

# Q in AI: Requirements for tools and processes for assurance

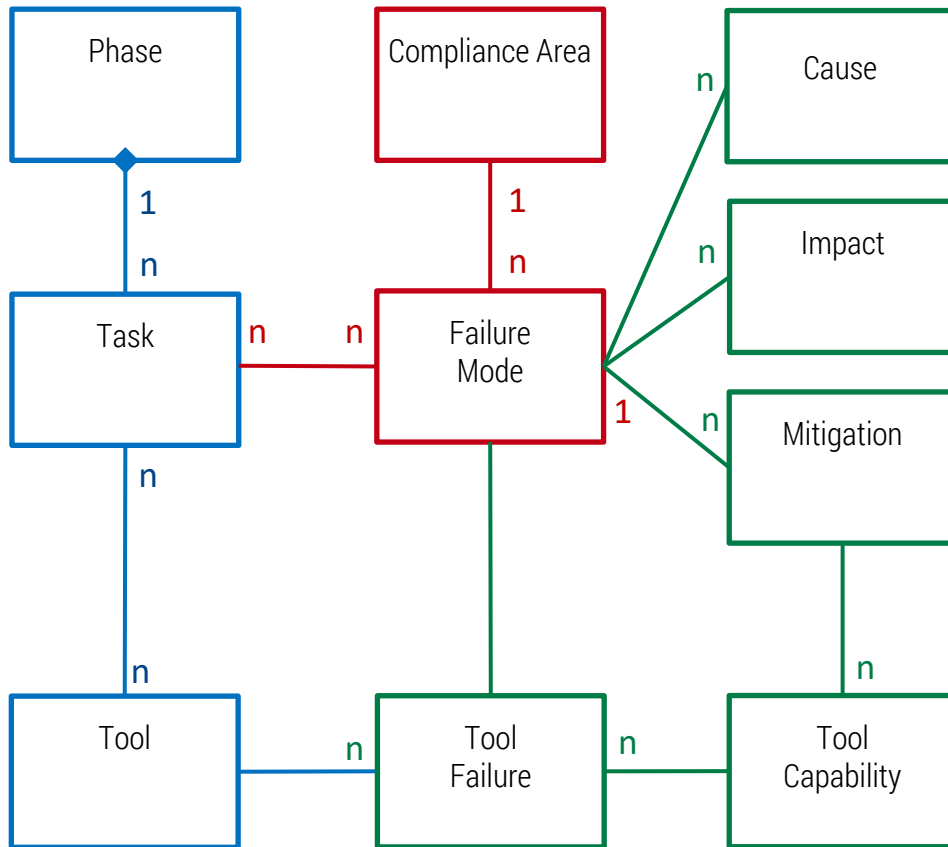
Instructions for qualifying tools AssessML

Dr. Björn Schünemann ([bjoern.schuenemann@aqigmbh.de](mailto:bjoern.schuenemann@aqigmbh.de))

Dr. Jürgen Großmann ([juergen.grossmann@fokus.fraunhofer.de](mailto:juergen.grossmann@fokus.fraunhofer.de))

# Relationship between concepts

Information model: Project procedure and terminology used



- Systematic listing of life cycle phases (phase) and AI tasks (task) as well as exemplary tools (tool).
- Identification of the failure mode (failure mode) along the AI tasks for each compliance area (compliance area).
- Identification of causes, effects (impact) and possible countermeasures (mitigation) for each failure mode.
- Identification of specific tool failures (tool failure) and required tool capabilities (tool capability) for a failure mode.

**Phase:** Phase from the AI life cycle

**Task:** Task during the AI life cycle that is supported by tools .

**Tool:** Tool that supports/executes a task.

**Compliance Area:** Area of compliance (functional safety, data protection, AI regulation).

**Failure Mode:** Hazard in the event of non-compliance.

**Tool Failure:** Tool error as the cause of an error condition.

**Cause:** Causes of the risk.

**Impact:** Effects of an error condition.

**Mitigation:** Measures to reduce the risk.

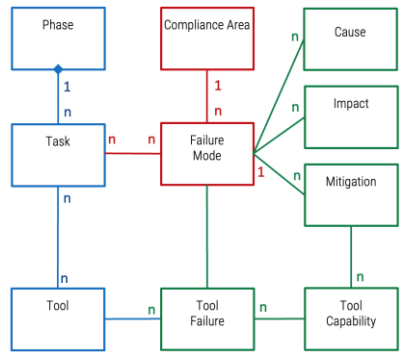
**Tool Capability:** Capabilities required of a tool to address the risk.

# Systematic presentation of the developed content

Result is a Python-based web application

## Terminology

The contents of the Excel spreadsheet are compiled in accordance with the terminology used



Compliance Area	Phase	Task	Failure Mode	Impact	Severity	Causes	Mitigation	Occurrence	Detectability	RPN	Potential Tool Failure	Recommended Tool Capability
Functional Safety	Training & Validation	Deep Learning & ML Frameworks	Lack of Real-Time Processing	Trained models and supporting frameworks are unable to deliver low-latency, real-time inference performance required for critical applications such as autonomous systems.		Model architectures are too large or complex for target real-time environments. Frameworks do not leverage hardware acceleration efficiently (e.g., GPU, TPU, Edge devices). Pipeline overhead (e.g., preprocessing, postprocessing) adds excessive latency during inference.	Apply model optimization techniques such as pruning, quantization, or knowledge distillation. Ensure deployment-ready models are benchmarked and optimized for specific hardware accelerators. Optimize full pipeline latency by batching operations and reducing unnecessary transformations.				May fail to guarantee deterministic training execution and runtime validation, resulting in trained models that do not meet real-time safety-critical performance requirements.	Framework support for model optimization (e.g., TensorRT, ONNX Runtime optimizations). Support for hardware-specific runtime optimizations and deployment targeting (e.g., EdgeTPU, CUDA kernels). Framework support for pipeline fusion, minimal preprocessing APIs, and asynchronous inference.
Functional Safety	Training & Validation	Deep Learning & ML Frameworks	Inadequate Error Handling & Logging	Frameworks and models lack robust error detection, recovery, and structured logging capabilities, making root cause analysis and system recovery difficult.		Training and inference exceptions are not properly caught or logged. Losses, gradients, or other training dynamics are not systematically monitored for anomalies. Training and deployment pipelines lack comprehensive logging and metadata tracking.	Integrate structured exception handling and mandatory logging for all critical training and inference operations. Implement runtime monitors for numerical instabilities (e.g., NaNs, divergence) and alert on thresholds. Embed pipeline-wide metadata capture and systematic experiment logging into the training framework.				May fail to provide comprehensive error logging and runtime monitoring, resulting in trained models with limited traceability for failure analysis in safety-critical systems.	Framework support for structured error reporting, custom callbacks, and logging integration. Built-in hooks for anomaly detection during training and validation (gradient checkers, divergence detectors). ML lifecycle management tools (e.g., MLflow, Weights & Biases) integration for metadata and error tracking.
Functional Safety	Training & Validation	Deep Learning & ML Frameworks	Model Instability & Poor Generalization	Inconsistent training leading to unreliable AI behaviour in critical applications.		Lack of deterministic/reproducible training settings due to parallel computation and floating-point arithmetic. Limited built-in validation or monitoring. No built-in numerical stability mechanisms.	Implement deterministic training pipelines with seed control and precision configuration to ensure repeatability. Integrate continuous performance monitoring during training with automatic validation checkpoints. Use frameworks that incorporate gradient clipping, normalization, and regularization for improved numerical robustness.				May fail to enforce deterministic training procedures and runtime validation, resulting in trained models with unpredictable behavior and poor generalization in safety-critical applications.	Explicit deterministic training support with reproducibility settings and controlled parallelism. Built-in validation dashboards and metric tracking APIs. Numerical stability modules with configurable regularization options.

## Dashboard created based on Excel spreadsheet

**Tool Qualification Recommendation App**

Choose Phase(s):  Choose Task(s):  Choose Tool(s):

Failure Mode Selection | Task Assessment | Tool Assessment | Assessment Summary

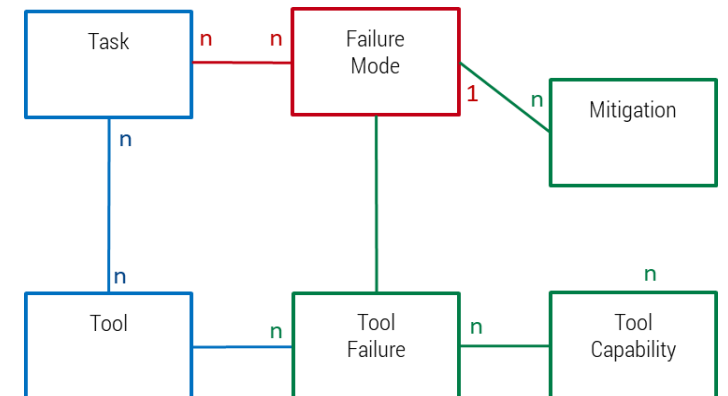
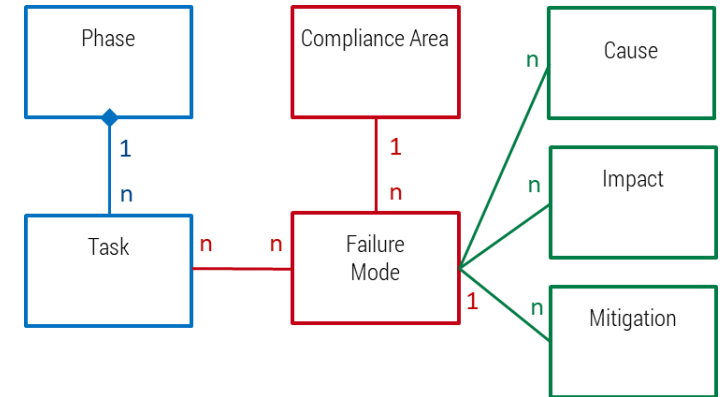
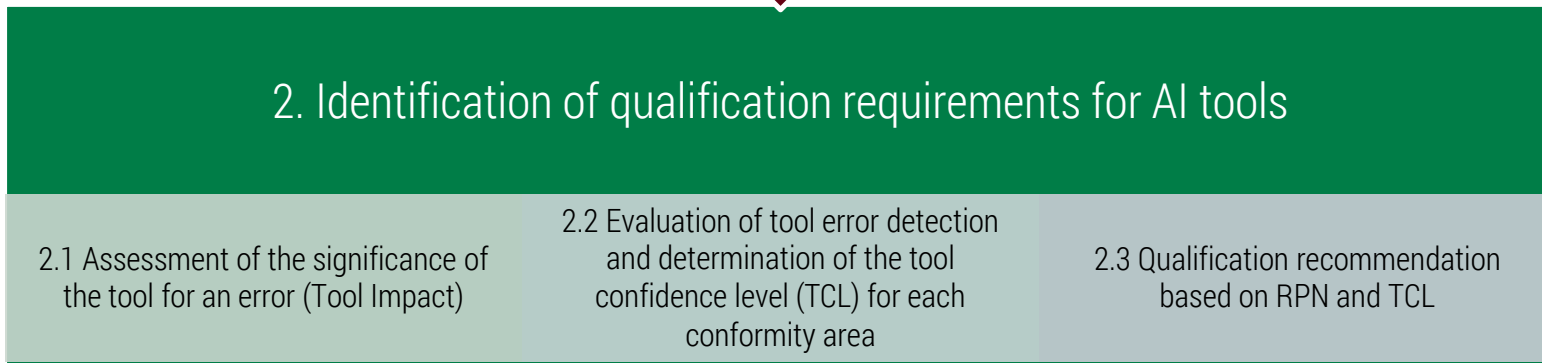
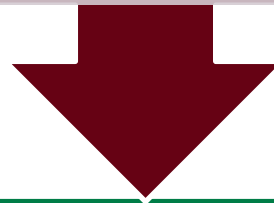
**Failure Mode and Compliance Area Selection**

