

DNN SAFETY CONCERNS AND THEIR MITIGATIONS

Dr. Shervin Raafatnia
Robert Bosch GmbH





SoTA Standards and Their Applicability to DNNs

- **ISO 26262:** focus on E/E malfunctions and HW failures
→ disregards many aspects relevant to assuring safety
DNN applicability: HW issues, tooling verification, ...

- **ISO/PAS 21448 (aka. SOTIF):** addresses functional insufficiencies leading to hazardous behavior when triggered
 - “A proper understanding of the function by the user, its behavior and its limitations (including the human/machine interface) is the key to ensuring safety”*DNN applicability:* appropriate logic

DNNs and Safety



Functional
Insufficiencies

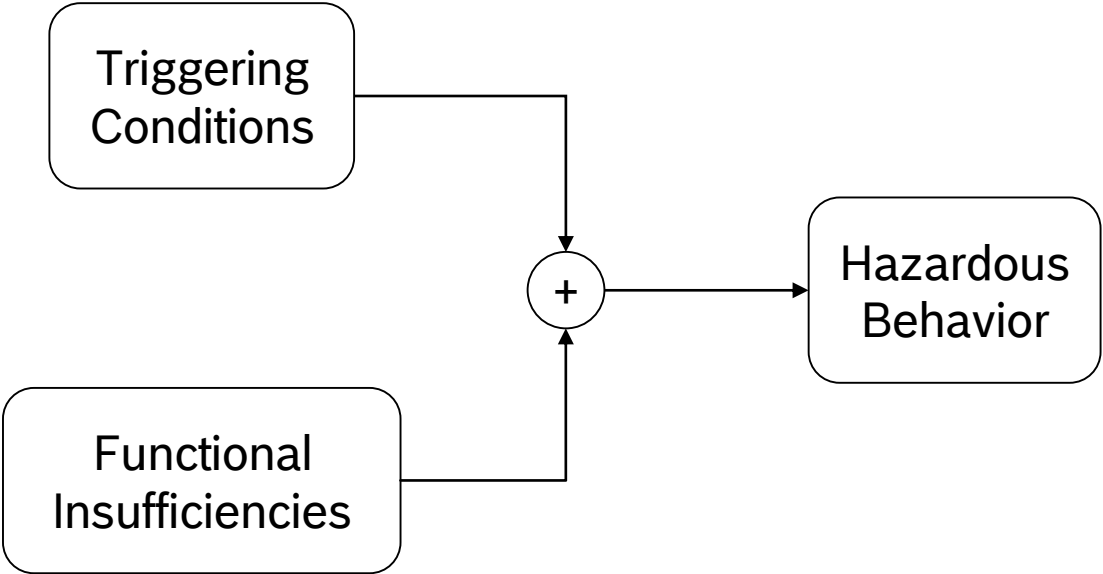
DNNs and Safety



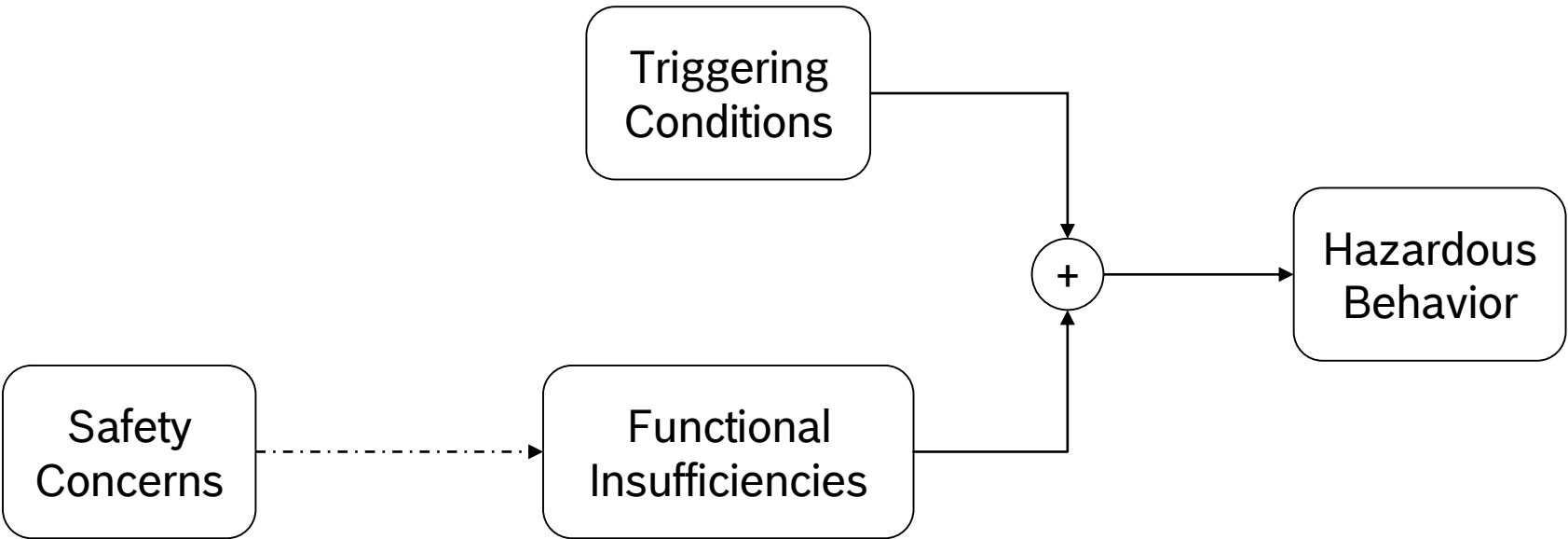
Triggering
Conditions

Functional
Insufficiencies

DNNs and Safety

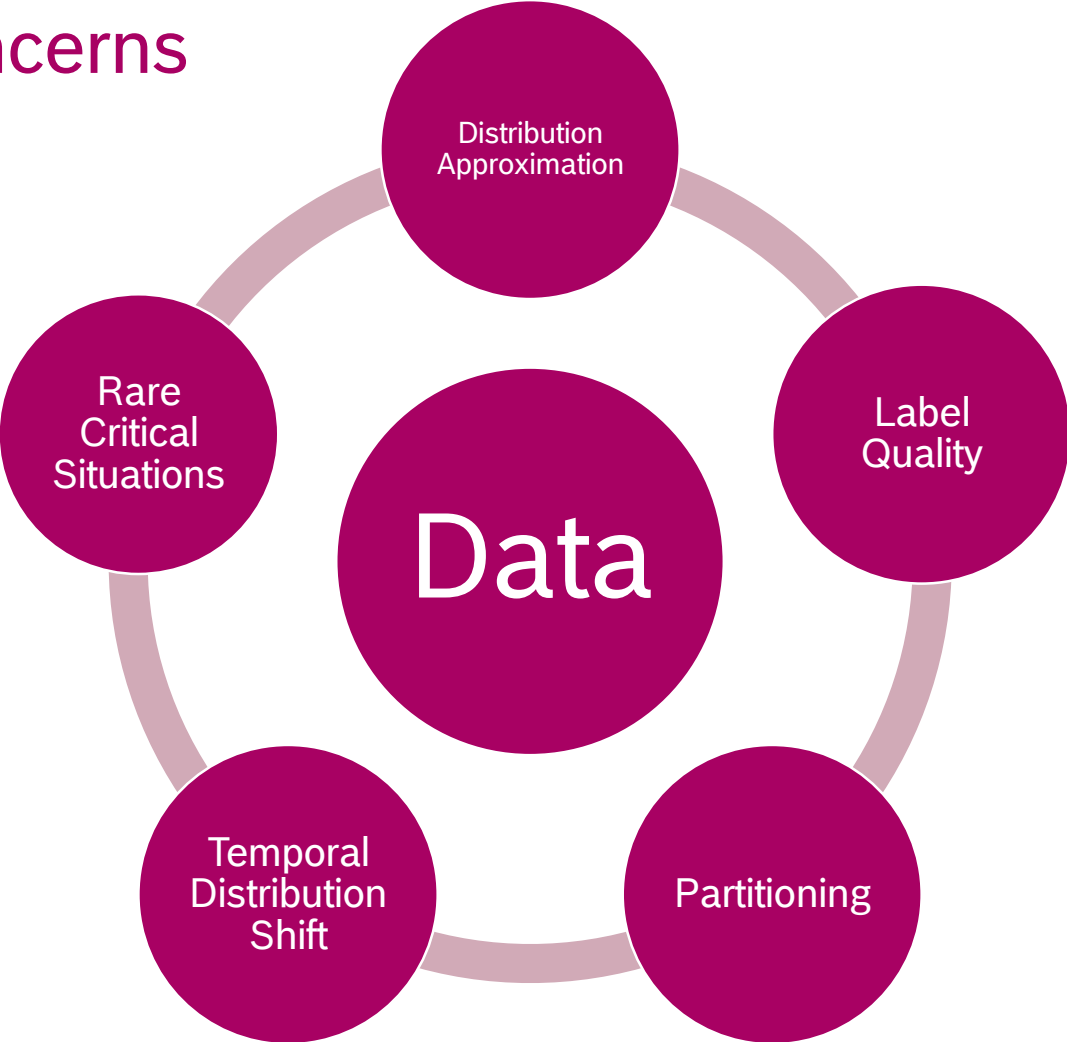


DNNs and Safety

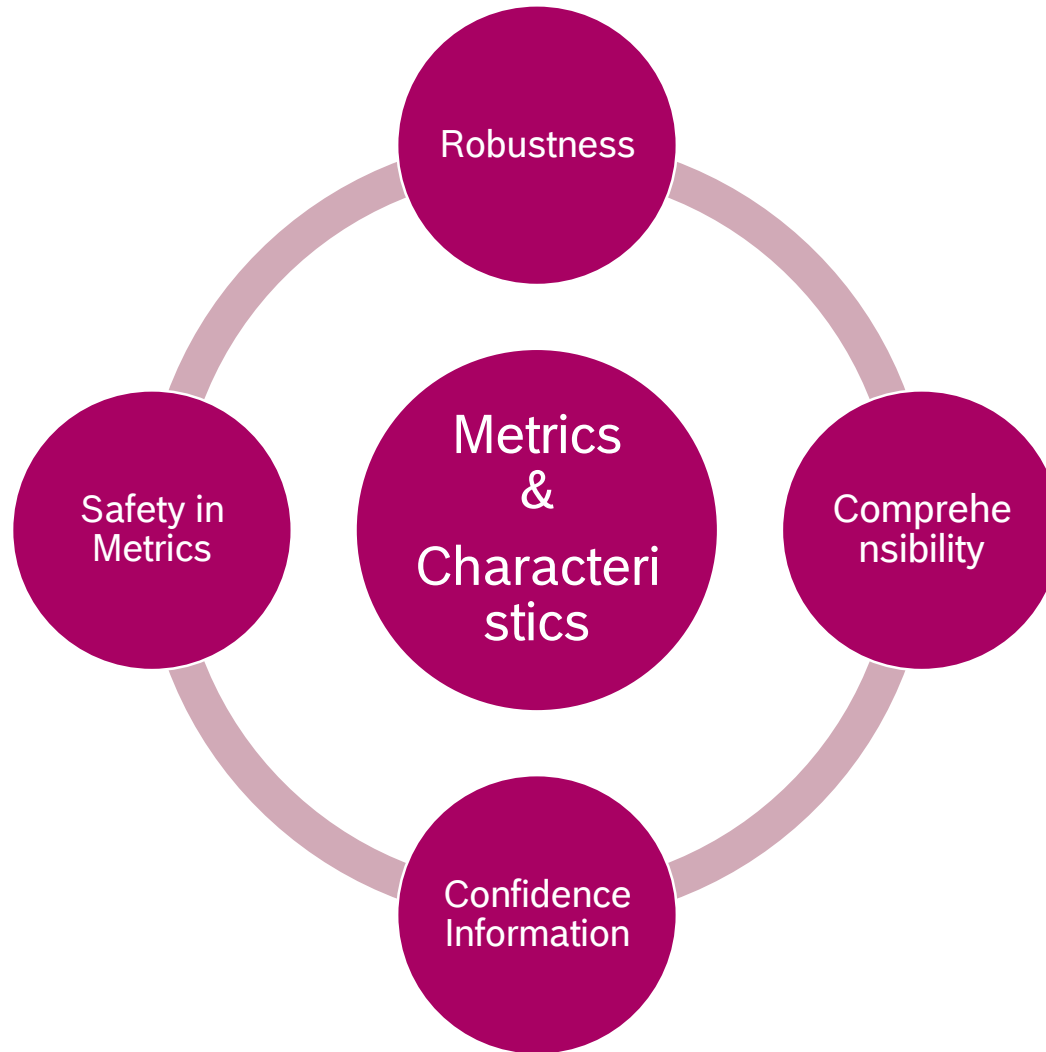




