

The logo consists of the letters 'DIN' in a bold, sans-serif font, enclosed within a white square. This square is positioned on a dark blue background that is part of a larger graphic element on the left side of the slide, which includes several overlapping squares in various shades of blue.

Zur Normungsroadmap KI Ausgabe 2 – 2023-07-07, Start um 13:30 Uhr

KI – bestens geschützt mit Normen und Standards

Workshop zu Normungs- und Standardisierungsbedarfen aus der Roadmap

Mit Annegrit Seyerlein-Klug und Dr. rer. nat. Henrich C. Pöhls
Moderiert von Jan Rösler

Agenda

1. Begrüßung

- Warmup
- Ziele dieses Workshops

2. Thematische Einführung

- Überblick zur Normungsroadmap KI
- Vorstellung der Speaker
- Impulsvorträge
- Vorstellung der Normungs- und Standardisierungsbedarfe
- Beantwortung von Fragen zu den Bedarfen

3. Interaktive Priorisierung der Bedarfe

- Interessierte für Bedarfsumsetzung
- Priorisierung der Bedarfe für den Workshop

KURZE PAUSE

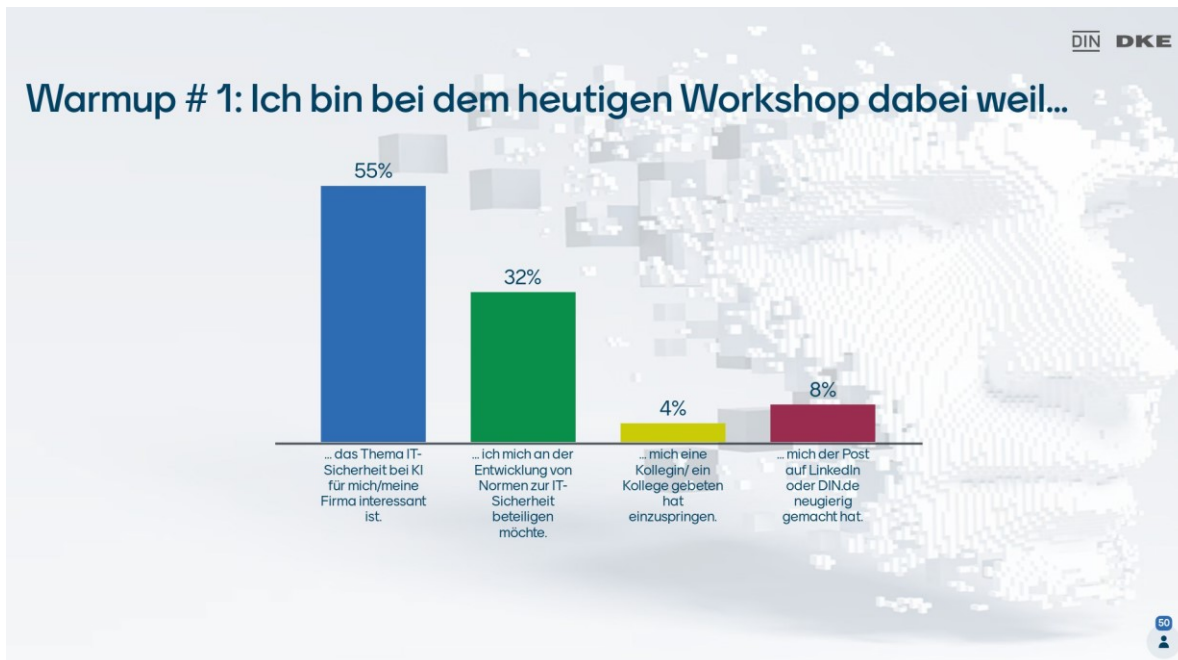
4. Diskussion der Bedarfe anhand der Priorisierung

- Konkretisierung der Bedarfsschwerpunkte
- Bedarfsabgrenzung durch Anwendungsbereiche (Scopes)
- Beziehung zwischen Bedarfen und laufenden Projekten oder veröffentlichten Normen und Standards

