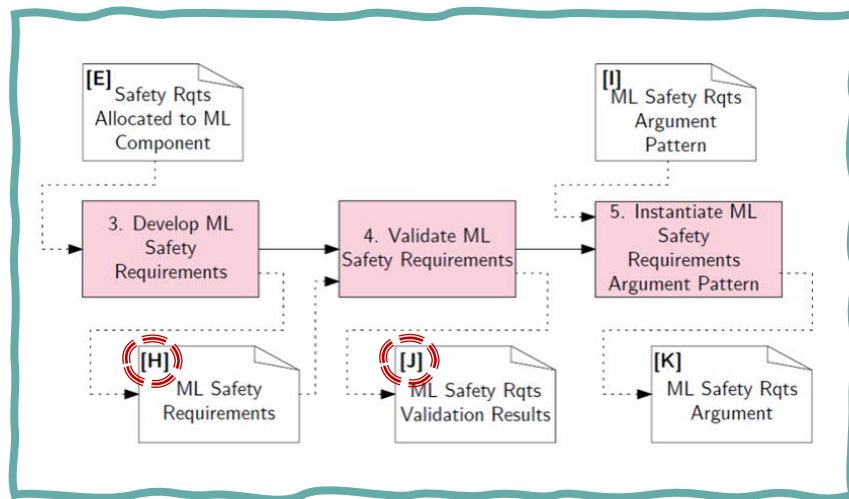
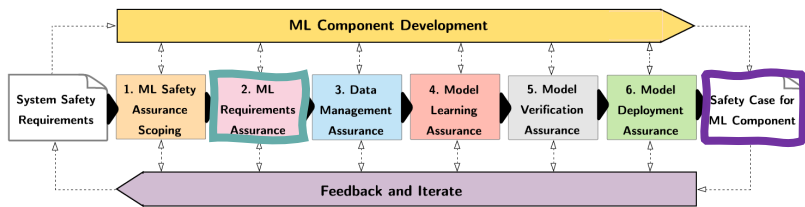
This work is licensed under the Creative Commons Attribution-NoDerivatives 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nd/4.0/> or send a letter to Creative Commons, PO Box 1866, Mountain View, CA 94042, USA. Requests for permission for wider use or dissemination should be made to the authors:-

Contact : firstname.lastname@york.ac.uk.

AMLAS Overview



- For each stage we provide
 - Process description
 - Safety argument pattern



Guidance Structure

Stage 1. ML Safety Assurance Scoping

Objectives

1. Define the scope of the safety assurance process for the ML component.
2. Define the scope of the safety case for the ML component.
3. Create the top-level safety assurance claim and specify the relevant contextual information for the ML safety argument.

Inputs to the Stage

- [A] : System Safety Requirements
- [B] : Description of Operating Environment of System
- [C] : System Description
- [D] : ML Component Description
- [F] : ML Assurance Scoping Argument Pattern

Outputs of the Stage

- [E] : Safety Requirements Allocated to ML Component
- [G] : ML Safety Assurance Scoping Argument

Description of the Stage

As shown in Figure 2¹, this stage consists of two activities that are performed to define the safety assurance scope for an ML component. The artefacts generated from this stage are used to instantiate the ML safety assurance scoping argument pattern as part of Activity 2. An ML component comprises an ML model, e.g. a neural network, that is deployed onto the intended computing platform (i.e. comprising both hardware and software).

Additional guidance on the use of ML for autonomous systems can be found at [9].

Activity 1: Define the Safety Assurance Scope for the ML Component [E]

This activity requires as input the system safety requirements ([A]), descriptions of the system and the operating environment ([B], [C]), and a description of the ML component that is being considered ([D]). These inputs shall be used to determine the safety requirements that are allocated to the ML component.

The safety requirements allocated to the ML component shall be defined to control the risk of the identified contributions of the ML component to system hazards. This shall take account of the defined system architecture and the operating environment. At this stage the requirement is independent of any ML technology or metric but instead reflects the need for the component to perform safely with the system regardless of the technology later deployed. The safety requirements allocated to the ML component generated from this activity shall be explicitly documented ([E]).

¹In the AMLAS process diagrams, rectangles represent activities. Document symbols represent input or output artefacts. Each document symbol has a unique ID (top left) that is used to refer to the artefact in the guidance text or the argument pattern, e.g. [A] is a reference to artefact A.

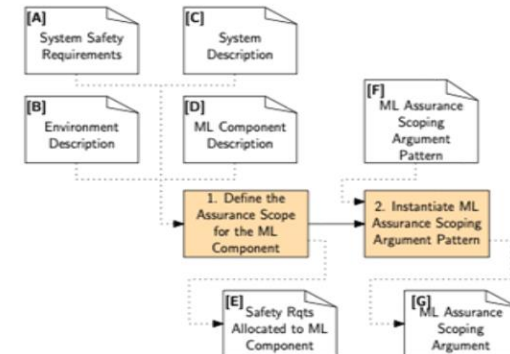


Figure 2: AMLAS ML Assurance Scoping Process

Example 1.

Consider an autonomous driving application in which a subsystem may be required to identify pedestrians at a crossing. A component within the perception pipeline may have a requirement of the form "When Ego is 50 metres from the crossing, the object detection component shall identify pedestrians that are on or close to the crossing in their correct position."

Note 1.

The allocation of safety requirements must consider architectural features such as redundancy when allocating the safety requirements to the ML component. Where redundancy is provided by other, non-machine-learned components, this may reduce the assurance burden on the ML component that should be reflected in the allocated safety requirements.

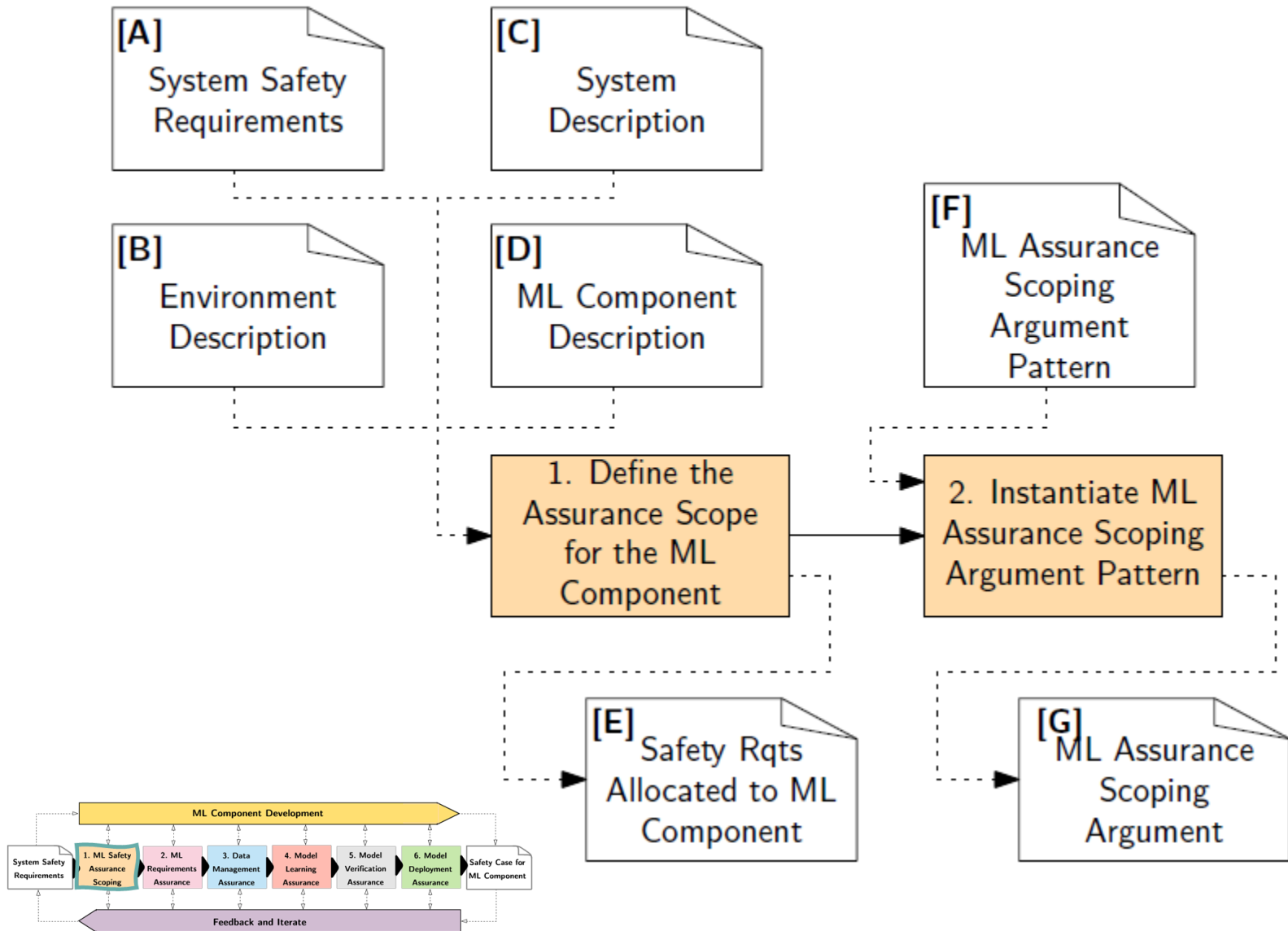
Note 2.