Select	Failure Mode	Description
<input checked="" type="checkbox"/>	<b>AI Regulation</b>	
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input checked="" type="checkbox"/>	Lack of Model Risk Assessment	Failure to categorize AI models under AI Act risk levels and document safety concerns.
<input checked="" type="checkbox"/>	Security Vulnerabilities in AI	Exposure to adversarial attacks compromising AI safety and trustworthiness.
<input checked="" type="checkbox"/>	Non-compliant Human Oversight Mechanisms	Failure to implement necessary human-in-the-loop controls in high-risk AI applications.
<input checked="" type="checkbox"/>	Insufficient AI Impact Assessment	Failure to conduct proper risk-benefit analysis for high-risk AI applications.
<input checked="" type="checkbox"/>	Lack of Ethical Safeguards in AI	Failure to ensure ethical principles (e.g., non-harm, fairness) in AI decision-making.
<input checked="" type="checkbox"/>	Inadequate Monitoring of Deployed AI Systems	Lack of continuous assessment of AI systems in real-world settings.
<input checked="" type="checkbox"/>	Lack of Auditability & Documentation	Inability to provide detailed logs of data usage, model training, validation etc. hindering compliance audits.
<input checked="" type="checkbox"/>	<b>Data Protection</b>	
<input checked="" type="checkbox"/>	Unsanctioned Data Access	Lack of strict access control leading to potential data breaches and privacy violations.
<input checked="" type="checkbox"/>	Failure to Handle Right to Erasure	Non-compliance with GDPR's Right to be Forgotten, leading to legal risks.
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.
<input checked="" type="checkbox"/>	Lack of Data Minimization	Collecting and storing excessive personal data, violating GDPR principles.
<input checked="" type="checkbox"/>	Failure in Data Anonymization & Pseudonymization	Inadequate de-identification techniques lead to re-identifiable personal data.
<input checked="" type="checkbox"/>	Non-Transparent User Consent Management	Lack of clear user consent tracking for data collection and processing.
<input checked="" type="checkbox"/>	Cross-border Data Transfer Violations	Failure to comply with international data transfer rules under GDPR/Schrems II.
<input checked="" type="checkbox"/>	<b>Functional Safety</b>	
<input checked="" type="checkbox"/>	Data Integrity Failure	Corrupt or inconsistent data leading to incorrect model predictions, endangering safety-critical systems.
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.
<input checked="" type="checkbox"/>	Failure Mode Undetection	Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.
<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility	Non-deterministic model outputs lead to unpredictable AI system behavior.
<input checked="" type="checkbox"/>	Insufficient Edge Case Handling	AI models fail to correctly handle rare but critical safety scenarios, leading to unreliable operation.
<input checked="" type="checkbox"/>	Undetected Model Performance Drift	Model performance degrades over time without triggering corrective actions, leading to unsafe decisions.
<input checked="" type="checkbox"/>	Insufficient Redundancy & Failover Mechanisms	Lack of backup AI models or failover strategies in critical systems.
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

Settings (Remember (e.g., with)) | Save Settings | Load Settings

# Risk-based tool qualification (details)

Two-stage process for identifying task-related risks and qualifying tools based on the risks



# AssessML: structure of the dashboard

Homepage

## Tool Qualification Recommendation App

Choose Phase(s):  Choose Task(s):  Choose Tool(s):

Failure Mode Selection	Task Assessment	Tool Assessment	Assessment Summary
<b>Failure Mode and Compliance Area Selection</b>			
<b>Select</b>	<b>Failure Mode</b>	<b>Description</b>	
<input checked="" type="checkbox"/> <b>AI Regulation</b>			
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.	
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.	
<input checked="" type="checkbox"/>	Lack of Model Risk Assessment	Failure to categorize AI models under AI Act risk levels and document safety concerns.	
<input checked="" type="checkbox"/>	Security Vulnerabilities in AI	Exposure to adversarial attacks compromising AI safety and trustworthiness.	
<input checked="" type="checkbox"/>	Non-compliant Human Oversight Mechanisms	Failure to implement necessary human-in-the-loop controls in high-risk AI applications.	
<input checked="" type="checkbox"/>	Insufficient AI Impact Assessment	Failure to conduct proper risk-benefit analysis for high-risk AI applications.	
<input checked="" type="checkbox"/>	Lack of Ethical Safeguards in AI	Failure to ensure ethical principles (e.g., non-harm, fairness) in AI decision-making.	
<input checked="" type="checkbox"/>	Inadequate Monitoring of Deployed AI Systems	Lack of continuous assessment of AI systems in real-world settings.	
<input checked="" type="checkbox"/>	Lack of Auditability & Documentation	Inability to provide detailed logs of data usage, model training, validation etc. hindering compliance audits.	
<input checked="" type="checkbox"/> <b>Data Protection</b>			
<input checked="" type="checkbox"/>	Unauthorized Data Access	Lack of strict access control leading to potential data breaches and privacy violations.	
<input checked="" type="checkbox"/>	Failure to Handle Right to Erasure	Non-compliance with GDPR's Right to be Forgotten, leading to legal risks.	
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.	
<input checked="" type="checkbox"/>	Lack of Data Minimization	Collecting and storing excessive personal data, violating GDPR principles.	
<input checked="" type="checkbox"/>	Failure in Data Anonymization & Pseudonymization	Inadequate de-identification techniques lead to re-identifiable personal data.	
<input checked="" type="checkbox"/>	Non-Transparent User Consent Management	Lack of clear user consent tracking for data collection and processing.	
<input checked="" type="checkbox"/>	Cross-border Data Transfer Violations	Failure to comply with international data transfer rules under GDPR/Schrems II.	
<input checked="" type="checkbox"/> <b>Functional Safety</b>			
<input checked="" type="checkbox"/>	Data Integrity Failure	Corrupt or inconsistent data leading to incorrect model predictions, endangering safety-critical systems.	
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.	
<input checked="" type="checkbox"/>	Failure Mode Undetection	Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.	
<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility	Non-deterministic model outputs lead to unpredictable AI system behavior.	
<input checked="" type="checkbox"/>	Insufficient Edge Case Handling	AI models fail to correctly handle rare but critical safety scenarios, leading to unreliable operation.	
<input checked="" type="checkbox"/>	Undetected Model Performance Drift	Model performance degrades over time without triggering corrective actions, leading to unsafe decisions.	
<input checked="" type="checkbox"/>	Insufficient Redundancy & Failover Mechanisms	Lack of backup AI models or failover strategies in critical systems.	
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.	
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.	

Settings filename (e.g. setti)

The homepage is divided into three sections:

[A: Selection of phases and tasks](#)

Selection tabs for phases, tasks or tools to be considered.

[B: Selection of error states and assessment implementation](#)

Used to select the error states to be considered, as well as to determine risks in the development tasks and the qualification requirements for AI tools.

[C: Saving and loading](#)

Preset or already completed assessments can be saved and retrieved.

# 1. Identification of risks in development tasks

## 1.1. Selection of relevant tasks phases

1. Identification of risks in development tasks (according to FMEA)

$$\text{RPN} = \text{Severity} * \text{Occurrence} * \text{Detection}$$

1.1. Selection of relevant tasks  
phases

1.2 Selection of area of compliance  
and applicable failure modes

1.3 Calculation of the RPN for  
selected failure modes in a task

# AssessML: selection area

What is the selection area used for?

The dashboard contains 25 error states that are selected from the start.

This automatically transfers all statuses to the task and tool assessment as well as to the assessment summary, which can quickly make the board confusing.

Select	Failure Mode	Description
<input checked="" type="checkbox"/>	<b>AI Regulation</b>	
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input checked="" type="checkbox"/>	Lack of Model Risk Assessment	Failure to categorize AI models under AI Act risk levels and document safety concerns.
<input checked="" type="checkbox"/>	Security Vulnerabilities in AI	Exposure to adversarial attacks compromising AI safety and trustworthiness.
<input checked="" type="checkbox"/>	Non-compliant Human Oversight Mechanisms	Failure to implement necessary human-in-the-loop controls in high-risk AI applications.
<input checked="" type="checkbox"/>	Insufficient AI Impact Assessment	Failure to conduct proper risk-benefit analysis for high-risk AI applications.
<input checked="" type="checkbox"/>	Lack of Ethical Safeguards in AI	Failure to ensure ethical principles (e.g., non-harm, fairness) in AI decision-making.
<input checked="" type="checkbox"/>	Inadequate Monitoring of Deployed AI Systems	Lack of continuous assessment of AI systems in real-world settings.
<input checked="" type="checkbox"/>	Lack of Auditability & Documentation	Inability to provide detailed logs of data usage, model training, validation etc. hindering compliance audits.
<input checked="" type="checkbox"/>	<b>Data Protection</b>	
<input checked="" type="checkbox"/>	Unauthorized Data Access	Lack of strict access control leading to potential data breaches and privacy violations.
<input checked="" type="checkbox"/>	Failure to Handle Right to Erasure	Non-compliance with GDPR's Right to be Forgotten, leading to legal risks.
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.
<input checked="" type="checkbox"/>	Lack of Data Minimization	Collecting and storing excessive personal data, violating GDPR principles.
<input checked="" type="checkbox"/>	Failure in Data Anonymization & Pseudonymization	Inadequate de-identification techniques lead to re-identifiable personal data.
<input checked="" type="checkbox"/>	Non-Transparent User Consent Management	Lack of clear user consent tracking for data collection and processing.
<input checked="" type="checkbox"/>	Cross-border Data Transfer Violations	Failure to comply with international data transfer rules under GDPR/Schrems II.
<input checked="" type="checkbox"/>	<b>Functional Safety</b>	
<input checked="" type="checkbox"/>	Data Integrity Failure	Corrupt or inconsistent data leading to incorrect model predictions, endangering safety-critical systems.
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.
<input checked="" type="checkbox"/>	Failure Mode Undetection	Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.
<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility	Non-deterministic model outputs lead to unpredictable AI system behavior.
<input checked="" type="checkbox"/>	Insufficient Edge Case Handling	AI models fail to correctly handle rare but critical safety scenarios, leading to unreliable operation.
<input checked="" type="checkbox"/>	Undetected Model Performance Drift	Model performance degrades over time without triggering corrective actions, leading to unsafe decisions.
<input checked="" type="checkbox"/>	Insufficient Redundancy & Failover Mechanisms	Lack of backup AI models or failover strategies in critical systems.
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

Task Assessment

Tool Assessment

Assessment Summary

# AssessML: selection area

What is the selection area used for?

To enable more targeted and efficient processing, the dashboard offers the option of filtering via the **selection area**.

Each selection made reduces the number of error states displayed, thereby simplifying the processing of the following steps in the assessments and in the Assessment Summary.

The screenshot shows the top navigation area of the AssessML dashboard. It features three dropdown menus labeled 'Choose Phase(s):', 'Choose Task(s):', and 'Choose Tool(s):', each with a 'Select...' option. Below these menus is a horizontal bar with four tabs: 'Failure Mode Selection', 'Task Assessment', 'Tool Assessment', and 'Assessment Summary'. A red box highlights the selection area, including the dropdown menus and the tabs.

**Failure Mode and Compliance Area Selection**

Select	Failure Mode	Description
<input checked="" type="checkbox"/>	<b>AI Regulation</b>	
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing bias, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input checked="" type="checkbox"/>	Lack of Model Risk Assessment	Failure to categorize AI models under AI Act risk levels and document safety concerns.
<input checked="" type="checkbox"/>	Security Vulnerabilities in AI	Exposure to adversarial attacks compromising AI safety and trustworthiness.
<input checked="" type="checkbox"/>	Non-compliant Human Oversight Mechanisms	Failure to implement necessary human-in-the-loop controls in high-risk AI applications.
<input checked="" type="checkbox"/>	Insufficient AI Impact Assessment	Failure to conduct proper risk-benefit analysis for high-risk AI applications.
<input checked="" type="checkbox"/>	Lack of Ethical Safeguards in AI	Failure to ensure ethical principles (e.g., non-harm, fairness) in AI decision-making.
<input checked="" type="checkbox"/>	Inadequate Monitoring of Deployed AI Systems	Lack of continuous assessment of AI systems in real-world settings.
<input checked="" type="checkbox"/>	Lack of Auditability & Documentation	Inability to provide detailed logs of data usage, model training, validation etc. hindering compliance audits.
<input checked="" type="checkbox"/>	<b>Data Protection</b>	
<input checked="" type="checkbox"/>	Unauthorized Data Access	Lack of strict access control leading to potential data breaches and privacy violations.
<input checked="" type="checkbox"/>	Failure to Handle Right to Erasure	Non-compliance with GDPR's Right to be Forgotten, leading to legal risks.
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.
<input checked="" type="checkbox"/>	Lack of Data Minimization	Collecting and storing excessive personal data, violating GDPR principles.
<input checked="" type="checkbox"/>	Failure in Data Anonymization & Pseudonymization	Inadequate de-identification techniques lead to re-identifiable personal data.
<input checked="" type="checkbox"/>	Non-Transparent User Consent Management	Lack of clear user consent tracking for data collection and processing.
<input checked="" type="checkbox"/>	Cross-border Data Transfer Violations	Failure to comply with international data transfer rules under GDPR, Schengen II.
<input checked="" type="checkbox"/>	<b>Functional Safety</b>	
<input checked="" type="checkbox"/>	Data Integrity Failure	Corrupt or inconsistent data leading to incorrect model predictions, endangering safety-critical systems.
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.
<input checked="" type="checkbox"/>	Failure Mode Underestimation	Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.
<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility	Non-deterministic model outputs lead to unpredictable AI system behavior.
<input checked="" type="checkbox"/>	Insufficient Edge Case Handling	AI models fail to correctly handle rare but critical safety scenarios, leading to unreliable operation.
<input checked="" type="checkbox"/>	Undetected Model Performance Drift	Model performance degrades over time without triggering corrective actions, leading to unsafe decisions.
<input checked="" type="checkbox"/>	Insufficient Redundancy & Failover Mechanisms	Lack of backup AI models or failover strategies in critical systems.
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

The image displays four screenshots of the AssessML dashboard, each representing a different stage of the assessment process. The screenshots are arranged in a 2x2 grid. The top-left screenshot is titled 'Failure Mode Overview' and shows a list of failure modes with their descriptions and risk levels. The top-right screenshot is titled 'Task Assessment' and shows a table of tasks with their associated failure modes and risk levels. The bottom-left screenshot is titled 'Tool Assessment' and shows a table of tools with their associated failure modes and risk levels. The bottom-right screenshot is titled 'Assessment Summary' and shows a comprehensive overview of the assessment results, including a summary of the number of failure modes, tasks, and tools, and a detailed breakdown of the results.