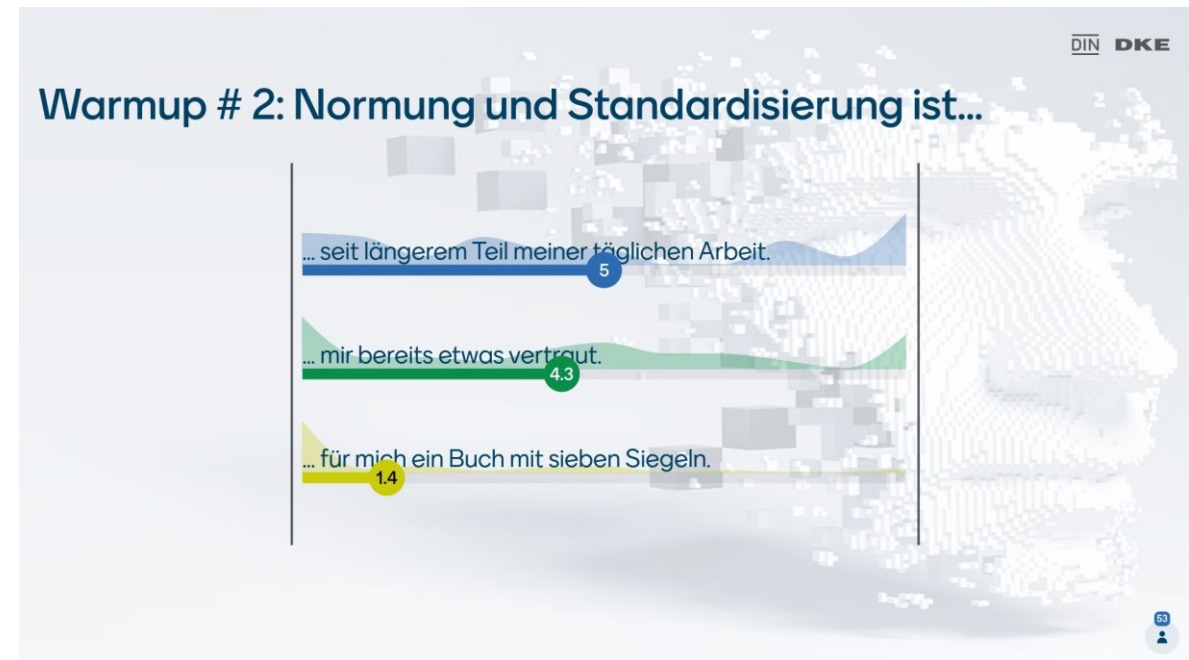
5. Nächste Schritte

1. Begrüßung

Warmup per Mentimeter-Umfrage – Ergebnisse



Mehrfachauswahl war möglich



Statements wurden von 1 bis 10 gewichtet

Ziele dieses Workshops

Recap

- Worum geht's überhaupt?
- Ergebnisse der Normungsroadmap KI
- Umsetzung der Handlungsbedarfe

Interaktion

- Vorstellung der Standardisierungsbedarfe
- Gemeinsames Verständnis schaffen
- Diskussion der Bedarfe

Beginn...

- der Bedarfsumsetzung und
- Vorbereitung für die Normung

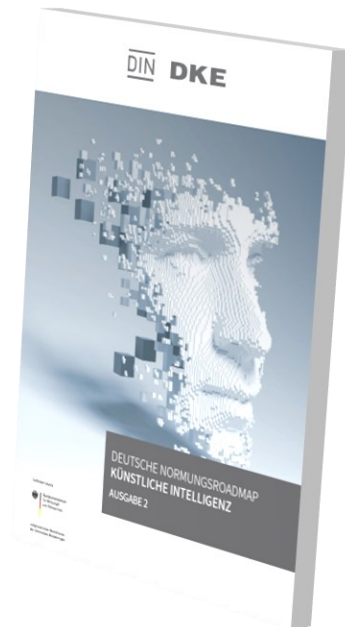
2. Thematische Einführung

Veröffentlichung der Normungsroadmap KI Ausgabe 2

- Maßnahme der KI Strategie der Bundesregierung
- Veröffentlicht am 9. Dezember 2022
- Fortschreibung der ersten Ausgabe (2020)
- Fokus: AI Act
- Kostenfreier Download:



www.din.de/go/normungsroadmapki



Filiz Elmas, Leiterin Strategische Entwicklung Künstliche Intelligenz bei DIN, Christoph Winterhalter, Vorsitzender des Vorstandes von DIN, Robert Habeck, Vizkanzler und Bundesminister für Wirtschaft und Klimaschutz, Prof. Dr. Wolfgang Wahlster, CEA des DFKI und Michael Teigeler, Geschäftsführer DKE (v.l.n.r.) © Stefan Zeitz

2. Thematische Einführung

Überblick zur Normungsroadmap KI Ausgabe 2



Normen und Standards ermöglichen eine **zuverlässige und sichere Anwendung** von KI-Technologien und tragen zur **Erklärbarkeit und Nachvollziehbarkeit** bei.

Sie sind dadurch Schlüsselfaktoren für die **Akzeptanz von KI-Anwendungen** und schaffen **Vertrauen** am Markt und bei Verbraucher*innen.



Die Normungsroadmap KI beschreibt das **Umfeld der KI-Standardisierung** und gibt einen Überblick über bestehende Normen und Standards.

Der Fokus ist das Aufzeigen von **Normungs- und Standardisierungsbedarfen**, sowie konkreter **Handlungsempfehlung**.