The contribution of the human as part of the broader system should be considered. A human may provide, for example, oversight or fallback in the case of failure of the ML component. These human contributions, and any associated human factors issues, e.g. automation bias [59], should be reflected when allocating safety requirements to the ML component.

Artefact [A]: System Safety Requirements

The safety requirements are generated from the system safety assessment process. Such a process covers hazard identification and risk analysis. Importantly, it shall determine the contribution, i.e. in the form of concrete failure conditions, that the output of the machine learning component makes to potential system hazards. A simplified linear chain of events that links a machine learning failure with a hazard is illustrated in Figure 3.

1. ML Safety Assurance Scoping

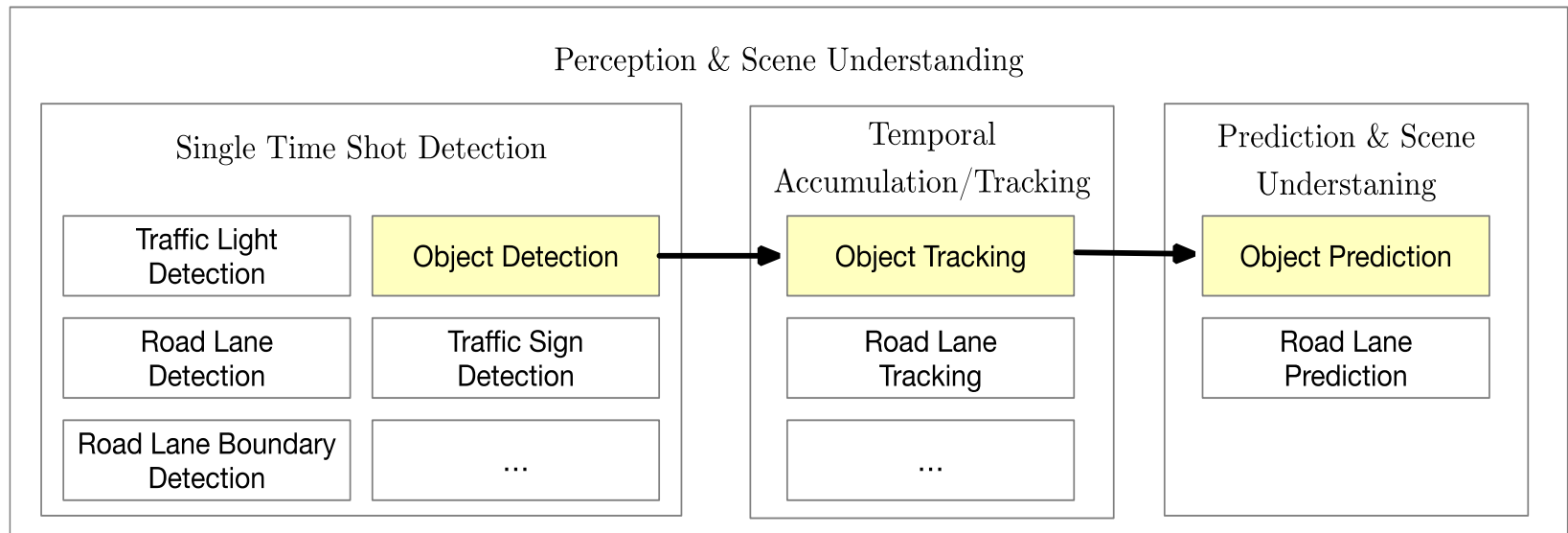
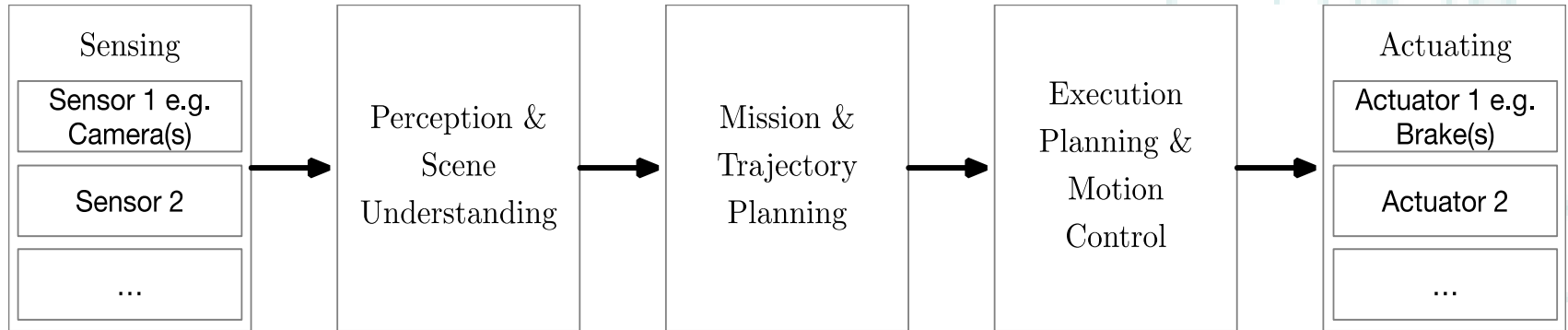


slide)
t (see below)

-
- The image shows a street scene with three pedestrians detected by a system. Each detection is represented by a green bounding box and a set of data points:
- Left Pedestrian:**
 - score: 0.99
 - coord: [50, 36, 820, 38, 100, 217, 743, 95]
 - class: 1
 - category: pedestrian
 - confidence: 0.99
 - Center Pedestrian:**
 - score: 0.99
 - coord: [148, 59, 633, 60, 574, 697, 823, 60]
 - class: 1
 - category: pedestrian
 - confidence: 0.99
 - Right Pedestrian:**
 - score: 0.99
 - coord: [612, 615, 659, 616, 659, 616, 715, 621]
 - class: 1
 - category: pedestrian
 - confidence: 0.99
- Additional data for the rightmost pedestrian group:
- score: 0.99
 - coord: [612, 615, 659, 616, 659, 616, 715, 621]
 - class: 1
 - category: pedestrian
 - confidence: 0.99