Instructions for assessments follow from slide 18 onwards

# AssessML: selection area

## 1. Selection of phases and tasks

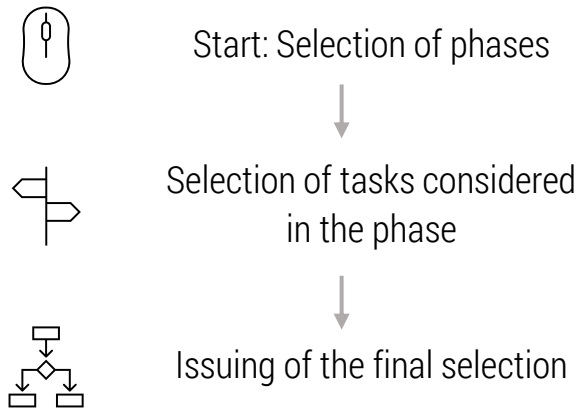
Choose Phase(s):

Choose Task(s):

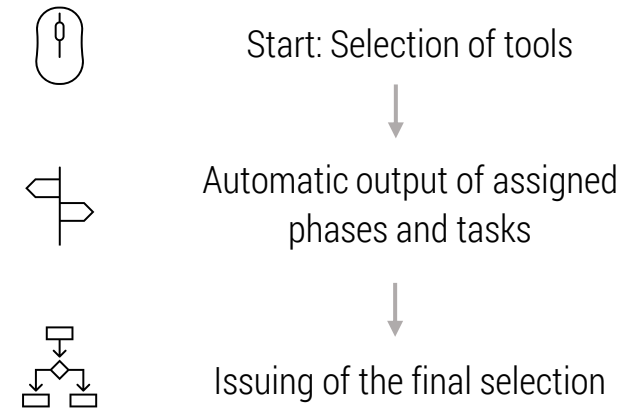
Choose Tool(s):

The selection area is divided into three **search tabs**: *Phase*, *Task* and *Tools*, allowing users to limit their search to their preferred phases and tasks. There are **two** ways to specify a phase and task:

### 1st approach: Phases and tasks



### 2nd approach: Tools



# AssessML: selection area

1st approach: Phases and tasks

Start: Selection of phases

Choose Phase(s):  Choose Task(s):  Choose Tool(s):

Task Assessment	Tool Assessment	Assessment Summary	
<b>Description</b>			
<input checked="" type="checkbox"/>		Lack of AI transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>		Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input checked="" type="checkbox"/>		Lack of Model Risk Assessment	Failure to categorize AI models under AI Act risk levels and document safety concerns.
<input checked="" type="checkbox"/>		Security Vulnerabilities in AI	Exposure to adversarial attacks compromising AI safety and trustworthiness.
<input checked="" type="checkbox"/>		Non-compliant Human Oversight Mechanisms	Failure to implement necessary human-in-the-loop controls in high-risk AI applications.
<input checked="" type="checkbox"/>		Insufficient AI Impact Assessment	Failure to conduct proper risk-benefit analysis for high-risk AI applications.
<input checked="" type="checkbox"/>		Lack of Ethical Safeguards in AI	Failure to ensure ethical principles (e.g., non-harm, fairness) in AI decision-making.
<input checked="" type="checkbox"/>		Inadequate Monitoring of Deployed AI Systems	Lack of continuous assessment of AI systems in real-world settings.
<input checked="" type="checkbox"/>		Lack of Auditability & Documentation	Inability to provide detailed logs of data usage, model training, validation etc. hindering compliance audits.
<b>Data Protection</b>			
<input checked="" type="checkbox"/>		Unauthorized Data Access	Lack of strict access control leading to potential data breaches and privacy violations.
<input checked="" type="checkbox"/>		Failure to Handle Right to Erasure	Non-compliance with GDPR's Right to be Forgotten, leading to legal risks.
<input checked="" type="checkbox"/>		Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.
<input checked="" type="checkbox"/>		Lack of Data Minimization	Collecting and storing excessive personal data, violating GDPR principles.
<input checked="" type="checkbox"/>		Failure in Data Anonymization & Pseudonymization	Inadequate de-identification techniques lead to re-identifiable personal data.
<input checked="" type="checkbox"/>		Non-Transparent User Consent Management	Lack of clear user consent tracking for data collection and processing.
<input checked="" type="checkbox"/>		Cross-border Data Transfer Violations	Failure to comply with international data transfer rules under GDPR/Schrems II.
<b>Functional Safety</b>			
<input checked="" type="checkbox"/>		Data Integrity Failure	Corrupt or inconsistent data leading to incorrect model predictions, endangering safety-critical systems.
<input checked="" type="checkbox"/>		Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.
<input checked="" type="checkbox"/>		Failure Mode Undetection	Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.
<input checked="" type="checkbox"/>		Inconsistent Model Reproducibility	Non-deterministic model outputs lead to unpredictable AI system behavior.
<input checked="" type="checkbox"/>		Insufficient Edge Case Handling	AI models fail to correctly handle rare but critical safety scenarios, leading to unreliable operation.
<input checked="" type="checkbox"/>		Undetected Model Performance Drift	Model performance degrades over time without triggering corrective actions, leading to unsafe decisions.
<input checked="" type="checkbox"/>		Insufficient Redundancy & Failover Mechanisms	Lack of backup AI models or failover strategies in critical systems.
<input checked="" type="checkbox"/>		Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.
<input checked="" type="checkbox"/>		Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

Settings filename (e.g. settj)

# AssessML: selection area

1st approach: Phases and tasks

Selection of tasks considered in the phase

The screenshot displays the AssessML interface with three main selection areas: 'Choose Phase(s)', 'Choose Task(s)', and 'Choose Tool(s)'. The 'Choose Phase(s)' dropdown is set to 'Training & Validation'. The 'Choose Task(s)' dropdown is open, showing a list of tasks with 'Deep Learning & ML Frameworks' selected. The 'Choose Tool(s)' dropdown is set to 'Assessment'. Below these, the 'Failure Mode and Compliance Area Selection' section is visible, featuring a table of failure modes and compliance areas, all of which are checked. At the bottom, there are buttons for 'Settings filename (e.g. setti)', 'Save Settings', and 'Load Settings'.

Choose Phase(s):  
× Training & Validation ×

Choose Task(s):  
Select...  
Bias & Fairness Audits  
Data Versioning & Management  
Deep Learning & ML Frameworks  
Distributed Training & Optimization  
Model Evaluation & Explainability  
Performance Metrics & Validation

Choose Tool(s):  
Select...  
Assessment  
Assessment Summary

**Failure Mode and Compliance Area Selection**

Select	Failure Mode
<input checked="" type="checkbox"/>	<b>AI Regulation</b>
<input checked="" type="checkbox"/>	Lack of AI Transparency Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input checked="" type="checkbox"/>	Lack of Ethical Safeguards in AI Failure to ensure ethical principles (e.g., non-harm, fairness) in AI decision-making.
<input checked="" type="checkbox"/>	Lack of Auditability & Documentation Inability to provide detailed logs of data usage, model training, validation etc. hindering compliance audits.
<input checked="" type="checkbox"/>	<b>Data Protection</b>
<input checked="" type="checkbox"/>	Personal Data Leakage ML models unintentionally exposing personal data during training or inference.
<input checked="" type="checkbox"/>	<b>Functional Safety</b>
<input checked="" type="checkbox"/>	Lack of Real-Time Processing Failure to meet real-time inference requirements in autonomous or safety-critical applications.
<input checked="" type="checkbox"/>	Failure Mode Undetection Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.
<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility Non-deterministic model outputs lead to unpredictable AI system behavior.
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging Poor logging and monitoring hinder failure analysis and root cause identification.
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

Settings filename (e.g. setti) Save Settings Load Settings

Note:

Selecting the phases automatically filters the fault conditions.

# AssessML: selection area

1st approach: Phases and tasks

Issuing of the final selection

Choose Phase(s):  Choose Task(s):  Choose Tool(s):

Failure Mode Selection	Task Assessment	Tool Assessment	Assessment Summary
------------------------	-----------------	-----------------	--------------------

## Failure Mode and Compliance Area Selection

Select	Failure Mode	Description
<input checked="" type="checkbox"/>	<b>AI Regulation</b>	
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input checked="" type="checkbox"/>	<b>Data Protection</b>	
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.
<input checked="" type="checkbox"/>	<b>Functional Safety</b>	
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

Settings filename (e.g. settii)

After going through the 'Phases' and 'Tasks' tabs from left to right, the dashboard displays all potential error states that remain after filtering. → The possible error states are restricted!

# AssessML: selection area

## 2nd approach: Tools

Start: Selection of tools used

Choose Phase(s):  Choose Task(s):  Choose Tool(s):

Failure Mode Selection Task Assessment Tool Assessment

### Failure Mode and Compliance Area Selection

Select	Failure Mode	Description
<input checked="" type="checkbox"/>	<b>AI Regulation</b>	
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input checked="" type="checkbox"/>	Lack of Model Risk Assessment	Failure to categorize AI models under AI Act risk levels and document safety concerns.
<input checked="" type="checkbox"/>	Security Vulnerabilities in AI	Exposure to adversarial attacks compromising AI safety and trustworthiness.
<input checked="" type="checkbox"/>	Non-compliant Human Oversight Mechanisms	Failure to implement necessary human-in-the-loop controls in high-risk AI applications.
<input checked="" type="checkbox"/>	Insufficient AI Impact Assessment	Failure to conduct proper risk-benefit analysis for high-risk AI applications.
<input checked="" type="checkbox"/>	Lack of Ethical Safeguards in AI	Failure to ensure ethical principles (e.g., non-harm, fairness) in AI decision-making.
<input checked="" type="checkbox"/>	Inadequate Monitoring of Deployed AI Systems	Lack of continuous assessment of AI systems in real-world settings.
<input checked="" type="checkbox"/>	Lack of Auditability & Documentation	Inability to provide detailed logs of data usage, model training, validation etc. hindering compliance audits.
<input checked="" type="checkbox"/>	<b>Data Protection</b>	
<input checked="" type="checkbox"/>	Unauthorized Data Access	Lack of strict access control leading to potential data breaches and privacy violations.
<input checked="" type="checkbox"/>	Failure to Handle Right to Erasure	Non-compliance with GDPR's Right to be Forgotten, leading to legal risks.
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.
<input checked="" type="checkbox"/>	Lack of Data Minimization	Collecting and storing excessive personal data, violating GDPR principles.
<input checked="" type="checkbox"/>	Failure in Data Anonymization & Pseudonymization	Inadequate de-identification techniques lead to re-identifiable personal data.
<input checked="" type="checkbox"/>	Non-Transparent User Consent Management	Lack of clear user consent tracking for data collection and processing.
<input checked="" type="checkbox"/>	Cross-border Data Transfer Violations	Failure to comply with international data transfer rules under GDPR/Schrems II.
<input checked="" type="checkbox"/>	<b>Functional Safety</b>	
<input checked="" type="checkbox"/>	Data Integrity Failure	Corrupt or inconsistent data leading to incorrect model predictions, endangering safety-critical systems.
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.
<input checked="" type="checkbox"/>	Failure Mode Undetection	Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.
<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility	Non-deterministic model outputs lead to unpredictable AI system behavior.
<input checked="" type="checkbox"/>	Insufficient Edge Case Handling	AI models fail to correctly handle rare but critical safety scenarios, leading to unreliable operation.
<input checked="" type="checkbox"/>	Undetected Model Performance Drift	Model performance degrades over time without triggering corrective actions, leading to unsafe decisions.
<input checked="" type="checkbox"/>	Insufficient Redundancy & Failover Mechanisms	Lack of backup AI models or failover strategies in critical systems.
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

TensorFlow Serving  
TensorRT  
Weights & Biases  
WhyLabs  
**XGBoost**  
tsfresh

Settings filename (e.g. setti) Save Settings Load Settings

# AssessML: selection area

2nd approach: Tools

Automatic output of assigned phases and tasks

Choose Phase(s):  Choose Task(s):  Choose Tool(s):

Failure Mode Selection	Task Assessment	Tool Assessment	Assessment Summary
------------------------	-----------------	-----------------	--------------------

## Failure Mode and Compliance Area Selection

Select	Failure Mode	Description
<input checked="" type="checkbox"/>	<b>AI Regulation</b>	
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input checked="" type="checkbox"/>	<b>Data Protection</b>	
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.
<input checked="" type="checkbox"/>	<b>Functional Safety</b>	
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

Settings filename (e.g. setti)

After going through the 'Tools' tab, the dashboard displays all error states associated with the selected tool.

→ There is a specific restriction of error states!

# AssessML: selection area

## Note

Multiple selections are also possible; this applies to **both** approaches!

### 1st approach: Phases and tasks

If several phases and tasks are selected, the error states that are relevant to all selected phases and tasks are displayed.