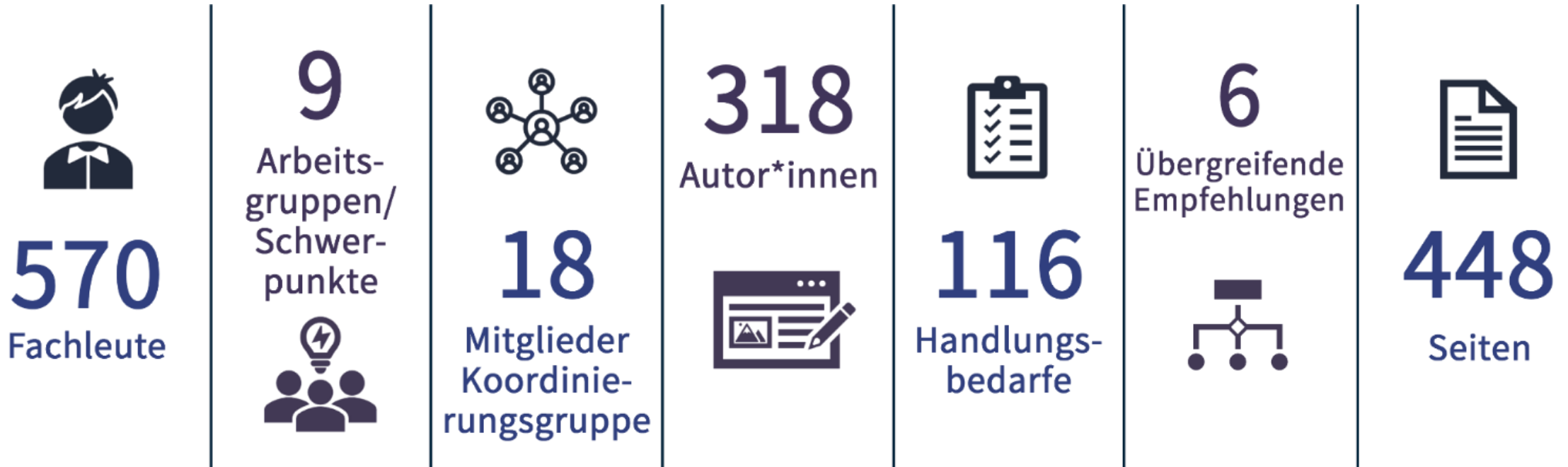


Die Erarbeitung erfolgte unter Mitwirkung von Expert*innen aller relevanten Kreise.

Dafür wurden Arbeitsgruppen zu verschiedenen Schwerpunktthemen gegründet.

Bearbeitungszeitraum war von Januar 2022 bis September 2022.

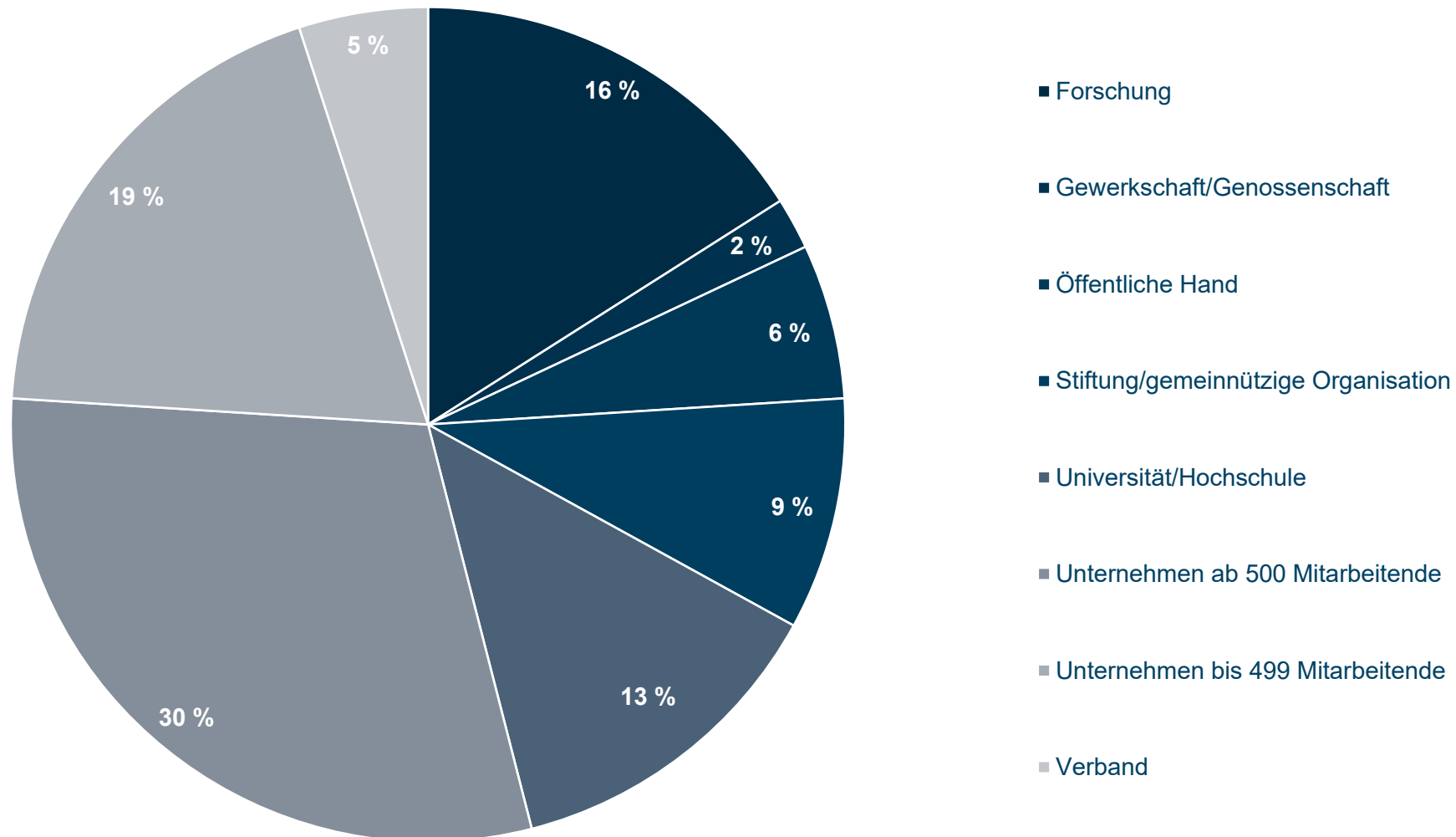
Überblick zur Normungsroadmap KI Ausgabe 2