Model Integration

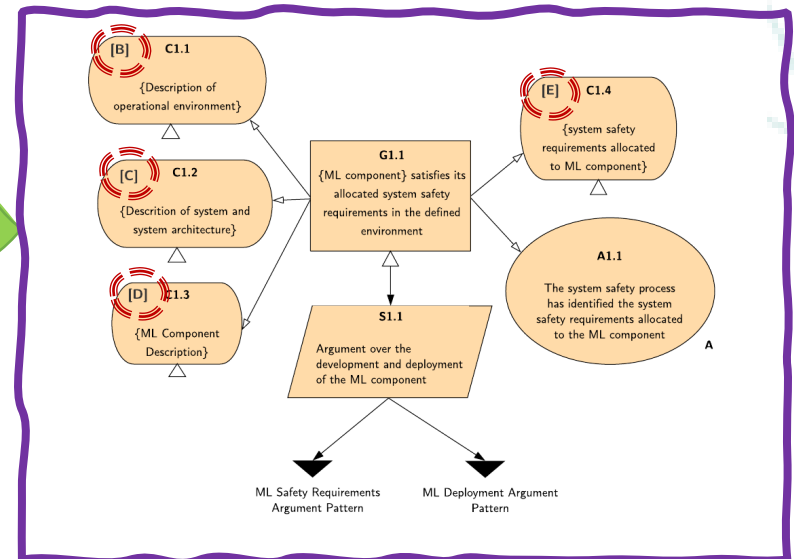
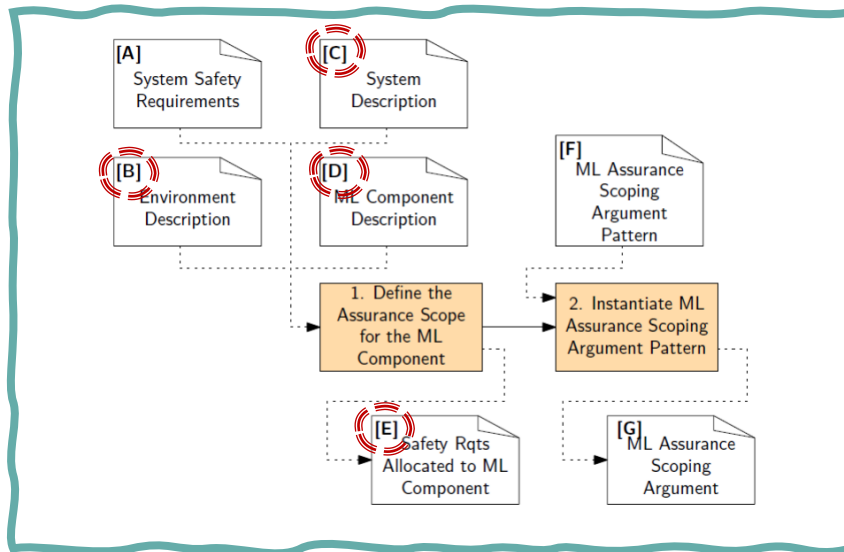
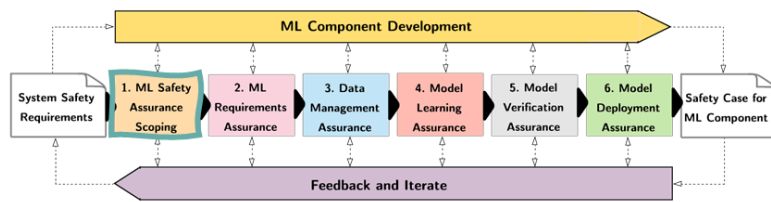


Allocating safety requirements

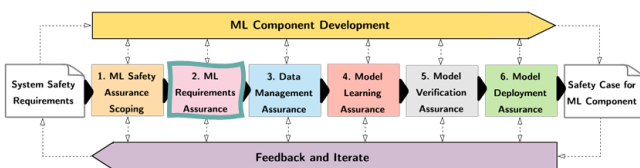
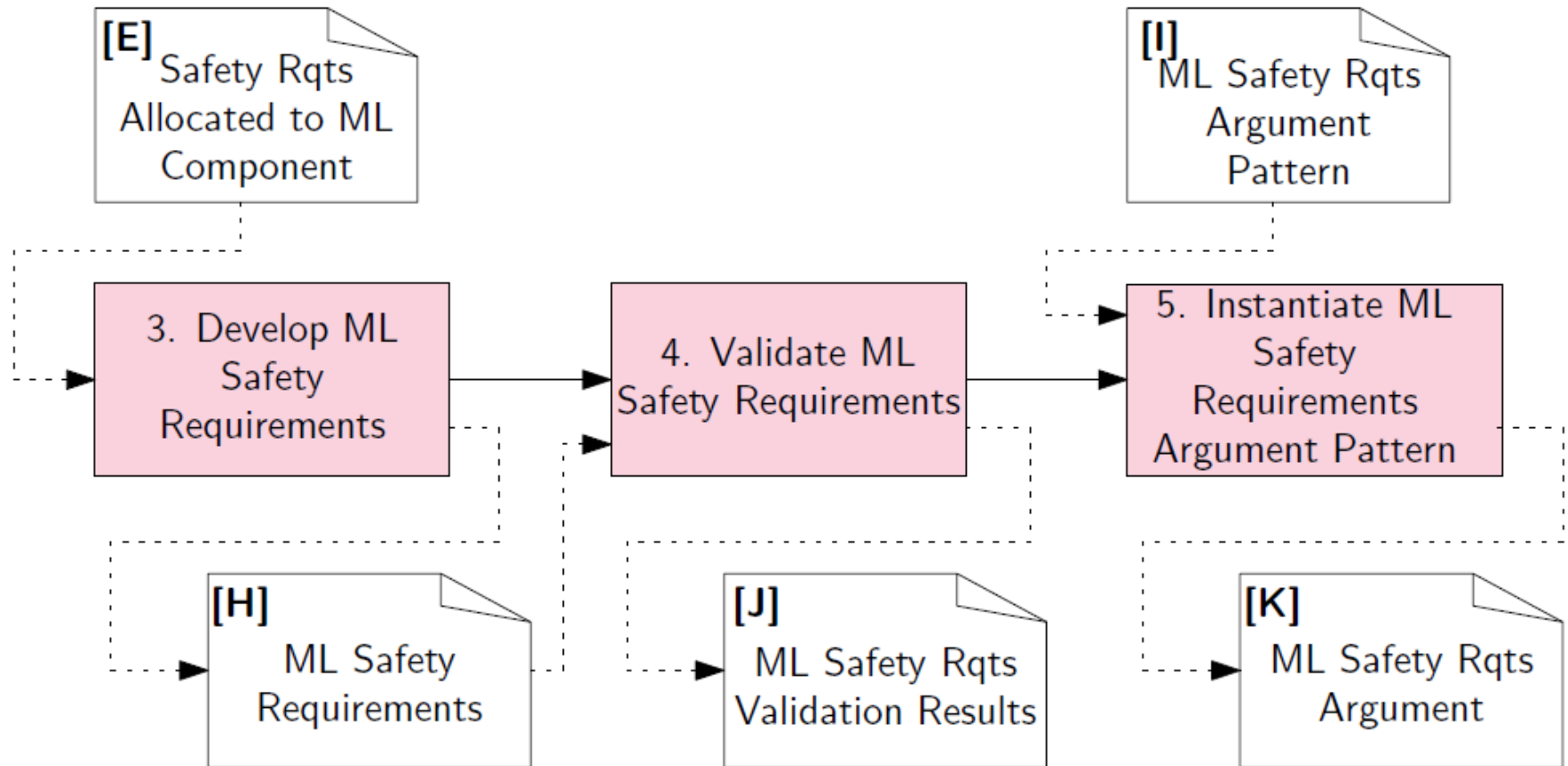
Example 1.

Consider an autonomous driving application in which a subsystem may be required to identify pedestrians at a crossing. A component within the perception pipeline may have a requirement of the form “When Ego is 50 metres from the crossing, the object detection component shall identify pedestrians that are on or close to the crossing in their correct position.”





2. ML Safety Rqts Assurance



Environment Context



What are the **key features** of the operating environment?