Failure Mode Selection	Task Assessment	Tool Assessment	Assessment Summary
<b>Failure Mode and Compliance Area Selection</b>			
Select	Failure Mode	Description	
<input checked="" type="checkbox"/>	<b>AI Regulation</b>		
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.	
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.	
<input checked="" type="checkbox"/>	Security Vulnerabilities in AI	Exposure to adversarial attacks compromising AI safety and trustworthiness.	
<input checked="" type="checkbox"/>	<b>Data Protection</b>		
<input checked="" type="checkbox"/>	Unauthorized Data Access	Lack of strict access control leading to potential data breaches and privacy violations.	
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.	
<input checked="" type="checkbox"/>	Failure in Data Anonymization & Pseudonymization	Inadequate de-identification techniques lead to re-identifiable personal data.	
<input checked="" type="checkbox"/>	Cross-border Data Transfer Violations	Failure to comply with international data transfer rules under GDPR Schema II.	
<input checked="" type="checkbox"/>	<b>Functional Safety</b>		
<input checked="" type="checkbox"/>	Data Integrity Failure	Corrupt or inconsistent data leading to incorrect model predictions, endangering safety-critical systems.	
<input checked="" type="checkbox"/>	Failure Mode Undetection	Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.	
<input checked="" type="checkbox"/>	Insufficient Edge Case Handling	AI models fail to correctly handle rare but critical safety scenarios, leading to unreliable operation.	
<input checked="" type="checkbox"/>	Insufficient Redundancy & Failover Mechanisms	Lack of backup AI models or failover strategies in critical systems.	
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.	

There is a restriction on the possible error states!

### 2. Approach: Tools

If several tools are selected, the error states that are relevant to all selected tools are displayed.

Failure Mode Selection	Task Assessment	Tool Assessment	Assessment Summary
<b>Failure Mode and Compliance Area Selection</b>			
Select	Failure Mode	Description	
<input checked="" type="checkbox"/>	<b>AI Regulation</b>		
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.	
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.	
<input checked="" type="checkbox"/>	<b>Data Protection</b>		
<input checked="" type="checkbox"/>	Unauthorized Data Access	Lack of strict access control leading to potential data breaches and privacy violations.	
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.	
<input checked="" type="checkbox"/>	Lack of Data Minimization	Collecting and storing excessive personal data, violating GDPR principles.	
<input checked="" type="checkbox"/>	Non-Transparent User Consent Management	Lack of clear user consent tracking for data collection and processing.	
<input checked="" type="checkbox"/>	<b>Functional Safety</b>		
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.	
<input checked="" type="checkbox"/>	Failure Mode Undetection	Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.	
<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility	Non-deterministic model outputs lead to unpredictable AI system behavior.	
<input checked="" type="checkbox"/>	Insufficient Edge Case Handling	AI models fail to correctly handle rare but critical safety scenarios, leading to unreliable operation.	
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.	
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.	

There is a specific restriction of error states!

# 1. Identification of risks in development tasks

## 1.2 Selection of area of compliance and applicable failure modes

1. Identification of risks in development tasks (according to FMEA)

$$\text{RPN} = \text{Severity} * \text{Occurrence} * \text{Detection}$$

1.1 Selection of relevant tasks  
phases

1.2 Selection of area of compliance  
and applicable failure modes

1.3 Calculation of the RPN for  
selected failure modes in a task

# AssessML: selection of failure modes

## Structure



### Failure Mode and Compliance Area Selection

Select	Failure Mode	Description
<input checked="" type="checkbox"/>	<b>AI Regulation</b>	
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input checked="" type="checkbox"/>	Lack of Model Risk Assessment	Failure to categorize AI models under AI Act risk levels and document safety concerns.
<input checked="" type="checkbox"/>	Security Vulnerabilities in AI	Exposure to adversarial attacks compromising AI safety and trustworthiness.
<input checked="" type="checkbox"/>	Non-compliant Human Oversight Mechanisms	Failure to implement necessary human-in-the-loop controls in high-risk AI applications.
<input checked="" type="checkbox"/>	Insufficient AI Impact Assessment	Failure to conduct proper risk-benefit analysis for high-risk AI applications.
<input checked="" type="checkbox"/>	Lack of Ethical Safeguards in AI	Failure to ensure ethical principles (e.g., non-harm, fairness) in AI decision-making.
<input checked="" type="checkbox"/>	Inadequate Monitoring of Deployed AI Systems	Lack of continuous assessment of AI systems in real-world settings.
<input checked="" type="checkbox"/>	Lack of Auditability & Documentation	Inability to provide detailed logs of data usage, model training, validation etc. hindering compliance audits.
<input checked="" type="checkbox"/>	<b>Data Protection</b>	
<input checked="" type="checkbox"/>	Unauthorized Data Access	Lack of strict access control leading to potential data breaches and privacy violations.
<input checked="" type="checkbox"/>	Failure to Handle Right to Erasure	Non-compliance with GDPR's Right to be Forgotten, leading to legal risks.
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.
<input checked="" type="checkbox"/>	Lack of Data Minimization	Collecting and storing excessive personal data, violating GDPR principles.
<input checked="" type="checkbox"/>	Failure in Data Anonymization & Pseudonymization	Inadequate de-identification techniques lead to re-identifiable personal data.
<input checked="" type="checkbox"/>	Non-Transparent User Consent Management	Lack of clear user consent tracking for data collection and processing.
<input checked="" type="checkbox"/>	Cross-border Data Transfer Violations	Failure to comply with international data transfer rules under GDPR/Schrems II.
<input checked="" type="checkbox"/>	<b>Functional Safety</b>	
<input checked="" type="checkbox"/>	Data Integrity Failure	Corrupt or inconsistent data leading to incorrect model predictions, endangering safety-critical systems.
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.
<input checked="" type="checkbox"/>	Failure Mode Undetection	Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.
<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility	Non-deterministic model outputs lead to unpredictable AI system behavior.
<input checked="" type="checkbox"/>	Insufficient Edge Case Handling	AI models fail to correctly handle rare but critical safety scenarios, leading to unreliable operation.
<input checked="" type="checkbox"/>	Undetected Model Performance Drift	Model performance degrades over time without triggering corrective actions, leading to unsafe decisions.
<input checked="" type="checkbox"/>	Insufficient Redundancy & Failover Mechanisms	Lack of backup AI models or failover strategies in critical systems.
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

The assessment section is divided into four slides:

- **Failure Mode Selection** (List of error conditions),
- **Task Assessment** (Identification of risks in development tasks),
- **Tool Assessment** (Training requirements for AI tools) and
- **Assessment Summary** (Qualification summary based on the RPN and the TCL)

# AssessML: selection of failure modes

## Fine selection - Failure Mode Overview

Failure Mode Selection	Task Assessment	Tool Assessment	Assessment Summary
------------------------	-----------------	-----------------	--------------------

### Failure Mode and Compliance Area Selection

Select	Failure Mode	Description
<input checked="" type="checkbox"/>	<b>AI Regulation</b>	
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input checked="" type="checkbox"/>	Lack of Model Risk Assessment	Failure to categorize AI models under AI Act risk levels and document safety concerns.
<input checked="" type="checkbox"/>	Security Vulnerabilities in AI	Exposure to adversarial attacks compromising AI safety and trustworthiness.
<input checked="" type="checkbox"/>	Non-compliant Human Oversight Mechanisms	Failure to implement necessary human-in-the-loop controls in high-risk AI applications.
<input checked="" type="checkbox"/>	Insufficient AI Impact Assessment	Failure to conduct proper risk-benefit analysis for high-risk AI applications.
<input checked="" type="checkbox"/>	Lack of Ethical Safeguards in AI	Failure to ensure ethical principles (e.g., non-harm, fairness) in AI decision-making.
<input checked="" type="checkbox"/>	Inadequate Monitoring of Deployed AI Systems	Lack of continuous assessment of AI systems in real-world settings.
<input checked="" type="checkbox"/>	Lack of Auditability & Documentation	Inability to provide detailed logs of data usage, model training, validation etc. hindering compliance audits.
<input checked="" type="checkbox"/>	<b>Data Protection</b>	
<input checked="" type="checkbox"/>	Unauthorized Data Access	Lack of strict access control leading to potential data breaches and privacy violations.
<input checked="" type="checkbox"/>	Failure to Handle Right to Erasure	Non-compliance with GDPR's Right to be Forgotten, leading to legal risks.
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.
<input checked="" type="checkbox"/>	Lack of Data Minimization	Collecting and storing excessive personal data, violating GDPR principles.
<input checked="" type="checkbox"/>	Failure in Data Anonymization & Pseudonymization	Inadequate de-identification techniques lead to re-identifiable personal data.
<input checked="" type="checkbox"/>	Non-Transparent User Consent Management	Lack of clear user consent tracking for data collection and processing.
<input checked="" type="checkbox"/>	Cross-border Data Transfer Violations	Failure to comply with international data transfer rules under GDPR/Schrems II.
<input checked="" type="checkbox"/>	<b>Functional Safety</b>	
<input checked="" type="checkbox"/>	Data Integrity Failure	Corrupt or inconsistent data leading to incorrect model predictions, endangering safety-critical systems.
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.
<input checked="" type="checkbox"/>	Failure Mode Undetection	Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.
<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility	Non-deterministic model outputs lead to unpredictable AI system behavior.
<input checked="" type="checkbox"/>	Insufficient Edge Case Handling	AI models fail to correctly handle rare but critical safety scenarios, leading to unreliable operation.
<input checked="" type="checkbox"/>	Undetected Model Performance Drift	Model performance degrades over time without triggering corrective actions, leading to unsafe decisions.
<input checked="" type="checkbox"/>	Insufficient Redundancy & Failover Mechanisms	Lack of backup AI models or failover strategies in critical systems.
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

The **Failure Mode Selection** section shows:

- List of all error states,
- Explanation for each error condition

and is used to select the desired/required error states of the three compatibility areas.

# AssessML: selection of failure modes

## Fine selection - Failure Mode Overview

For a detailed search, the user can apply further filters using the 'Select boxes':

- If you make a selection for the individual **failure modes**, the modes that are not taken into account are highlighted in **white**.
- If a selection is made for an entire **compatibility area**, the area that is not taken into account is highlighted in **grey**.

Select the failure modes by clicking on the 'Select' boxes for each failure mode and/or for each compliance area.

Select	Failure Mode	Select	Failure Mode
<input checked="" type="checkbox"/> <b>AI Regulation</b>		<input type="checkbox"/> <b>AI Regulation</b>	
<input checked="" type="checkbox"/>	Lack of AI Transparency	<input checked="" type="checkbox"/>	Lack of AI Transparency
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance
<input checked="" type="checkbox"/>	Lack of Model Risk Assessment	<input checked="" type="checkbox"/>	Lack of Model Risk Assessment
<input checked="" type="checkbox"/>	Security Vulnerabilities in AI	<input checked="" type="checkbox"/>	Security Vulnerabilities in AI
<input type="checkbox"/>	Non-compliant Human Oversight Mechanisms	<input checked="" type="checkbox"/>	Non-compliant Human Oversight Mechanisms
<input type="checkbox"/>	Insufficient AI Impact Assessment	<input checked="" type="checkbox"/>	Insufficient AI Impact Assessment
<input checked="" type="checkbox"/>	Lack of Ethical Safeguards in AI	<input checked="" type="checkbox"/>	Lack of Ethical Safeguards in AI
<input checked="" type="checkbox"/>	Inadequate Monitoring of Deployed AI Systems	<input checked="" type="checkbox"/>	Inadequate Monitoring of Deployed AI Systems
<input type="checkbox"/>	Lack of Auditability & Documentation	<input checked="" type="checkbox"/>	Lack of Auditability & Documentation
<input checked="" type="checkbox"/> <b>Data Protection</b>		<input checked="" type="checkbox"/> <b>Data Protection</b>	
<input checked="" type="checkbox"/>	Unauthorized Data Access	<input checked="" type="checkbox"/>	Unauthorized Data Access
<input type="checkbox"/>	Failure to Handle Right to Erasure	<input checked="" type="checkbox"/>	Failure to Handle Right to Erasure
<input checked="" type="checkbox"/>	Personal Data Leakage	<input checked="" type="checkbox"/>	Personal Data Leakage
<input checked="" type="checkbox"/>	Lack of Data Minimization	<input checked="" type="checkbox"/>	Lack of Data Minimization
<input type="checkbox"/>	Failure in Data Anonymization & Pseudonymization	<input checked="" type="checkbox"/>	Failure in Data Anonymization & Pseudonymization
<input type="checkbox"/>	Non-Transparent User Consent Management	<input checked="" type="checkbox"/>	Non-Transparent User Consent Management
<input checked="" type="checkbox"/>	Cross-border Data Transfer Violations	<input checked="" type="checkbox"/>	Cross-border Data Transfer Violations
<input checked="" type="checkbox"/> <b>Functional Safety</b>		<input type="checkbox"/> <b>Functional Safety</b>	
<input checked="" type="checkbox"/>	Data Integrity Failure	<input checked="" type="checkbox"/>	Data Integrity Failure
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	<input checked="" type="checkbox"/>	Lack of Real-Time Processing
<input type="checkbox"/>	Failure Mode Undetection	<input checked="" type="checkbox"/>	Failure Mode Undetection
<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility	<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility

# 1. Identification of risks in development tasks

## 1.3 Calculation of the RPN for selected failure modes in a task

1. Identification of risks in development tasks (according to FMEA)

$$\text{RPN} = \text{Severity} * \text{Occurrence} * \text{Detection}$$

1.1 Selection of relevant tasks  
phases

1.2 Selection of area of compliance  
and applicable failure modes

1.3 Calculation of the RPN for  
selected failure modes in a task

# AssessML: evaluation of failure modes

## Task Assessment - Calculation of the RPN for selected failure modes

The **Task Assessment** area is used to calculate the risk priority number (RPN) for each selected fault condition and/or compatibility area.