2. Thematische Einführung

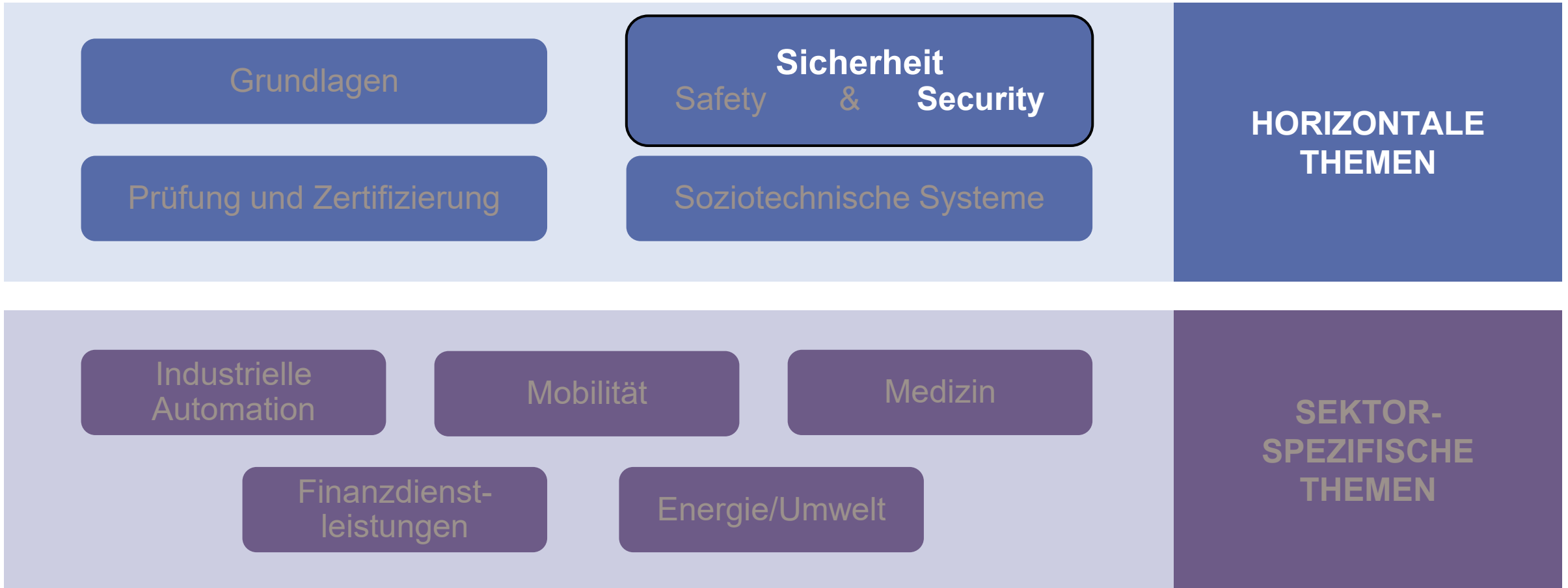
Überblick zur Normungsroadmap KI Ausgabe 2

Anzahl: 570 Fachleute



2. Thematische Einführung


Überblick zur Normungsroadmap KI Ausgabe 2




2. Thematische Einführung

Überblick zur Normungsroadmap KI Ausgabe 2


Insgesamt 116 Handlungsbedarfe identifiziert, unterteilt in 3 Kategorien:



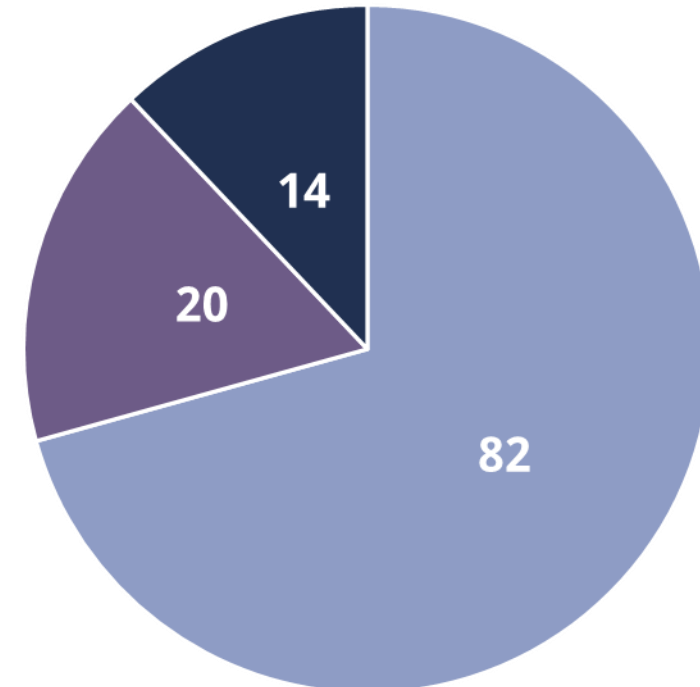
KATEGORIE 1:
Bedarf adressiert Normung und Standardisierung