- **People, Vehicles, weather conditions, lighting, road type, etc., etc.**

ML Safety Requirements

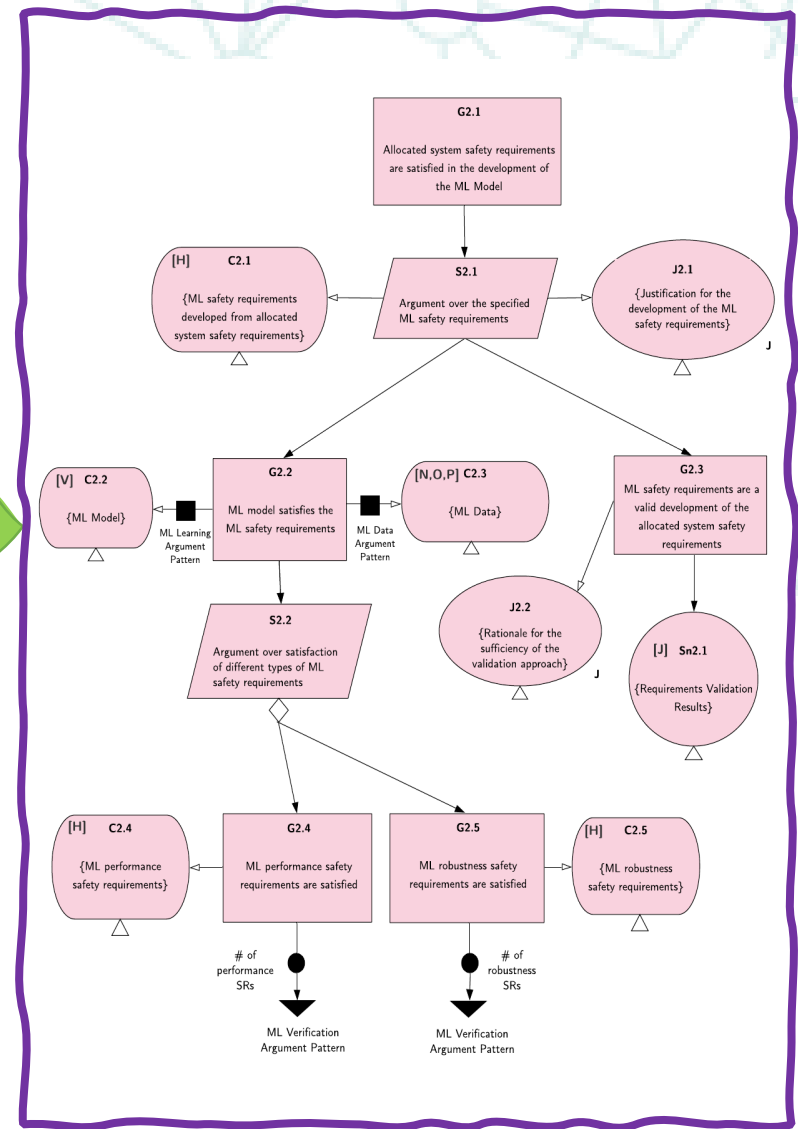
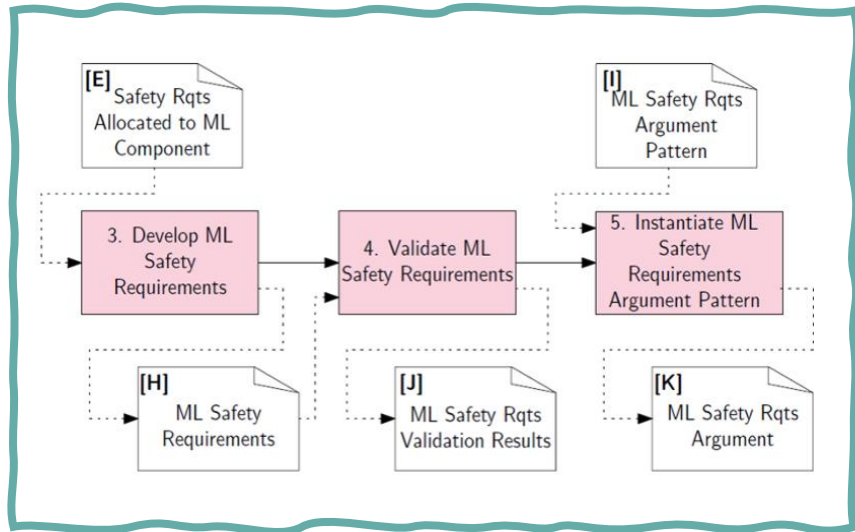
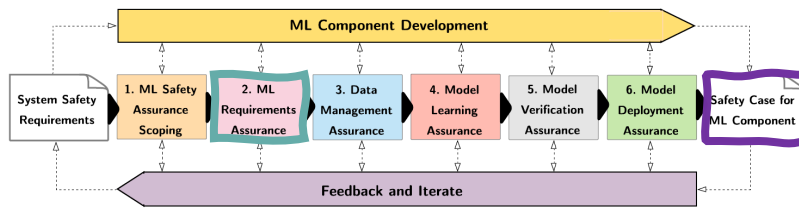
Example 7.

The ML safety requirement presented in Example 1 may now be refined into performance and robustness requirements [22]. Example performance requirements may include:

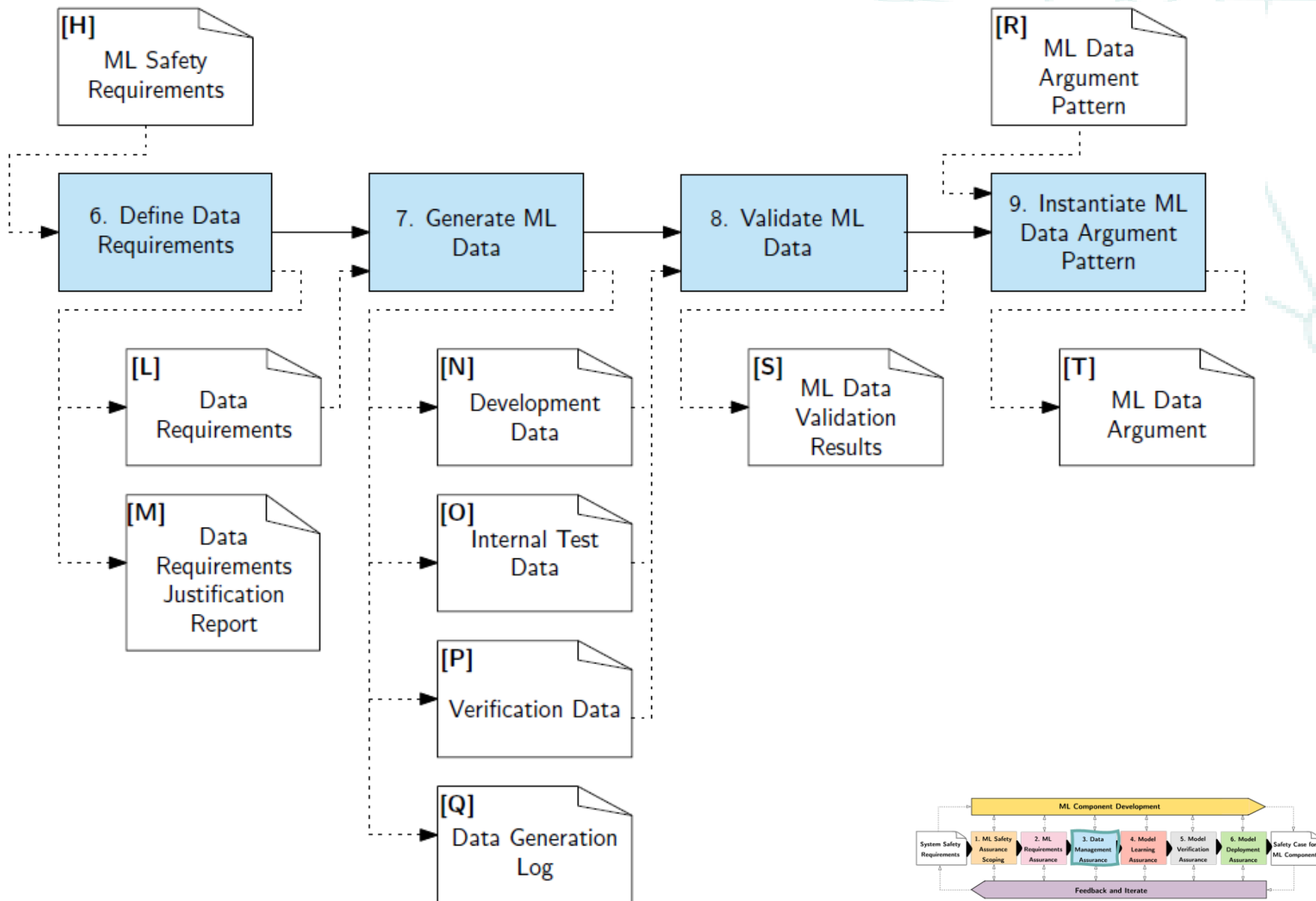
- The ML component shall determine the position of the specified feature in each input frame within 5 pixels of actual position.
- The ML component shall identify the presence of any person present in the defined area with an accuracy of at least 0.93

Example robustness requirements may include:

- The ML component shall perform as required in the defined range of lighting conditions experienced during operation of the system.
- The ML component shall identify a person irrespective of their pose with respect to the camera.