### Failure Mode Analysis for Individual ML Task

#### Phase: Training & Validation

##### Task: Deep Learning & ML Frameworks

*Optimizing model parameters based on training data.*

Artifact: Trained model

#### AI Regulation

**Failure mode:** Lack of AI Transparency

**Impact:** Failure to provide interpretability and documentation of model decisions.

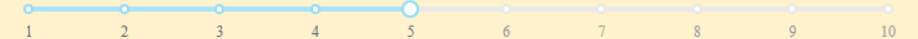
**Potential Causes:**

- No built-in interpretability or explainability methods.
- Poor model decision-tracking and logging support.
- Limited compatibility with external explainability libraries

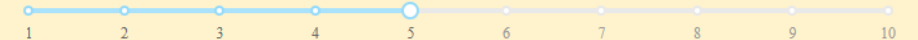
**Possible Mitigations:**

- Utilize ML frameworks that offer native support for model interpretation techniques like SHAP and LIME.
- Establish structured logging of intermediate decisions and model outputs for traceability.
- Ensure framework compatibility with transparency libraries via plugin architecture or open APIs.

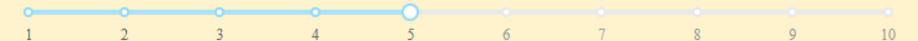
Severity of failure mode



Occurrence of failure mode



Detection of failure mode



RPN:

# AssessML: evaluation of failure modes

## Task Assessment - Calculation of the RPN for selected failure modes

The **RPN** is the sum of the three indicators:

- **Severity** (Significance of the error sequence),
- **Occurrence** (Probability of occurrence of the cause of the fault) and
- **Detection** (Probability of discovering the error or its cause)

$$\text{RPN} = \text{Severity} * \text{Occurrence} * \text{Detection}$$

### Phase: Training & Validation

#### Task: Deep Learning & ML Frameworks

*Optimizing model parameters based on training data.*

**Artifact:** Trained model

### Risk of Deep Learning & ML Frameworks in AI Regulation

**Failure mode:** Lack of AI Transparency: Failure to provide interpretability and documentation of model decisions.

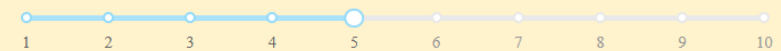
#### Potential Causes:

- No built-in interpretability or explainability methods.
- Poor model decision-tracking and logging support.
- Limited compatibility with external explainability libraries

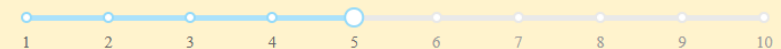
#### Possible Mitigations:

- Utilize ML frameworks that offer native support for model interpretation techniques like SHAP and LIME.
- Establish structured logging of intermediate decisions and model outputs for traceability.
- Ensure framework compatibility with transparency libraries via plugin architecture or open APIs.

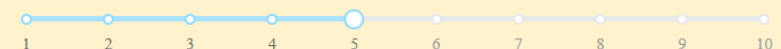
#### Severity of failure mode



#### Occurrence of failure mode



#### Detection of failure mode



**RPN:**

# AssessML: evaluation of failure modes

Significance of the failure mode (Severity)

## Action:

Assess the **severity** for the failure modes assigned to the selected ML tasks in the tool.

### Phase: Training & Validation

#### Task: Deep Learning & ML Frameworks

*Optimizing model parameters based on training data.*

Artifact: Trained model

### Risk of Deep Learning & ML Frameworks in AI Regulation

**Failure mode:** Lack of AI Transparency: Failure to provide interpretability and documentation of model decisions.

#### Potential Causes:

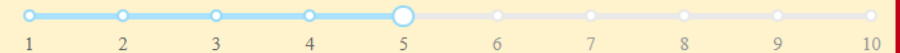
- No built-in interpretability or explainability methods.
- Poor model decision-tracking and logging support.
- Limited compatibility with external explainability libraries

#### Possible Mitigations:

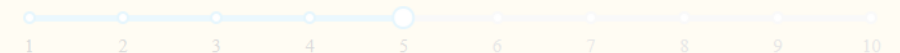
- Utilize ML frameworks that offer native support for model interpretation techniques like SHAP and LIME.
- Establish structured logging of intermediate decisions and model outputs for traceability.
- Ensure framework compatibility with transparency libraries via plugin architecture or open APIs.

By moving the slider, the values  
1 – 3; 4 – 6; 7 – 10 can be set.

#### Severity of failure mode



#### Occurrence of failure mode



#### Detection of failure mode



**RPN:**

# AssessML: evaluation of failure modes

Probability of occurrence of the cause of the failure mode (Occurrence)

## Action:

Assess the **probability of occurrence** for the failure modes assigned to the selected ML tasks in the tool.

### Phase: Training & Validation

#### Task: Deep Learning & ML Frameworks

*Optimizing model parameters based on training data.*

Artifact: Trained model

### Risk of Deep Learning & ML Frameworks in AI Regulation

**Failure mode:** Lack of AI Transparency: Failure to provide interpretability and documentation of model decisions.

#### Potential Causes:

- No built-in interpretability or explainability methods.
- Poor model decision-tracking and logging support.
- Limited compatibility with external explainability libraries

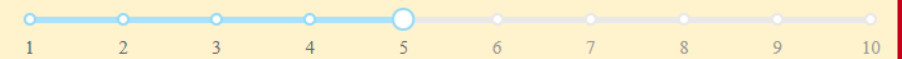
#### Possible Mitigations:

- Utilize ML frameworks that offer native support for model interpretation techniques like SHAP and LIME.
- Establish structured logging of intermediate decisions and model outputs for traceability.
- Ensure framework compatibility with transparency libraries via plugin architecture or open APIs.

Severity of failure



Occurrence of failure mode



Detection of failure mode



RPN:

# AssessML: evaluation of failure modes

Probability of discovering the failure mode or its cause (Detectability)

## Action:

Assess the **detectability** for the failure modes assigned to the selected ML tasks in the tool.

### Phase: Training & Validation

#### Task: Deep Learning & ML Frameworks

*Optimizing model parameters based on training data.*

Artifact: Trained model

### Risk of Deep Learning & ML Frameworks in AI Regulation

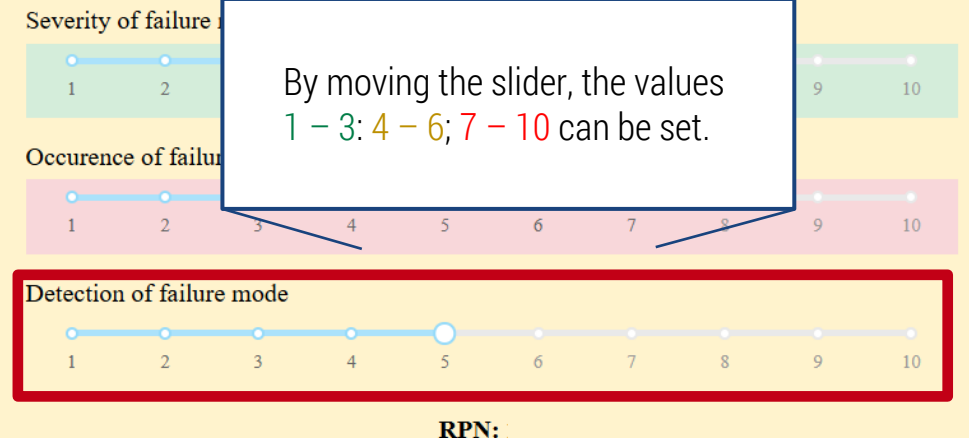
**Failure mode:** Lack of AI Transparency: Failure to provide interpretability and documentation of model decisions.

#### Potential Causes:

- No built-in interpretability or explainability methods.
- Poor model decision-tracking and logging support.
- Limited compatibility with external explainability libraries

#### Possible Mitigations:

- Utilize ML frameworks that offer native support for model interpretation techniques like SHAP and LIME.
- Establish structured logging of intermediate decisions and model outputs for traceability.
- Ensure framework compatibility with transparency libraries via plugin architecture or open APIs.



# AssessML: evaluation of failure modes

## Calculation of the RPN

### Action:

Calculation of the **RPN** for the failure modes assigned to the selected ML tasks in the tool.

#### Phase: Training & Validation

##### Task: Deep Learning & ML Frameworks

*Optimizing model parameters based on training data.*

Artifact: Trained model

#### Risk of Deep Learning & ML Frameworks in AI Regulation

**Failure mode:** Lack of AI Transparency: Failure to provide interpretability and documentation of model decisions.

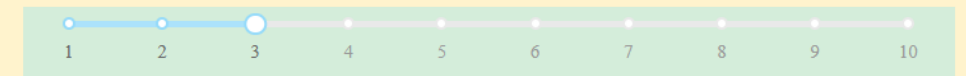
##### Potential Causes:

- No built-in interpretability or explainability methods.
- Poor model decision-tracking and logging support.
- Limited compatibility with external explainability libraries

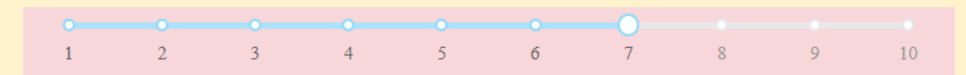
##### Possible Mitigations:

- Utilize ML frameworks that offer native support for model interpretation techniques like SHAP and LIME.
- Establish structured logging of intermediate decisions and model outputs for traceability.
- Ensure framework compatibility with transparency libraries via plugin architecture or open APIs.

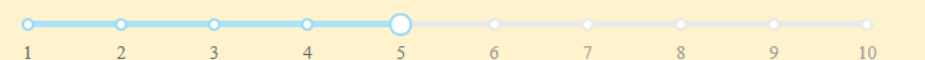
##### Severity of failure mode



##### Occurrence of failure mode



##### Detection of failure mode



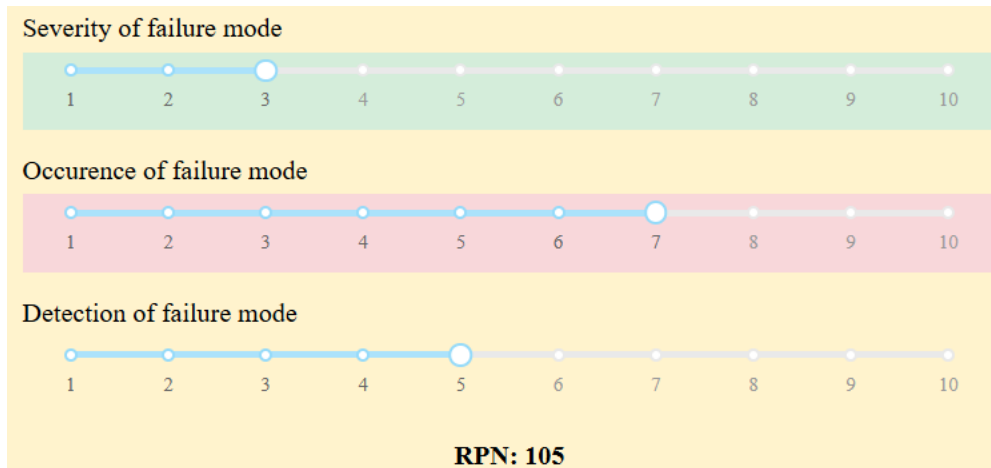
Calculation of the RPN (1 – 29; 30 – 299; 300 -1000).

**RPN: 105**

# AssessML: evaluation of failure modes

## Task Assessment - Calculation of the RPN for selected failure modes in a task

$$\text{RPN} = \text{Severity} * \text{Occurrence} * \text{Detection}$$



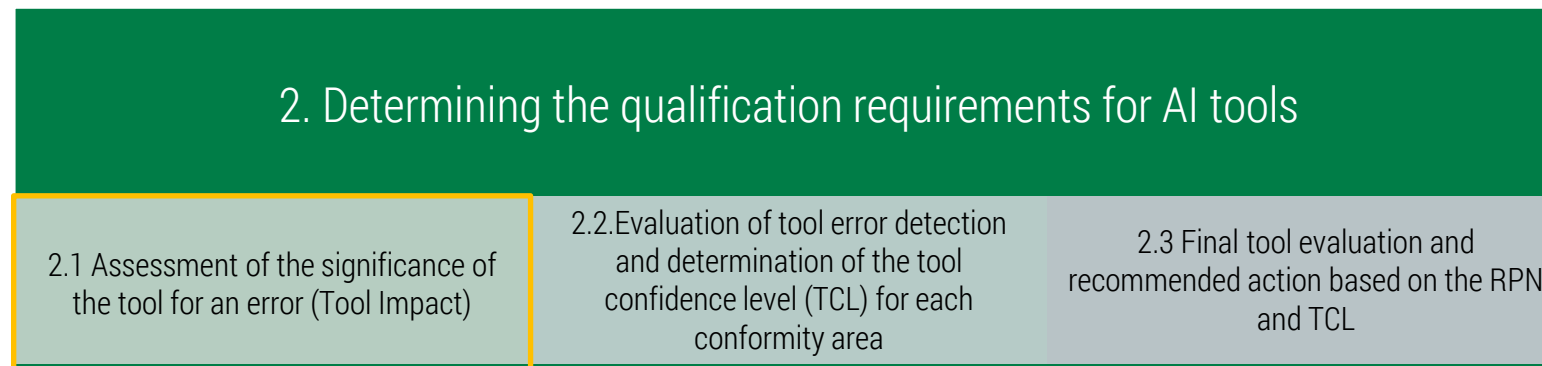
### Evaluation of the risk priority number (RPN):

- Low RPN values (RPN 1–30) indicate low criticality and no need for action
- Medium RPN values (RPN 31–299) indicate medium criticality and thus a need for action
- High values (RPN 300–1000) indicate high criticality and a greater need for action.