Data-related Concerns












Other Concerns





Potential Mitigation Approaches

 <p>Real-World Distribution Approximation</p>	<ul style="list-style-type: none">• <i>Well-justified data acquisition strategy</i>• <i>Data analysis</i>• <i>Reliable confidence information</i>• <i>Testing and detailed results analysis</i>	 <p>Distributional shift over time</p>	<ul style="list-style-type: none">• <i>Reliable confidence information</i>• <i>Continuous learning and updating</i>	 <p>Unreliable confidence information</p>	<ul style="list-style-type: none">• <i>Reliable confidence information</i>• <i>Using Gray-box methods</i>• <i>Testing and detailed results analysis</i>
 <p>Data Partitioning</p>	<ul style="list-style-type: none">• <i>Data partitioning guidelines</i>• <i>Assessing the adherence to partitioning guidelines</i>	 <p>Unknown behavior in rare critical situations</p>	<ul style="list-style-type: none">• <i>Well-justified data acquisition strategy</i>• <i>Reliable confidence information</i>• <i>Random and dedicated testing</i>• <i>Continuous learning and updating</i>	 <p>Incomprehensible Behavior</p>	<ul style="list-style-type: none">• <i>Testing and detailed results analysis</i>• <i>Using gray-box methods</i>
 <p>Dependence on labeling quality</p>	<ul style="list-style-type: none">• <i>Labeling guidelines</i>• <i>Continuous quality checks</i>	 <p>Brittleness of DNNs</p>	<ul style="list-style-type: none">• <i>Reliable confidence information</i>• <i>Adversarial attacks defenses</i>• <i>Targeted testing and detailed analysis development process</i>• <i>Continuous learning and updating</i>	 <p>Insufficient consideration of safety in metrics</p>	<ul style="list-style-type: none">• <i>Defining heuristics to assess criticality/relevance of different objects</i>• <i>Performance evaluation with respect to safety</i>

**THANKS FOR YOUR
ATTENTION**