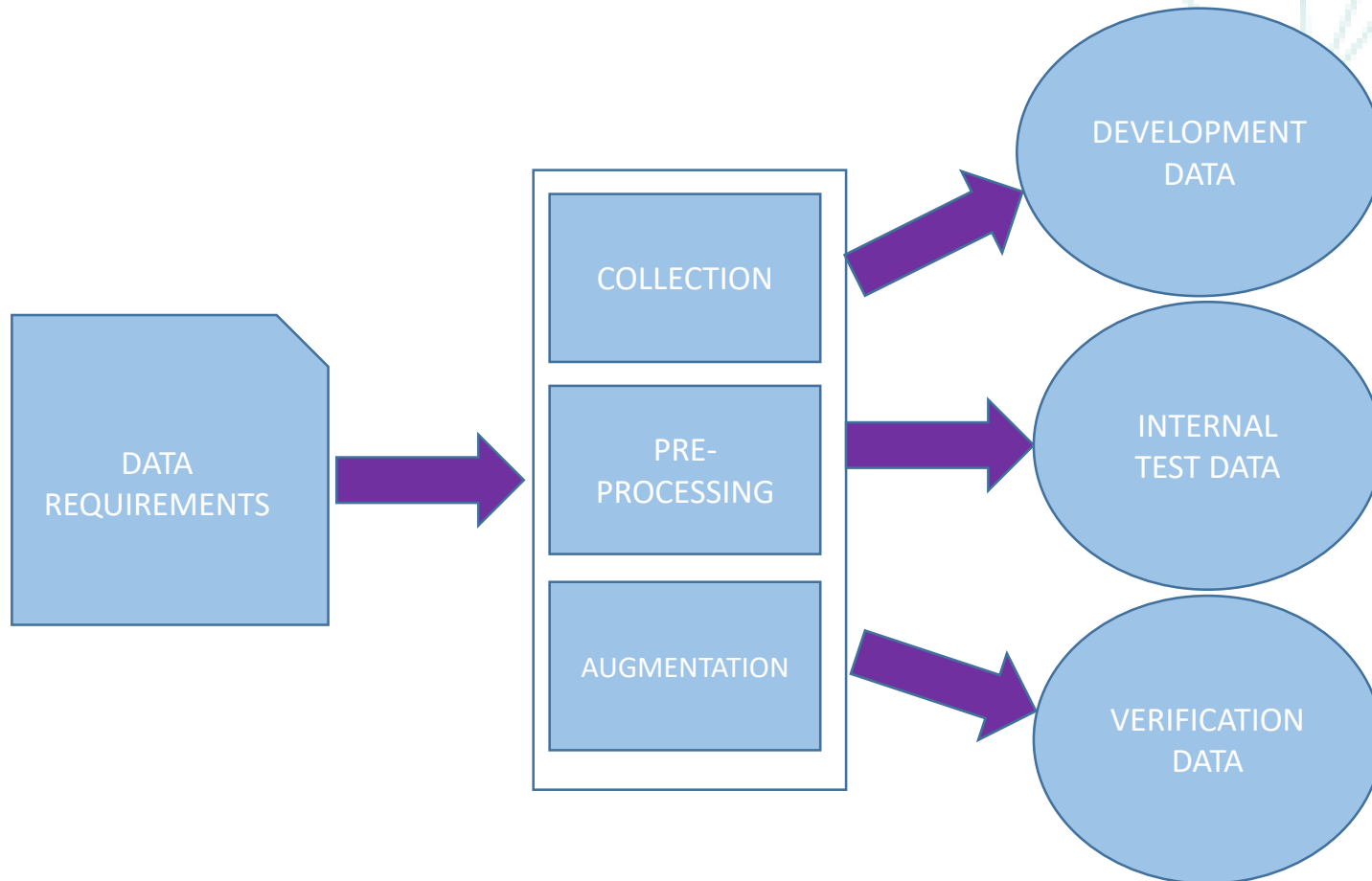
3. Data Management

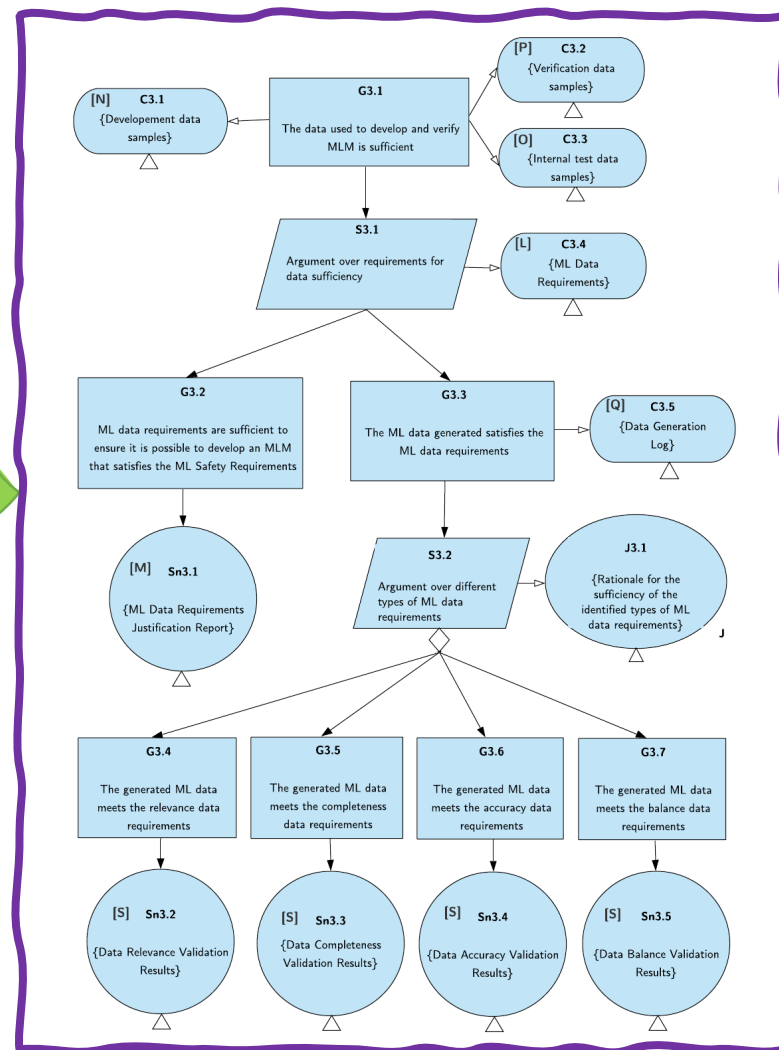
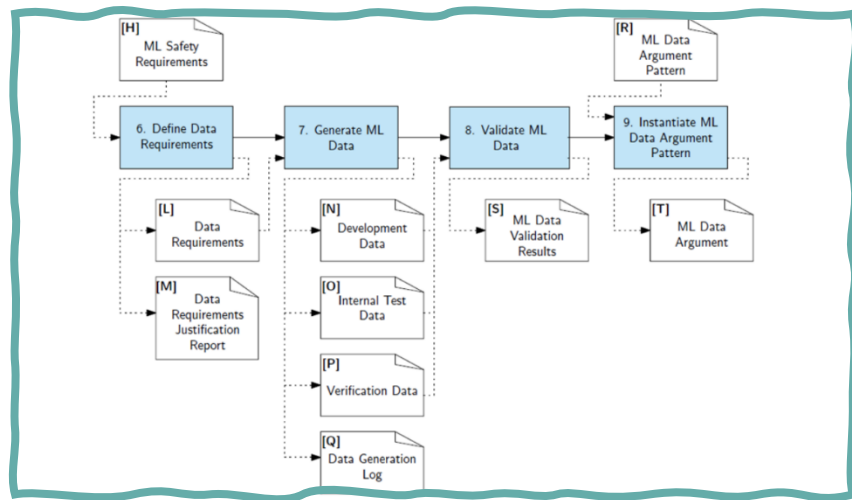
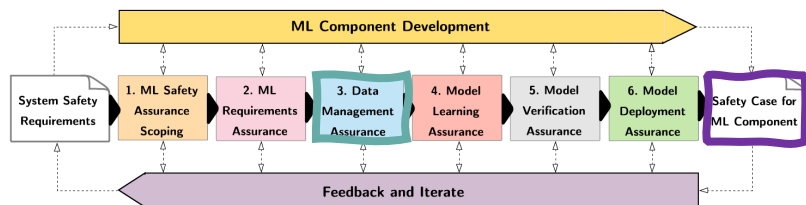