→ The determination of risks (RPN) in the development tasks (according to FMEA) focuses on the determination per use case

## 2. Determining the qualification requirements for AI tools

### 2.1 Assessment of the significance of the tool for an error (Tool Impact)



# AssessML: impact of a tool failure

Tool Assessment - Identifying training needs for AI tools

The **Tool Assessment** area is used to determine the qualification requirements for AI tools for each selected error state and/or compatibility area.



## Tool Impact and Tool Error Detection Analysis

### Phase: Training & Validation

#### Task: Deep Learning & ML Frameworks

*Optimizing model parameters based on training data.*

Artifact: Trained model

### AI Regulation

#### Tool Confidence Level:

Failure mode: Lack of AI Transparency

RPN: 105

Potential tool failure: May fail to generate or link detailed training documentation and validation reports, resulting in trained models lacking required transparency and conformity evidence.

#### Tool Impact Value



#### Tool Error Detection



# AssessML: impact of a tool failure

Effectiveness of the tool



## Tool Impact and Tool Error Detection Analysis

### Phase: Training & Validation

#### Task: Deep Learning & ML Frameworks

Optimizing model parameters based on training data.

Artifact: Trained model

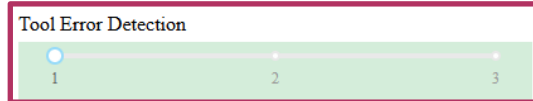
### AI Regulation

#### Tool Confidence Level:

Failure mode: Lack of AI Transparency

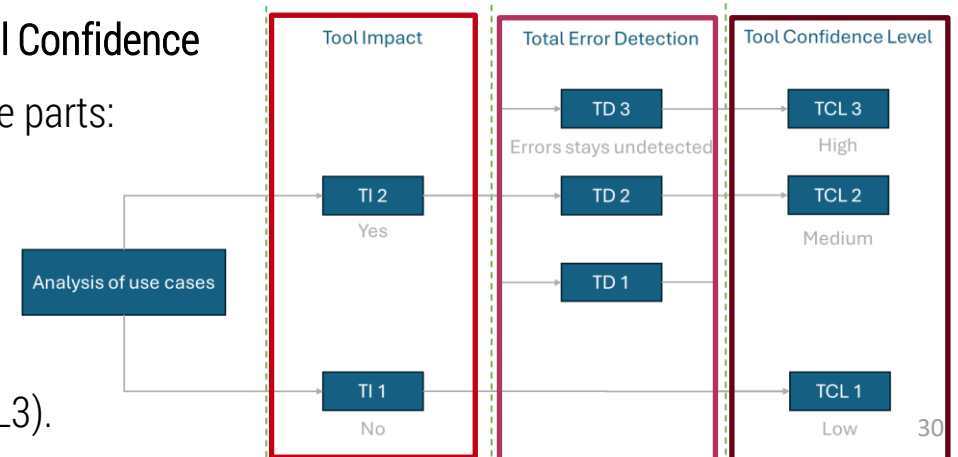
RPN: 105

Potential tool failure: May fail to generate or link detailed training documentation and validation reports, resulting in trained models lacking required transparency and conformity evidence.



Based on section 11.4.5.2 of ISO 26262:2018, the Tool Confidence Level (TCL) is determined based on the following three parts:

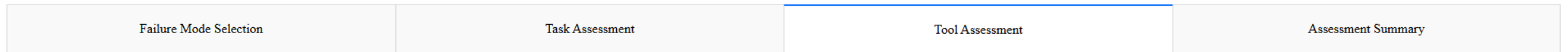
- **Tool Impact (TI)**, assesses whether a tool error could be safety-relevant,
- **Tool Error Detection (TD)**, assesses the probability that an error will be detected and results in the
- **Tool Confidence Level (TCL)**, the required level of confidence in the tool (TCL1 to TCL3).



# AssessML: impact of a tool failure

Assessment of the significance of the tool for an error (Tool Impact)

**Action:** Does a potential error in the tool affect the functionality of safety-related systems?



## Tool Impact and Tool Error Detection Analysis

### Phase: Training & Validation

#### Task: Deep Learning & ML Frameworks

*Optimizing model parameters based on training data.*

Artifact: Trained model

### AI Regulation

#### Tool Confidence Level:

Failure mode: Lack of AI Transparency

RPN: 105

Potential tool failure: May fail to generate or link detailed training documentation and conformity evidence.

Tool Impact Value

TI1  TI2

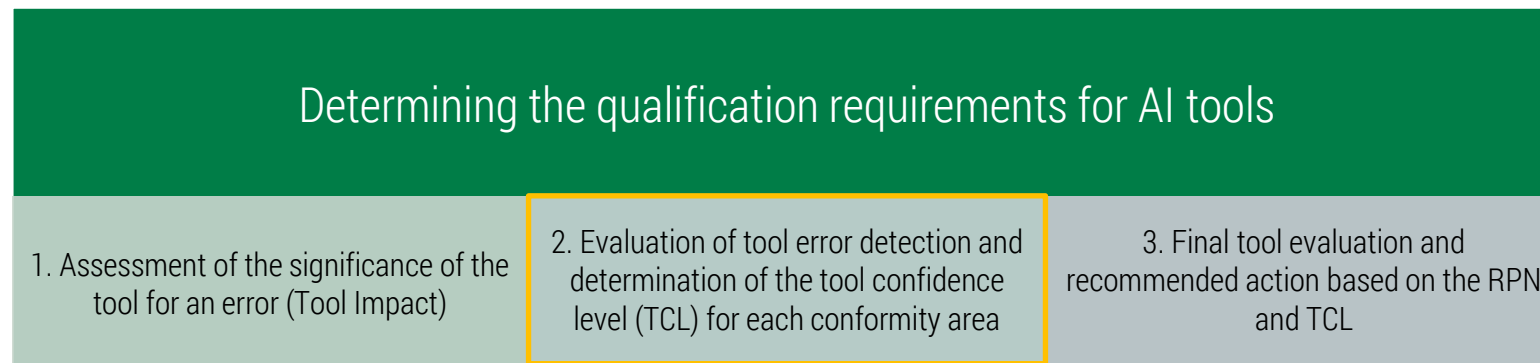
Moving the slider allows you to assess the presence of safety-related effects (1 (yes), 2 (no))

Tool Error Detection

1  2  3

## 2. Determining the qualification requirements for AI tools

### 2.2. Evaluation of tool error detection and TCL calculation



# AssessML: detectability of tool error

Assessment of tool error detectability (Tool Error Detection)

**Action:** How well can errors caused by the tool be detected in the further development process?



## Tool Impact and Tool Error Detection Analysis

### Phase: Training & Validation

#### Task: Deep Learning & ML Frameworks

*Optimizing model parameters based on training data.*

Artifact: Trained model

### AI Regulation

#### Tool Confidence Level:

Failure mode: Lack of AI Transparency

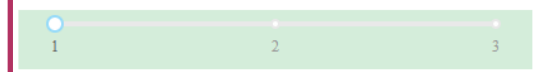
RPN: 105

Potential tool failure: May fail to generate or link detailed training documentation and validation reports, resulting in trained models lacking required transparency and conformity evidence.

Tool Impact Value



Tool Error Detection



Moving the slider allows you to assess the likelihood of security-related impacts  
(1 = wird erkannt, 2 = könnten unentdeckt bleiben, 3 = bleiben unentdeckt)

# Background: calculation of TCL

Need for trust in tool usage

**Starting point:** How high must the level of trust in the tool be (or is tool qualification necessary)?

**Table 3 — Determination of the Tool Confidence Level (TCL)**

		Tool error detection		
		TD1	TD2	TD3
Tool impact	TI1	TCL1	TCL1	TCL1
	TI2	TCL1	TCL2	TCL3

Source: ISO 26262:2018, Section. 11.4.5.4



Aspect	Picture	Description
Target	<i>Assignment of the appropriate Tool Confidence Level (TCL) based on Tool Impact (TI) and Tool Error Detection (TD)</i>	Determine <b>how high the level of trust</b> in the tool must be (or whether tool qualification is necessary)
TCL1 (low)	<ul style="list-style-type: none"> <li>Tool has no safety-related influence (TI1) or</li> <li>Errors are reliably detected (TI2 + TD1)</li> </ul> <p>→ <b>No tool qualification necessary</b></p>	<ul style="list-style-type: none"> <li><b>Low significance</b> for product quality</li> <li>No tool qualification necessary</li> </ul> <p>→ Confidence in tool behaviour from an ISO 26262 perspective <b>not critical</b></p>
TCL2 (medium)	<ul style="list-style-type: none"> <li>Tool has safety-related influence (TI2)</li> <li>Error <b>not reliably detectable</b> (TD2)</li> </ul>	<ul style="list-style-type: none"> <li>Tool <b>important for product quality</b></li> <li>Tool qualification required</li> <li>Requirements depend on ASIL level</li> </ul>
TCL3 (high)	<ul style="list-style-type: none"> <li>Requirements depend on ASIL level (TI2)</li> <li>Errors <b>remain undetected</b> (TD3)</li> </ul>	<ul style="list-style-type: none"> <li><b>High importance</b> for product quality</li> <li><b>Tool qualification absolutely necessary</b>, high requirements</li> <li>Differences between TCL2 and TCL3 are methodologically relevant but not dramatic (depending on ASIL)</li> </ul>

# Background: calculation of TCL

Need for trust in tool usage

**Starting point:** How high must the level of trust in the tool be (or is tool qualification necessary)?

**Table 3 — Determination of the Tool Confidence Level (TCL)**

		Tool error detection		
		TD1	TD2	TD3
Tool impact	T11	TCL1	TCL1	TCL1
	T12	TCL1	TCL2	TCL3

Source: ISO 26262:2018, Section. 11.4.5.4



Aspect	Picture	Description
Target	<i>Assignment of the appropriate Tool Confidence Level (TCL) based on Tool Impact (TI) and Tool Error Detection (TD)</i>	Determine <b>how high the level of trust</b> in the tool must be (or whether tool qualification is necessary)
TCL1 (low)	<ul style="list-style-type: none"> <li>Tool has no safety-related influence (T11) or</li> <li>Errors are reliably detected (T12 + TD1)</li> </ul> <p>→ <b>No tool qualification necessary</b></p>	<ul style="list-style-type: none"> <li><b>Low significance</b> for product quality</li> <li>No tool qualification necessary</li> </ul> <p>→ Confidence in tool behaviour from an ISO 26262 perspective <b>not critical</b></p>
TCL2 (medium)	<ul style="list-style-type: none"> <li>Tool has safety-related influence (T12)</li> <li>Error not reliably detectable (TD2)</li> </ul>	<ul style="list-style-type: none"> <li>Tool important for product quality</li> <li>Tool qualification required</li> <li>Requirements depend on ASIL level</li> </ul>
TCL3 (high)	<ul style="list-style-type: none"> <li>Requirements depend on ASIL level (T12)</li> <li>Errors remain undetected (TD3)</li> </ul>	<ul style="list-style-type: none"> <li><b>High importance</b> for product quality</li> <li>Tool qualification <b>absolutely necessary</b>, high requirements</li> <li>Differences between TCL2 and TCL3 are methodologically relevant but not dramatic (depending on ASIL)</li> </ul>

# Background: calculation of TCL

Need for trust in tool usage

**Starting point:** How high must the level of trust in the tool be (or is tool qualification necessary)?

**Table 3 — Determination of the Tool Confidence Level (TCL)**

		Tool error detection		
		TD1	TD2	TD3
Tool impact	T11	TCL1	TCL1	TCL1
	T12	TCL1	TCL2	TCL3

Source: ISO 26262:2018, Section. 11.4.5.4



Aspect	Picture	Description
Target	<b>Assignment of the appropriate Tool Confidence Level (TCL) based on Tool Impact (TI) and Tool Error Detection (TD)</b>	Determine <b>how high the level of trust</b> in the tool must be (or whether tool qualification is necessary)
TCL1 (low)	<ul style="list-style-type: none"> <li>Tool has no safety-related influence (T11) or</li> <li>Errors are reliably detected (T12 + TD1)</li> </ul> <p>→ No tool qualification necessary</p>	<ul style="list-style-type: none"> <li>Low significance for product quality</li> <li>No tool qualification necessary</li> </ul> <p>→ Confidence in tool behaviour from an ISO 26262 perspective not critical</p>
TCL2 (medium)	<ul style="list-style-type: none"> <li>Tool has safety-related influence (T12)</li> <li>Error <b>not reliably detectable</b> (TD2)</li> </ul>	<ul style="list-style-type: none"> <li>Tool <b>important for product quality</b></li> <li>Tool qualification required</li> <li>Requirements depend on ASIL level</li> </ul>
TCL3 (high)	<ul style="list-style-type: none"> <li>Requirements depend on ASIL level (T12)</li> <li>Errors <b>remain undetected</b> (TD3)</li> </ul>	<ul style="list-style-type: none"> <li>High importance for product quality</li> <li>Tool qualification <b>absolutely necessary</b>, high requirements</li> <li>Differences between TCL2 and TCL3 are methodologically relevant but not dramatic (depending on ASIL)</li> </ul>