KATEGORIE 2:
Bedarf adressiert Forschung



KATEGORIE 3:
Bedarf adressiert Politik/Gesetzgeber



■ Normung und Standardisierung ■ Forschung ■ Gesetzgeber

Stand: Mai 2023

2. Thematische Einführung

Conceptboard

Während des Workshops wurde Conceptboard verwendet. Das Board ist nun geschlossen. Sie können sich aber noch per E-Mail an Jan.Roesler@din.de für die Umsetzung der Bedarfe melden – bitte Bedarfs-Code mit angeben.

2. Thematische Einführung

Vorstellung der Speaker



Annegrit Seyerlein-Klug

Technische Hochschule Brandenburg
neurocat GmbH



Dr. rer. nat. Henrich C. Pöhls

Universität Passau

2. Thematische Einführung Themenschwerpunkt: KI-Sicherheit bei IT-Systemen

Vertrauen in KI durch Prüfung/Zertifizierung

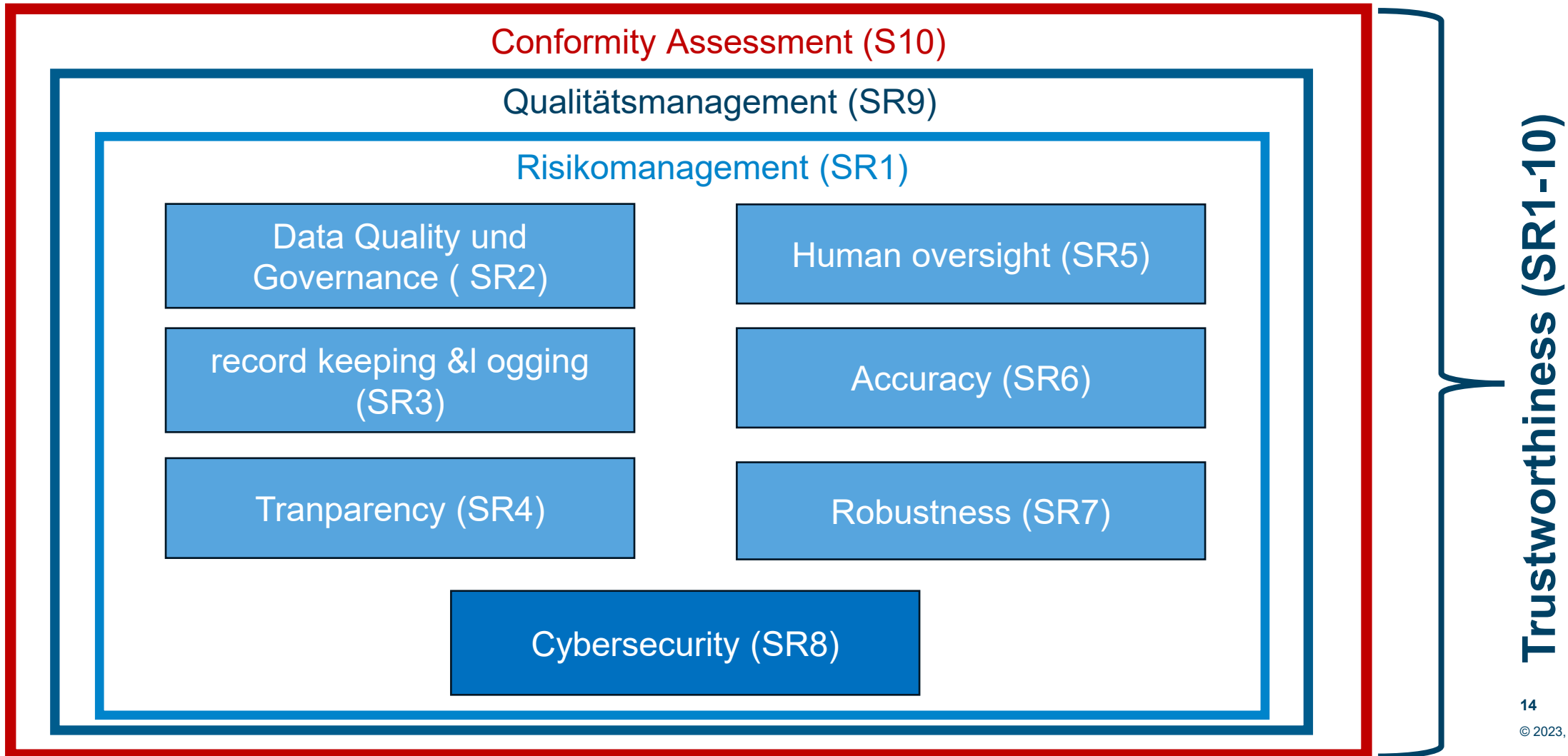
Informations-Sicherheit bei KI und zugehörige Schutzziele, 3 Perspektiven:

- Security - Vertraulichkeit, Integrität, Verfügbarkeit, Authentizität
- Safety - Leib und Leben
- Privacy - personenbezogener Daten

Konformität mit dem geplanten EU AI Act (Trilog) bei Hochrisikosysteme

2. Thematische Einführung Themenschwerpunkt: KI-Sicherheit bei IT-Systemen

Standardisierungsanforderung des EU AI Act



2. Thematische Einführung Themenschwerpunkt: KI-Sicherheit bei IT-Systemen

Cybersecurity specifications AI Systems (Software)

„This European standard or European standardisation deliverable **shall provide suitable organisational and technical solutions**, to ensure that AI systems are resilient against attempts to alter their use, behaviour, or performance or to compromise their security properties by malicious third parties exploiting the AI systems’ vulnerabilities.

Organisational and technical solutions shall therefore include, where appropriate, **measures to prevent and control cyberattacks trying to manipulate AI specific assets**, such as

- **training data sets** (e.g. data poisoning) or
- **trained models** (e.g. adversarial examples), or
- **trying to exploit vulnerabilities** in an AI system’s digital assets or the
- **underlying ICT infrastructure**.

These technical solutions shall be appropriate to the relevant circumstances and risks.

This European standard or European standardisation deliverable shall take due account of the **essential requirements for products with digital elements**

as listed in Sections 1 and 2 of Annex I to the proposal for a Regulation of the European Parliament and the Council on **horizontal cybersecurity requirements** for products with digital elements.“ – Standardization Request zum AI-Act

2. Thematische Einführung Themenschwerpunkt: KI-Sicherheit bei IT-Systemen

AI Act – Conformity Assessment for AI systems

*„This European standard or European standardisation deliverable shall provide **procedures and processes** for conformity assessment activities related to AI systems and **quality management systems** of AI providers.*

*This European standard or European standardisation deliverable shall also provide **criteria for assessing the competence of persons** tasked with the conformity assessment activities.“*

*„This European standard or European standardisation deliverable shall consider both the scenarios whereby the conformity assessment is carried out by the **provider itself** or with the **involvement of a professional external third-party organisation**.“*

2. Thematische Einführung

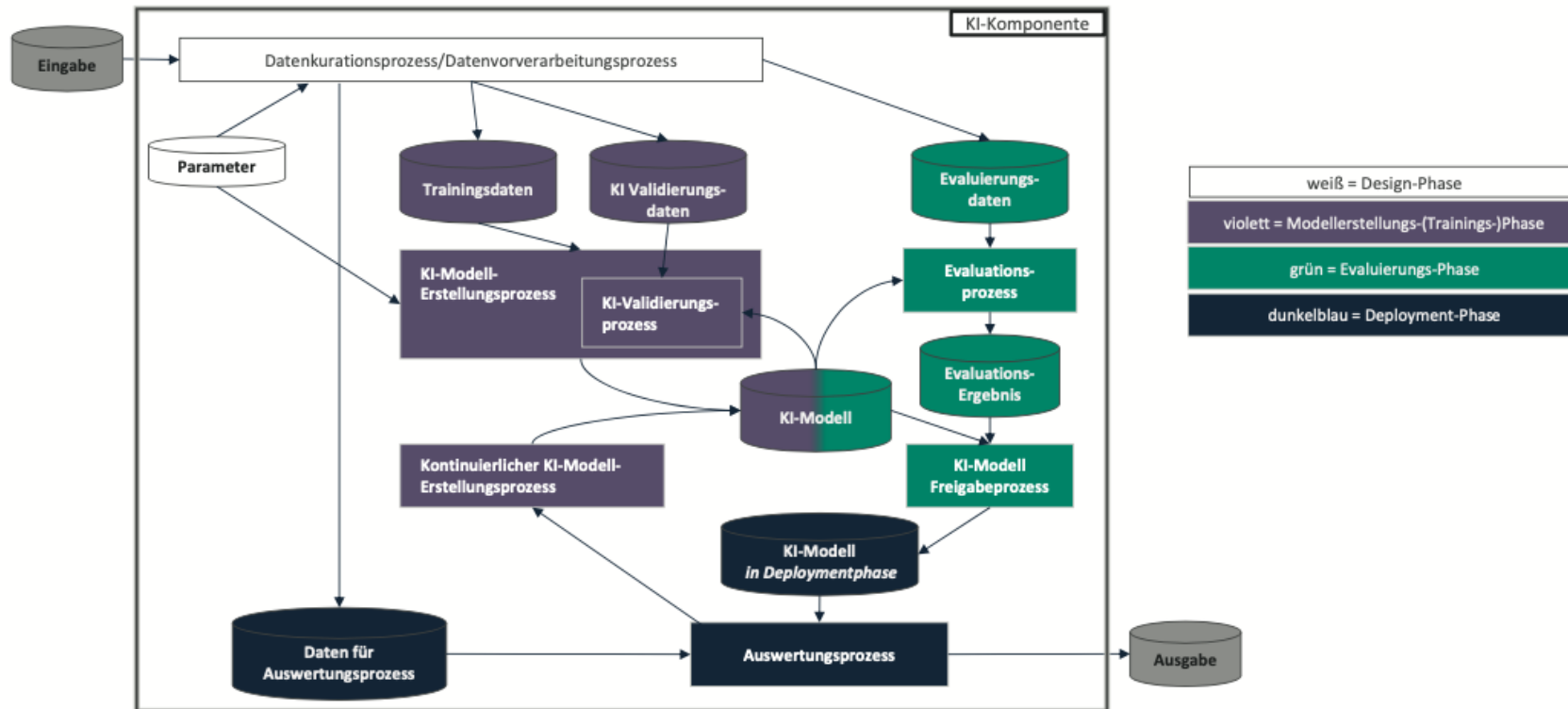
Security für KI – neue Herausforderungen