Data Requirements

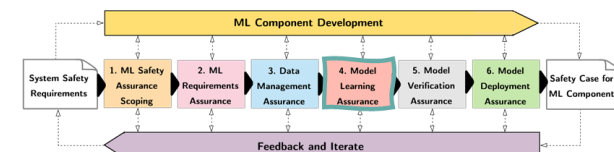
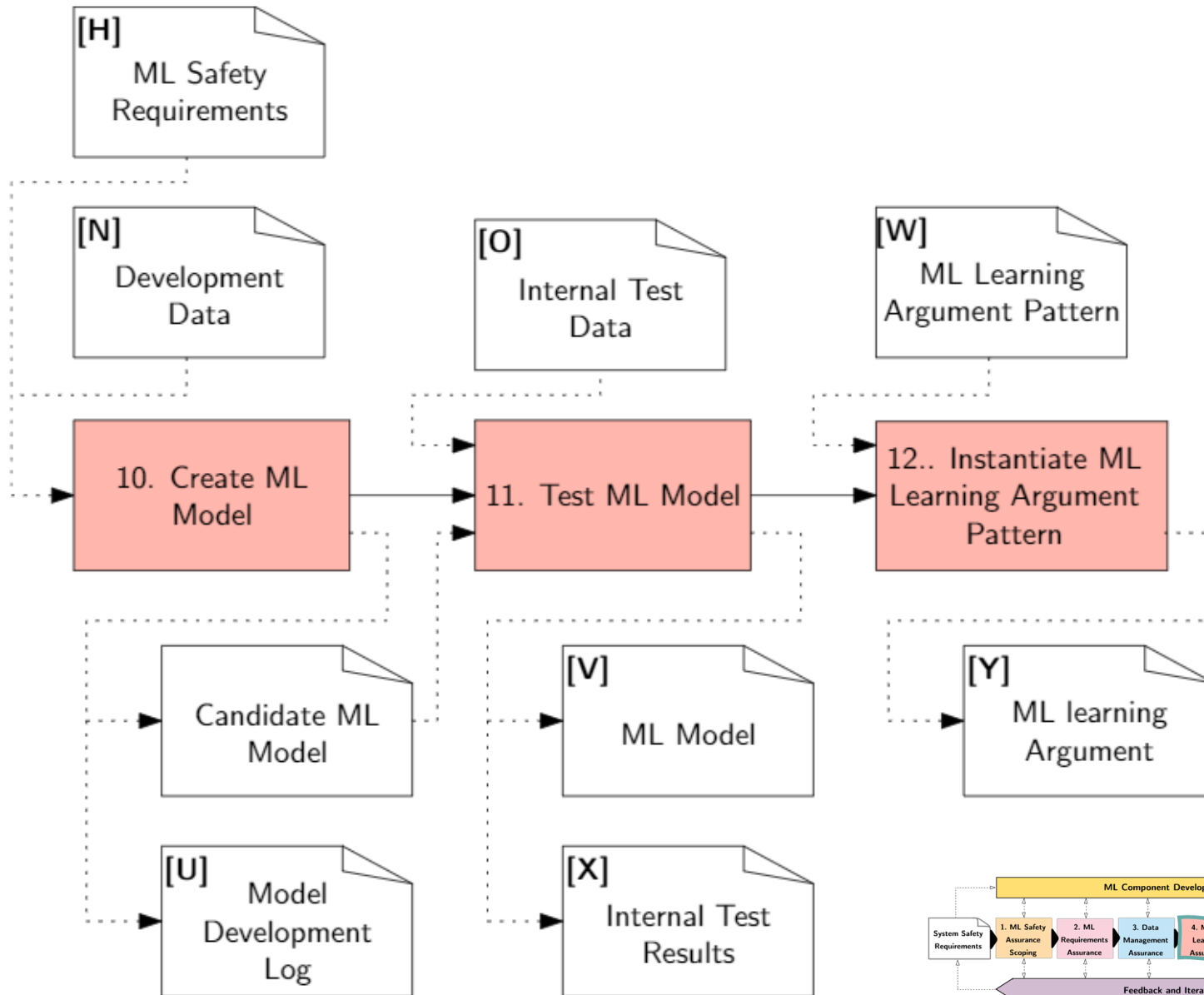
- ML data requirements shall include consideration of:
 - **Relevance**
 - **Completeness**
 - **Accuracy**
 - **Balance**

Generate ML data





4. Model Learning



4. Model Learning Assurance

Key Assurance Artefacts

Model Development Log

Q1. What forms of model were considered and on what basis was the model type selected?

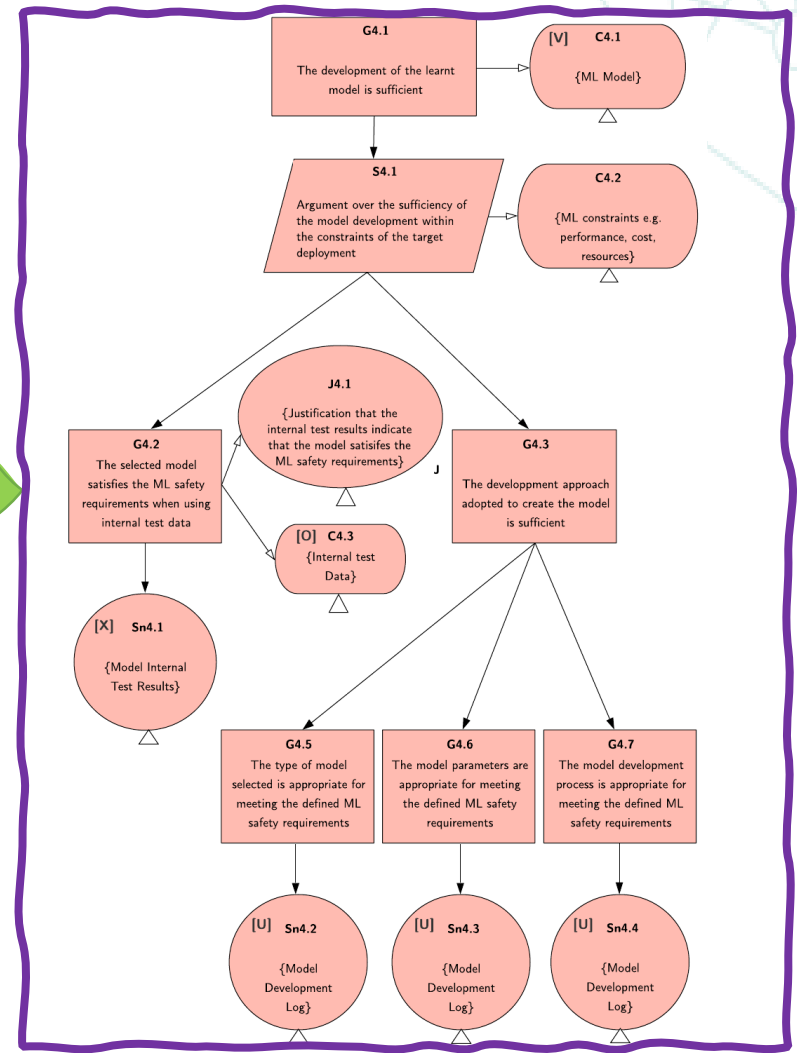
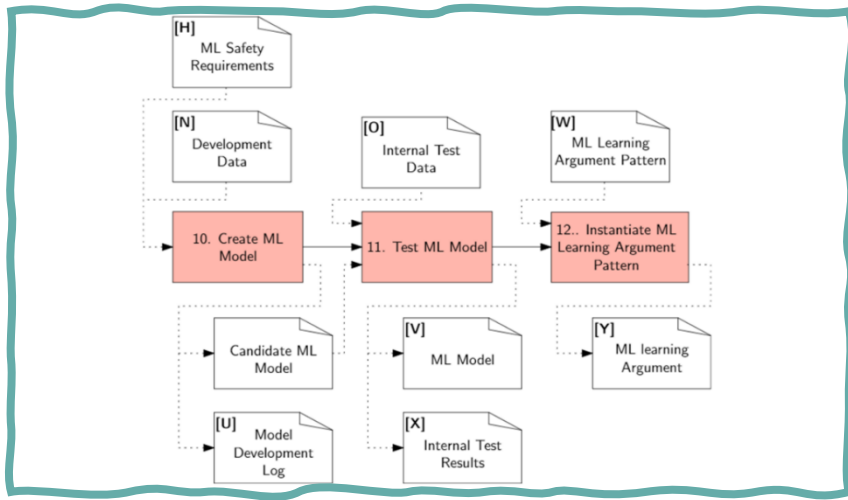
Q2. What approaches were used for hyper parameter tuning e.g. meta learning?