# Background: calculation of TCL

Need for trust in tool usage

**Starting point:** How high must the level of trust in the tool be (or is tool qualification necessary)?

**Table 3 — Determination of the Tool Confidence Level (TCL)**

		Tool error detection		
		TD1	TD2	TD3
Tool impact	T11	TCL1	TCL1	TCL1
	T12	TCL1	TCL2	TCL3

Source: ISO 26262:2018, Section. 11.4.5.4

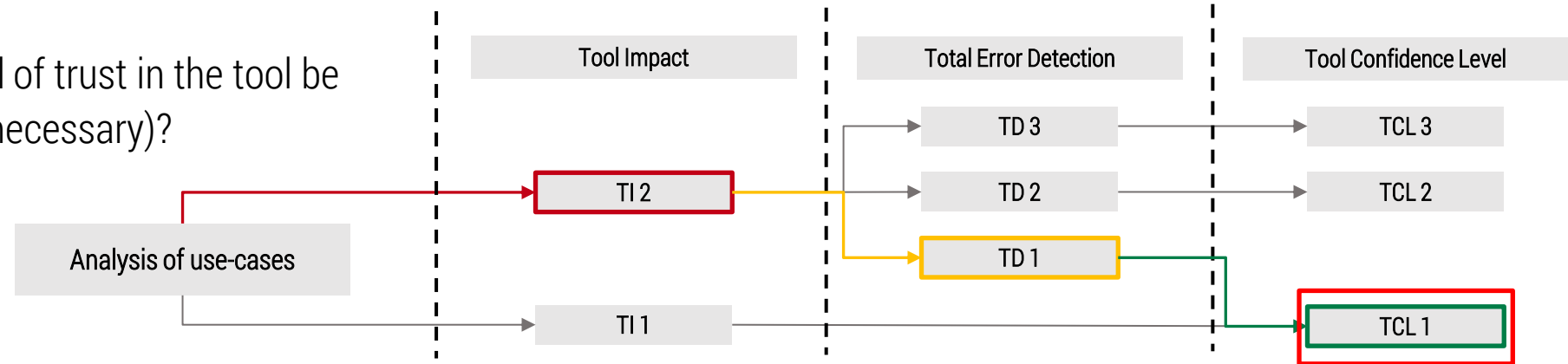


Aspect	Picture	Description
Target	<i>Assignment of the appropriate Tool Confidence Level (TCL) based on Tool Impact (TI) and Tool Error Detection (TD)</i>	Determine <b>how high the level of trust</b> in the tool must be (or whether tool qualification is necessary)
TCL1 (low)	<ul style="list-style-type: none"> <li>Tool has no safety-related influence (T11) or</li> <li>Errors are reliably detected (T12 + TD1)</li> </ul> <p>→ No tool qualification necessary</p>	<ul style="list-style-type: none"> <li>Low significance for product quality</li> <li>No tool qualification necessary</li> </ul> <p>→ Confidence in tool behaviour from an ISO 26262 perspective not critical</p>
TCL2 (medium)	<ul style="list-style-type: none"> <li>Tool has safety-related influence (T12)</li> <li>Error not reliably detectable (TD2)</li> </ul>	<ul style="list-style-type: none"> <li>Tool important for product quality</li> <li>Tool qualification required</li> <li>Requirements depend on ASIL level</li> </ul>
TCL3 (high)	<ul style="list-style-type: none"> <li>Requirements depend on ASIL level (T12)</li> <li>Errors <b>remain undetected</b> (TD3)</li> </ul>	<ul style="list-style-type: none"> <li>High importance for product quality</li> <li>Tool qualification <b>absolutely necessary</b>, high requirements</li> <li>Differences between TCL2 and TCL3 are methodologically relevant but not dramatic (depending on ASIL)</li> </ul>

# Background: calculation of TCL

## Determination of the Tool Confidence Level (TCL)

**Action:** How high must the level of trust in the tool be (or is tool qualification necessary)?



### Tool Impact and Tool Error Detection Analysis

#### Phase: Training & Validation

#### Task: Deep Learning & ML Frameworks

Optimizing model parameters based on training data.

Artifact: Trained model

#### AI Regulation

**Tool Confidence Level: 1**

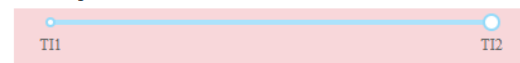
Failure mode: Lack of AI Transparency

RPN: 105

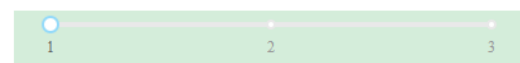
Potential tool failure: May fail to generate or link detailed training documentation and validation reports, resulting in trained models lacking required transparency and conformity evidence.

Calculation of confidence level  
(1 = low, 2 = medium, 3 = high)

#### Tool Impact Value

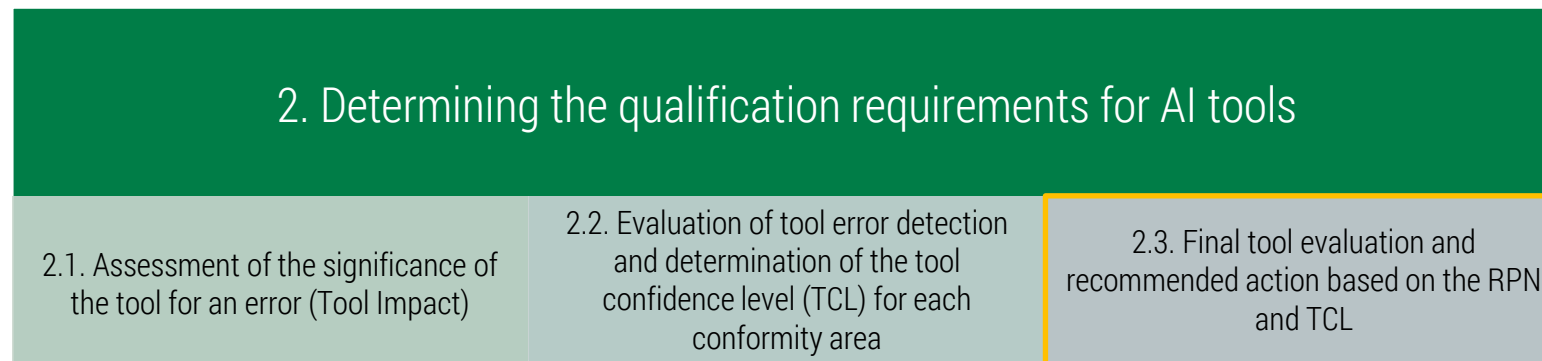


#### Tool Error Detection



## 2. Determining the qualification requirements for AI tools

### 2.3. Tool assessment and recommended actions based on RPN and TCL



# AssessML: evaluation and recommended action

## Assessment Summary - Qualification summary

The **Assessment Summary** section is used to evaluate and prioritise the qualification summary based on the RPN and TCL for selected fault conditions.

Failure Mode Selection	Task Assessment	Tool Assessment	<b>Assessment Summary</b>
------------------------	-----------------	-----------------	---------------------------

### Recommendations for Tool Qualification and Process Improvement

RPN threshold  TCL threshold

#### Phase: Training & Validation

##### Task: Deep Learning & ML Frameworks

*Optimizing model parameters based on training data.*

Artifact: Trained model

#### AI Regulation

Tool Confidence Level: 1

##### Potential tool failure: Lack of AI Transparency

RPN: 105

**Potential tool failure:** May fail to generate or link detailed training documentation and validation reports, resulting in trained models lacking required transparency and conformity evidence.

#### Possible Process Mitigations:

- Utilize ML frameworks that offer native support for model interpretation techniques like SHAP and LIME.
- Establish structured logging of intermediate decisions and model outputs for traceability.
- Ensure framework compatibility with transparency libraries via plugin architecture or open APIs.

#### Testable Tool Capabilities:

- Integrated explainability methods (e.g., SHAP, LIME) within model training environments
- Comprehensive model decision logging and training metadata tracking
- Easy integration with external explainability libraries and visualization tools

# AssessML: evaluation and recommended action

Overview of the specific risk priority score and tool confidence level

The error states can be filtered using the inputs for **thresholds**. Two different specifications can be weighted:

- **RPN threshold** (limit value for the risk priority number)
- **TCL threshold** (limit value for the tool confidence level)

Failure Mode Selection	Task Assessment	Tool Assessment	Assessment Summary
------------------------	-----------------	-----------------	--------------------

## Recommendations for Tool Qualification and Process Improvement

RPN threshold  TCL threshold

### Phase: Training & Validation

#### Task: Deep Learning & ML Frameworks

*Optimizing model parameters based on training data.*

Artifact: Trained model

### AI Regulation

#### Tool Confidence Level: 1

**Potential tool failure: Lack of AI Transparency**

RPN: 105

**Potential tool failure:** May fail to generate or link detailed training documentation and validation reports, resulting in trained models lacking required transparency and conformity evidence.

#### Possible Process Mitigations:

- Utilize ML frameworks that offer native support for model interpretation techniques like SHAP and LIME.
- Establish structured logging of intermediate decisions and model outputs for traceability.
- Ensure framework compatibility with transparency libraries via plugin architecture or open APIs.

#### Testable Tool Capabilities:

- Integrated explainability methods (e.g., SHAP, LIME) within model training environments
- Comprehensive model decision logging and training metadata tracking
- Easy integration with external explainability libraries and visualization tools

# AssessML: evaluation and recommended action

What is the purpose of thresholds?

All statuses are automatically transferred to the assessment summary, which can quickly make the board confusing.

Threshold values can be used to filter out **critical error states** or **unsafe tools** by setting them appropriately. Setting the filter makes the analysis clearer.

RPN threshold 1 TCL threshold 1

Failure Mode Selection Task Assessment Tool Assessment Assessment Summary

Recommendations for Tool Qualification and Process Improvement

RPN threshold 1 TCL threshold 1

Phase: Training & Validation

Task: Deep Learning & ML Frameworks  
Optimizing model parameters based on training data.  
Artifact: Trained model

Functional Safety

Tool Confidence Level: 2

Potential tool failure: Lack of Real-Time Processing  
RPN: 504  
Potential tool failure: May fail to guarantee deterministic training execution and runtime validation, resulting in trained models that do not meet real-time safety-critical performance requirements.

Possible Process Mitigations:

- Apply model optimization techniques such as pruning, quantization, or knowledge distillation.
- Ensure deployment-ready models are benchmarked and optimized for specific hardware accelerators.
- Optimize full pipeline latency by batching operations and reducing unnecessary transformations.

Testable Tool Capabilities:

- Framework support for model optimization (e.g., TensorRT, ONNX Runtime optimizations).
- Support for hardware-specific runtime optimizations and deployment targeting (e.g., EdgeTPU, CUDA kernels).
- Framework support for pipeline fusion, minimal preprocessing APIs, and asynchronous inference.

Tool Confidence Level: 3

Potential tool failure: Inadequate Error Handling & Logging  
RPN: 50

Potential tool failure: May fail to provide comprehensive error logging and runtime monitoring, resulting in trained models with limited traceability for failure analysis in safety-critical systems.

Possible Process Mitigations:

- Integrate structured exception handling and mandatory logging for all critical training and inference operations.
- Implement runtime monitors for numerical instabilities (e.g., NaNs, divergence) and alert on thresholds.
- Embed pipeline-wide metadata capture and systematic experiment logging into the training framework.

Testable Tool Capabilities:

- Framework support for structured error reporting, custom callbacks, and logging integration.
- Built-in hooks for anomaly detection during training and validation (gradient checkers, divergence detectors).
- ML lifecycle management tools (e.g., MLflow, Weights & Biases) integration for metadata and error tracking.

Tool Confidence Level: 1

Potential tool failure: Model Instability & Poor Generalization  
RPN: 60

Potential tool failure: May fail to enforce deterministic training procedures and runtime validation, resulting in trained models with unpredictable behavior and poor generalization in safety-critical applications.

Possible Process Mitigations:

- Implement deterministic training pipelines with seed control and precision configuration to ensure repeatability.
- Integrate continuous performance monitoring during training with automatic validation checkpoints.
- Use frameworks that incorporate gradient clipping, normalization, and regularization for improved numerical robustness.

Testable Tool Capabilities:

- Explicit deterministic training support with reproducibility settings and controlled parallelism.
- Built-in validation dashboards and metric tracking APIs.
- Numerical stability modules with configurable regularization options.



RPN threshold 100 TCL threshold 2

Failure Mode Selection Task Assessment Tool Assessment Assessment Summary

Recommendations for Tool Qualification and Process Improvement

RPN threshold 100 TCL threshold 2

Phase: Training & Validation

Task: Deep Learning & ML Frameworks  
Optimizing model parameters based on training data.  
Artifact: Trained model

Functional Safety

Tool Confidence Level: 2

Potential tool failure: Lack of Real-Time Processing  
RPN: 504  
Potential tool failure: May fail to guarantee deterministic training execution and runtime validation, resulting in trained models that do not meet real-time safety-critical performance requirements.