- Grundsätzlich können bei KI-Systemen, wie bei anderen IT-Systemen auch, alle üblichen Sicherheitsschutzziele wie Vertraulichkeit, Verfügbarkeit, Integrität oder auch Belastbarkeit durch Angriffe gefährdet sein.
- Die ‚Unzulänglichkeiten‘, welche spätere Angriffe ermöglichen, sind vielfältig
 - z.B. Fehler in den Daten durch menschliches Fehlverhalten
 - z.B. gezielte Manipulation von Daten (z.B. Data Poisoning von Trainingsdaten)
 - z.B. unzureichender Trainingsprozess beim Maschinellen Lernen
- Grundsätzlich drei Fragestellungen:
 1. Welche **bereits existierenden Normen** kann man zur Erhöhung (und zum Nachweis dieser Erhöhung) der IT-Sicherheit von KI-Systemen sinnvoll **direkt nutzen**?
 2. Welche **bereits existierenden Normen** muss man um KI-spezifische IT-Sicherheitsprobleme **erweitern**?
 3. Welche **neuen Normen** muss man für KI-spezifische Sicherheitsprobleme **schaffen**?

2. Thematische Einführung

Security für KI – neue Herausforderungen

1. Herausforderung: Definition von Schutzzielen auf der Ebene von Prozessen und Daten innerhalb der KI-Komponente



NRM KI v2 – Seite 122 – Abbildung 25

2. Thematische Einführung – Security bei KI – neue Herausforderungen

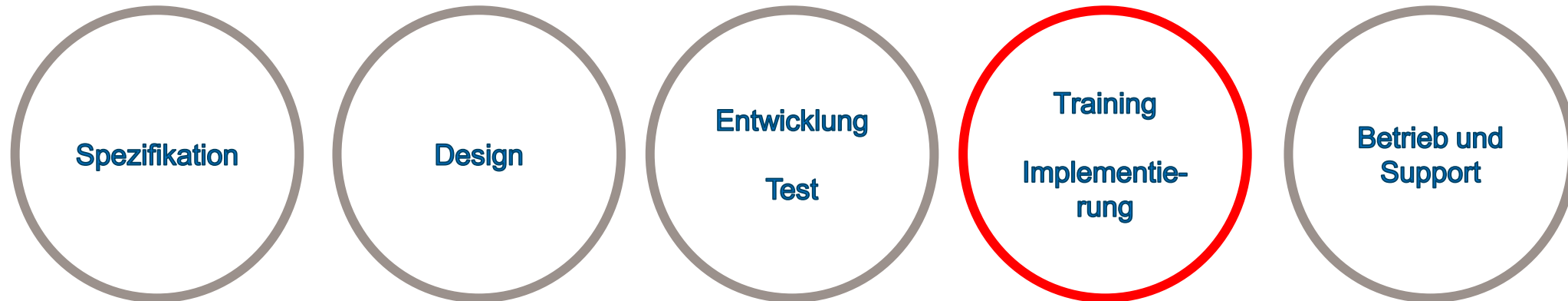
Ein Beispiel: Sicherung der Trainingsdaten

Angriff:

Poisoning (=Manipulation) der Datensätze für das Training eines Machine Learning Verfahrens