Q3. What techniques were applied to increase the generalizability of the model e.g. drop out?

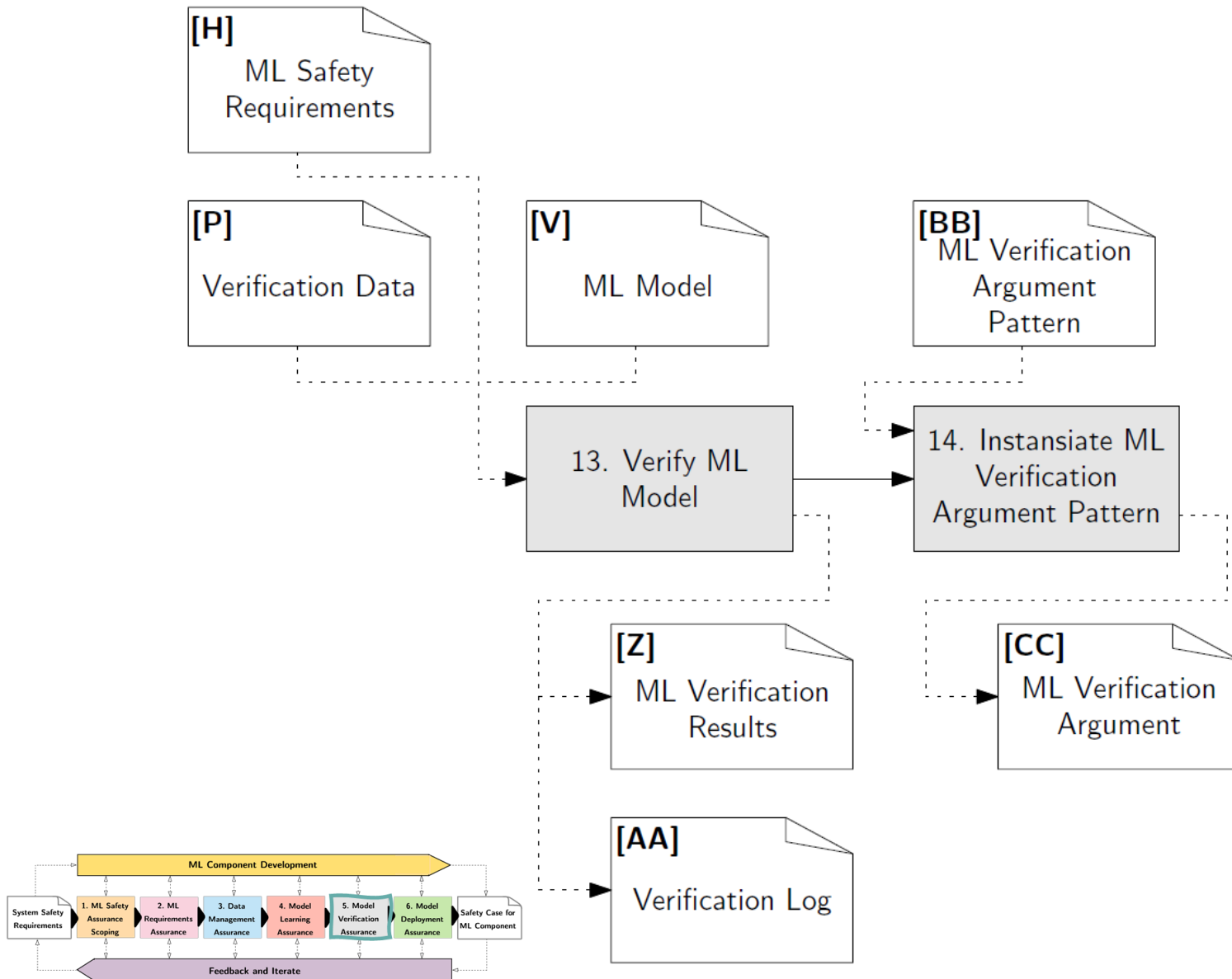
Q4. ...

Internal Test Results

Why was this particular model chosen?



5. Model Verification



Test-based Verification

- Use verification data to demonstrate model generalises to cases not present in model learning stage.
- Examine cases on boundaries or which are known to be problematic within the deployment context

Example 34.

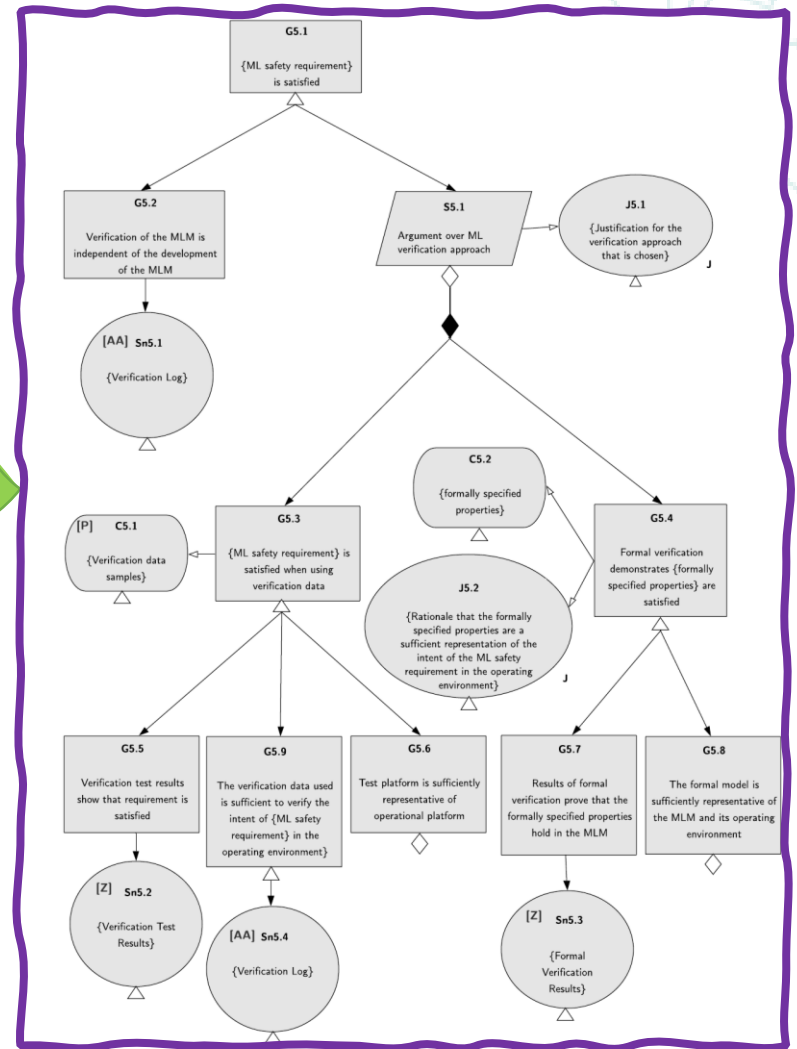
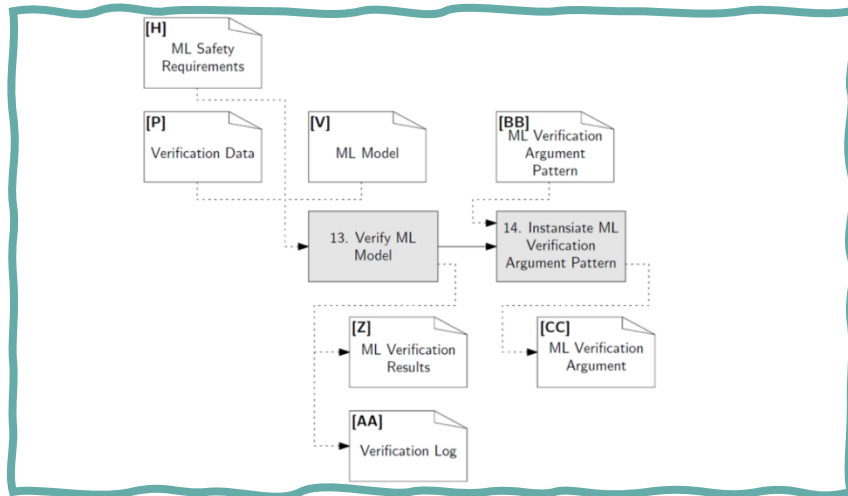
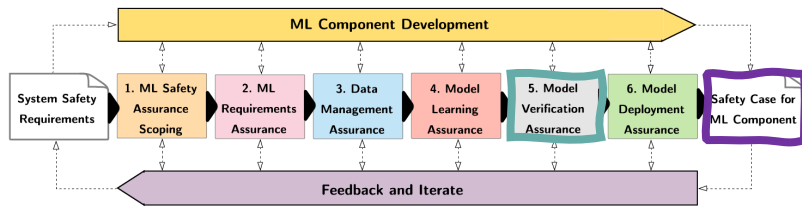
Since we know that material on a camera lens can lead to blurring in regions of an image, we may make use of 'contextual mutators' [52] to assess the robustness of a neural network with respect to levels of blur. In this way the level of blur which can be accommodated can be assessed and related to contextually meaningful measures.