Possible Process Mitigations:

- Apply model optimization techniques such as pruning, quantization, or knowledge distillation.
- Ensure deployment-ready models are benchmarked and optimized for specific hardware accelerators.
- Optimize full pipeline latency by batching operations and reducing unnecessary transformations.

Testable Tool Capabilities:

- Framework support for model optimization (e.g., TensorRT, ONNX Runtime optimizations).
- Support for hardware-specific runtime optimizations and deployment targeting (e.g., EdgeTPU, CUDA kernels).
- Framework support for pipeline fusion, minimal preprocessing APIs, and asynchronous inference.

Tool Confidence Level: 3

Potential tool failure: Inadequate Error Handling & Logging  
RPN: 50

Potential tool failure: May fail to provide comprehensive error logging and runtime monitoring, resulting in trained models with limited traceability for failure analysis in safety-critical systems.

Possible Process Mitigations:

- Integrate structured exception handling and mandatory logging for all critical training and inference operations.
- Implement runtime monitors for numerical instabilities (e.g., NaNs, divergence) and alert on thresholds.
- Embed pipeline-wide metadata capture and systematic experiment logging into the training framework.

Testable Tool Capabilities:

- Framework support for structured error reporting, custom callbacks, and logging integration.
- Built-in hooks for anomaly detection during training and validation (gradient checkers, divergence detectors).
- ML lifecycle management tools (e.g., MLflow, Weights & Biases) integration for metadata and error tracking.

Tool Confidence Level: 1

Potential tool failure: Model Instability & Poor Generalization  
RPN: 60

Potential tool failure: May fail to enforce deterministic training procedures and runtime validation, resulting in trained models with unpredictable behavior and poor generalization in safety-critical applications.

Possible Process Mitigations:

- Implement deterministic training pipelines with seed control and precision configuration to ensure repeatability.
- Integrate continuous performance monitoring during training with automatic validation checkpoints.
- Use frameworks that incorporate gradient clipping, normalization, and regularization for improved numerical robustness.

Testable Tool Capabilities:

- Explicit deterministic training support with reproducibility settings and controlled parallelism.
- Built-in validation dashboards and metric tracking APIs.
- Numerical stability modules with configurable regularization options.

# AssessML: evaluation and recommended action

Overview of the specific risk priority score and tool confidence level

The **select boxes (Checkboxes)** listed below are used to select specific tool capabilities or possible process measures:

## Recommendations for Tool Qualification and Process Improvement

RPN threshold  TCL threshold

### Phase: Training & Validation

Task: Deep Learning & ML Frameworks

*Optimizing model parameters based on training data.*

Artifact: Trained model

### AI Regulation

Tool Confidence Level: 1

Potential tool failure: Lack of AI Transparency

RPN: 105

Potential tool failure: May fail to generate or link detailed training documentation and validation reports, resulting in trained models lacking required transparency and conformity evidence.

#### Possible Process Mitigations:

- Utilize ML frameworks that offer native support for model interpretation techniques like SHAP and LIME.
- Establish structured logging of intermediate decisions and model outputs for traceability.
- Ensure framework compatibility with transparency libraries via plugin architecture or open APIs.

#### Testable Tool Capabilities:

- Integrated explainability methods (e.g., SHAP, LIME) within model training environments
- Comprehensive model decision logging and training metadata tracking
- Easy integration with external explainability libraries and visualization tools

## Possible Process Mitigations:

Here, methods or procedures can be selected that are intended to reduce the risk of a specific error condition.

## Testable Tool Capabilities:

Here, specific features or functions of tools that are to be tested or reviewed in the assessment can be marked.

# AssessML: export and print

The results in the Assessment Summary can be exported as a PDF and printed by pressing the 'Export PDF' button.

The screenshot displays the 'Assessment Summary' page of the AssessML tool. The page is divided into four main sections: 'Failure Mode Selection', 'Task Assessment', 'Tool Assessment', and 'Assessment Summary'. The 'Assessment Summary' section is currently active and shows the following details:

- Phase:** Experimentation
- Task:** Experiment Tracking & Management
- Artifact:** Data storage system
- Functional Safety**
- Tool Confidence Level: 2**
- Potential tool failure: Inconsistent Model Reproducibility**
- RPN: 315**
- Possible Process Mitigations:**
  - Mandate seed fixing across frameworks and environment metadata capture for every run.
  - Store all experiment artifacts, parameters, and outputs with unique, traceable version identifiers.
  - Capture and lock dependency environments during each experiment run.

The 'Export PDF' button is highlighted with a red box, and a red arrow points to it. A print dialog is also open, showing the 'Drucken' (Print) options, including the target printer (prt-6058), page count (3 Blatt Papier), and various print settings like 'Hochformat' (Portrait) and 'An Seitenbreite anpassen' (Fit to page width).

## Additional features of the AssessML

# AssessML: saving and loading

Why should one save and load?

The *Save and Load Settings* option allows you to save your settings and assessment settings and reload them as needed.

The screenshot displays the AssessML interface with a table of failure modes. A settings overlay is positioned over the table, highlighting the 'Model Instability & Poor Generalization' row. The overlay contains a text input field for the settings filename, and 'Save Settings' and 'Load Settings' buttons.

Select	Failure Mode	Description
<input checked="" type="checkbox"/>	<b>AI Regulation</b>	
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input checked="" type="checkbox"/>	Lack of Model Risk Assessment	Failure to categorize AI models under AI Act risk levels and document safety concerns.
<input checked="" type="checkbox"/>	Security Vulnerabilities in AI	Exposure to adversarial attacks compromising AI safety and trustworthiness.
<input checked="" type="checkbox"/>	Non-compliant Human Oversight Mechanisms	Failure to implement necessary human-in-the-loop controls in high-risk AI applications.
<input checked="" type="checkbox"/>	Insufficient AI Impact Assessment	Failure to conduct proper risk-benefit analysis for high-risk AI applications.
<input checked="" type="checkbox"/>	Lack of Ethical Safeguards in AI	Failure to ensure ethical principles (e.g., non-harm, fairness) in AI decision-making.
<input checked="" type="checkbox"/>	Inadequate Monitoring of Deployed AI Systems	Lack of continuous assessment of AI systems in real-world settings.
<input checked="" type="checkbox"/>	Lack of Auditability & Documentation	Inability to provide detailed logs of data usage, model training, validation etc. hindering compliance audits.
<input checked="" type="checkbox"/>	<b>Data Protection</b>	
<input checked="" type="checkbox"/>	Unauthorized Data Access	Lack of strict access control leading to potential data breaches and privacy violations.
<input checked="" type="checkbox"/>	Failure to Handle Right to Erasure	Non-compliance with GDPR's Right to be Forgotten, leading to legal risks.
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.
<input checked="" type="checkbox"/>	Lack of	
<input checked="" type="checkbox"/>	Failure	
<input checked="" type="checkbox"/>	Non-Tr	
<input checked="" type="checkbox"/>	Cross-b	
<input checked="" type="checkbox"/>	<b>Function</b>	
<input checked="" type="checkbox"/>	Data Int	
<input checked="" type="checkbox"/>	Lack of	
<input checked="" type="checkbox"/>	Failure Mode Undetection	Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.
<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility	Non-deterministic model outputs lead to unpredictable AI system behavior.
<input checked="" type="checkbox"/>	Insufficient Edge Case Handling	AI models fail to correctly handle rare but critical safety scenarios, leading to unreliable operation.
<input checked="" type="checkbox"/>	Undetected Model Performance Drift	Model performance degrades over time without triggering corrective actions, leading to unsafe decisions.
<input checked="" type="checkbox"/>	Insufficient Redundancy & Failover Mechanisms	Lack of backup AI models or failover strategies in critical systems.
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinders failure analysis and root cause identification.
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

**Model Instability & Poor Generalization** **Model**

Settings filename (e.g. setti | Save Settings | Load Settings

Settings filename (e.g. setti | Save Settings | Load Settings

# AssessML: saving and loading

## Explanation of the storage process

Choose Phase(s):  Choose Task(s):  Choose Tool(s):

Failure Mode Selection	Task Assessment	Tool Assessment	Assessment Summary
------------------------	-----------------	-----------------	--------------------

### Failure Mode and Compliance Area Selection

Select	Failure Mode	Description
<input checked="" type="checkbox"/> <b>AI Regulation</b>		
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input type="checkbox"/> <b>Data Protection</b>		
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.
<input type="checkbox"/> <b>Functional Safety</b>		
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

Model instability & Poor Generalization model may pe

Search and assessment settings that have already been made can be saved by clicking 'Save Settings'. To do this:

1. Select the check box



2. Enter a name



3. Save by clicking 'Save Settings'

# AssessML: saving and loading

## Explanation of the loading process

Saved settings are listed below. To retrieve saved settings:

Choose Phase(s):

Choose Task(s):

Failure Mode Selection	Task Assessment	
<b>Failure Mode and Compliance Area Selection</b>		
Select	Failure Mode	Description
<input checked="" type="checkbox"/>	<b>AI Regulation</b>	
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.
<input checked="" type="checkbox"/>	Lack of Model Risk Assessment	Failure to categorize AI models under AI Act risk levels and document safety concerns.
<input checked="" type="checkbox"/>	Security Vulnerabilities in AI	Exposure to adversarial attacks compromising AI safety and trustworthiness.
<input checked="" type="checkbox"/>	Non-compliant Human Oversight Mechanisms	Failure to implement necessary human-in-the-loop controls in high-risk AI applications.
<input checked="" type="checkbox"/>	Insufficient AI Impact Assessment	Failure to conduct proper risk-benefit analysis for high-risk AI applications.
<input checked="" type="checkbox"/>	Lack of Ethical Safeguards in AI	Failure to ensure ethical principles (e.g., non-harm, fairness) in AI decision-making.
<input checked="" type="checkbox"/>	Inadequate Monitoring of Deployed AI Systems	Lack of continuous assessment of AI systems in real-world settings.
<input checked="" type="checkbox"/>	Lack of Auditability & Documentation	Inability to provide detailed logs of data usage, model training, validation etc. hindering compliance audits.
<input checked="" type="checkbox"/>	<b>Data Protection</b>	
<input checked="" type="checkbox"/>	Unauthorized Data Access	Lack of strict access control leading to potential data breaches and privacy violations.
<input checked="" type="checkbox"/>	Failure to Handle Right to Erasure	Non-compliance with GDPR's Right to be Forgotten, leading to legal risks.
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.
<input checked="" type="checkbox"/>	Lack of Data Minimization	Collecting and storing excessive personal data, violating GDPR principles.
<input checked="" type="checkbox"/>	Failure in Data Anonymization & Pseudonymization	Inadequate de-identification techniques lead to re-identifiable personal data.
<input checked="" type="checkbox"/>	Non-Transparent User Consent Management	Lack of clear user consent tracking for data collection and processing.
<input checked="" type="checkbox"/>	Cross-border Data Transfer Violations	Failure to comply with international data transfer rules under GDPR/Schrems II.
<input checked="" type="checkbox"/>	<b>Functional Safety</b>	
<input checked="" type="checkbox"/>	Data Integrity Failure	Corrupt or inconsistent data leading to incorrect model predictions, endangering safety-critical systems.
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.
<input checked="" type="checkbox"/>	Failure Mode Undetection	Inability to detect and mitigate failure conditions in deployed AI models, leading to high-risk scenarios.
<input checked="" type="checkbox"/>	Inconsistent Model Reproducibility	Non-deterministic model outputs lead to unpredictable AI system behavior.
<input checked="" type="checkbox"/>	Insufficient Edge Case Handling	AI models fail to correctly handle rare but critical safety scenarios, leading to unreliable operation.
<input checked="" type="checkbox"/>	Undetected Model Performance Drift	Model performance degrades over time without triggering corrective actions, leading to unsafe decisions.
<input checked="" type="checkbox"/>	Insufficient Redundancy & Failover Mechanisms	Lack of backup AI models or failover strategies in critical systems.
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.

Setting\_1

Settings saved to Setting\_1

Setting\_1

Settings saved to Setting\_1

1. Enter setting names
2. Click Load Settings
3. Dashboard loads configuration

Choose Phase(s):

Choose Task(s):

Choose Tool(s):

Failure Mode Selection	Task Assessment	Tool Assessment	Assessment Summary
<b>Failure Mode and Compliance Area Selection</b>			
Select	Failure Mode	Description	
<input checked="" type="checkbox"/>	<b>AI Regulation</b>		
<input checked="" type="checkbox"/>	Lack of AI Transparency	Failure to provide explanations for AI decisions, violating AI Act transparency requirements.	
<input checked="" type="checkbox"/>	Bias & Fairness Non-Compliance	High-risk AI models reinforcing biases, leading to discrimination and non-compliance with AI Act fairness guidelines.	
<input type="checkbox"/>	<b>Data Protection</b>		
<input checked="" type="checkbox"/>	Personal Data Leakage	ML models unintentionally exposing personal data during training or inference.	
<input type="checkbox"/>	<b>Functional Safety</b>		
<input checked="" type="checkbox"/>	Lack of Real-Time Processing	Failure to meet real-time inference requirements in autonomous or safety-critical applications.	
<input checked="" type="checkbox"/>	Inadequate Error Handling & Logging	Poor logging and monitoring hinder failure analysis and root cause identification.	
<input checked="" type="checkbox"/>	Model Instability & Poor Generalization	Model may perform unpredictably or fail to generalize to unseen conditions, undermining functional safety.	

Setting\_1