Formal Verification

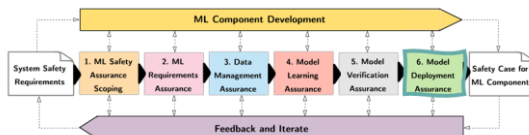
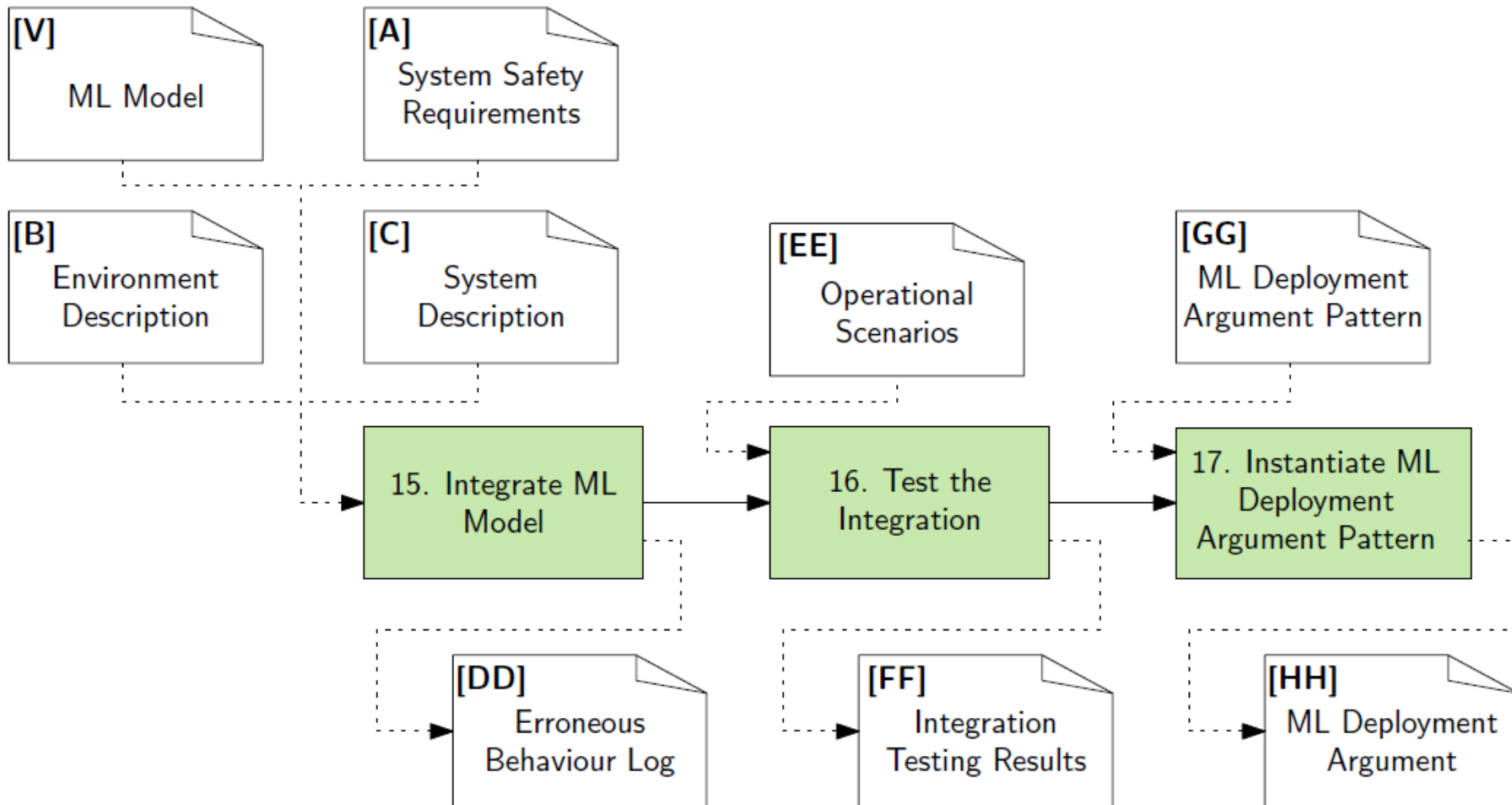
- Use mathematical techniques to prove that the learnt model satisfies formally specified properties
- Formal properties must be derived from the ML safety requirements.
 - Crucial to justify the sufficiency of the translation
 - Formal results must be mapped to operating context

Example 36.

Formal verification of neural networks are able to demonstrate that a perturbation of up to threshold value, ϵ , on all the inputs leads to a classification that is the same as the original sample class [36, 35]. Such a result is only meaningful when the value of ϵ is translated into contextually meaningful values such as steering angles. It is of little value when it simply provides a range of variation in pixel values for a single input sample (point wise verification).



6. Model Deployment



Integrate ML Model

- Monitor
 - Validity of key system and environment assumptions
 - output and internal states of ML model during operation
 - erroneous inputs to the ML model
 - noise and uncertainty in other components
 - environmental uncertainty
 - adversarial behaviours
- Can use statistical techniques to check inputs are close to the training data distributions

Integrate ML Model

- System must be designed to maintain safety even in face of predicted erroneous outputs from ML model

Example 46.

The ML model for pedestrian detection deployed in a self driving car has a performance requirement of 80% accuracy. Due to uncertainty in the model this performance cannot be achieved for every frame. The model uses as inputs a series of multiple images derived from consecutive image frames obtained from a camera. The presence of a pedestrian is determined by considering the result in the majority of the frames in the series. In this way the system compensates for the possible error of the model for any single image used.

Test the integration

- Following integration must check the system safety requirements are satisfied
- Requires operational scenarios against which behaviour implemented in ML can be tested
- Integration testing may include:
 - Simulation
 - Hardware in the loop
 - Shadow deployment

```

graph TD
    C1.1([C1.1  
Description of  
operational environment]) --> G1.1
    C1.2([C1.2  
Description of system and  
system architecture]) --> G1.1
    C1.3([C1.3  
ML Component  
Description]) --> G1.1
    G1.1[G1.1  
ML component's status is  
allocated system safety  
requirements in the defined  
environment] --> E1.4([E1.4  
System safety  
requirements allocated to  
ML component])
    E1.4 --> A1.1a([A1.1  
The system safety process  
has identified the system  
safety requirements allocated  
to the ML component])
    A1.1b([A1.1  
The system safety process  
has identified the system  
safety requirements allocated  
to the ML component])
  
```

On-Going Work



- Case studies
 - Wildfire Detection
 - Sepsis Diagnosis
- Interactive Website
- Tooling
- Similar Guidance for System level



**ASSURING
AUTONOMY**
INTERNATIONAL PROGRAMME

Funded by



Lloyd's Register
Foundation



UNIVERSITY
of York