Bild:
<https://owasp.org/www-project-ai-security-and-privacy-guide/>



2. Thematische Einführung – Security bei KI – neue Herausforderungen

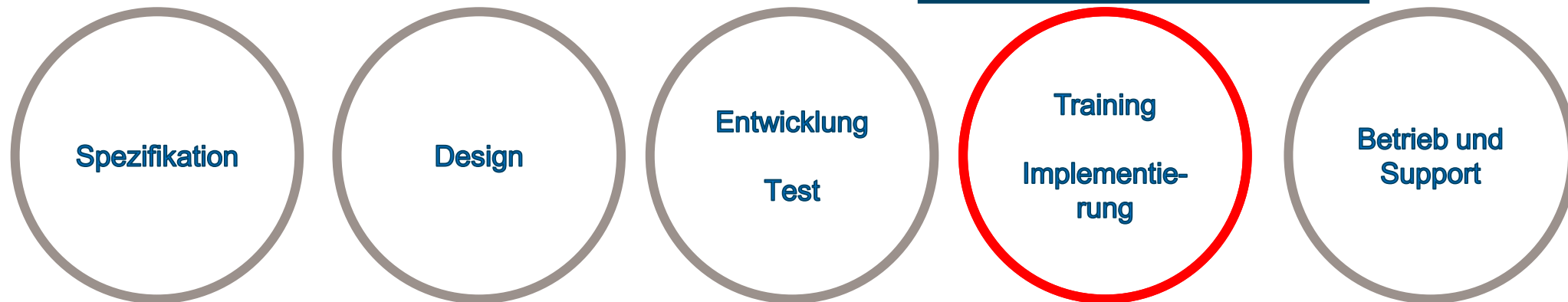
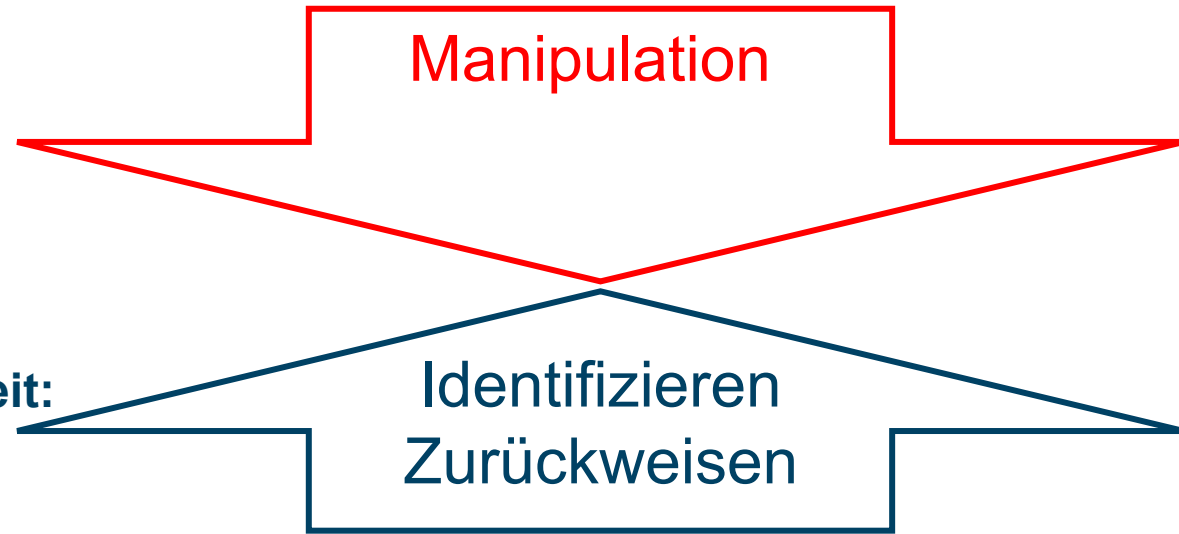
Ein Beispiel: Sicherung der Trainingsdaten

Angriff:

Poisoning (=Manipulation) der Datensätze für das Training eines Machine Learning Verfahrens

Verteidigung

mittels „klassischer“ IT-Sicherheit:
Datenintegrität, Datenauthenzizität



2. Thematische Einführung – Security bei KI – neue Herausforderungen

Ein Beispiel: Sicherung der Trainingsdaten

Angriff:

Poisoning (=Manipulation) der Datensätze für das Training eines Machine Learning Verfahrens

Verteidigung

mittels “sicherer“ Algorithmen:

Manipulation

Identifizieren
Zurückweisen

Design

Entwicklung
Test

Training
Implementierung

Betrieb und
Support

Manipulating Machine Learning: Poisoning Attacks and Countermeasures for Regression Learning

Matthew Jagielski*, Alina Oprea[†], Battista Biggio[‡], Chang Liu[§], Cristina Nita-Rotaru^{*}, and Bo Li[§]

*Northeastern University, Boston, MA [†]University of Cagliari, Italy [‡]Pluribus One, Italy [§]UC Berkeley, Berkeley, CA

Abstract—As machine learning becomes widely used for automated decisions, attackers have strong incentives to manipulate the results and models generated by machine learning algorithms. In this paper, we perform the first systematic study of poisoning attacks and their countermeasures for linear regression models. In poisoning attacks, attackers deliberately influence the training data to manipulate the results of a predictive model. We propose a theoretically-grounded optimization framework specifically designed for linear regression and demonstrate its effectiveness on a range of datasets and models. We also introduce a fast statistical attack that requires limited knowledge of the training process. Finally, we design a new principled defense method that is highly resilient against all poisoning attacks. We provide formal guarantees about its convergence and an upper bound on the effect of poisoning attacks when the defense is deployed. We evaluate extensively our attacks and defenses on three realistic datasets from health care, loan assessment, and real estate domains.¹

training process. Such poisoning attacks have been practically demonstrated in worm signature generation [47], [48], spam filters [49], DNS attack detection [49], PDF malware classification [50], handwritten digit recognition [5], and sentiment analysis [41]. We argue that these attacks become easier to mount today as many machine learning models need to be updated regularly to account for continuously-generated data. Such scenarios require *online training*, in which machine learning models are updated based on new incoming training data. For instance, in cyber-security analytics, new Indicators of Compromise (IoC) rise due to the natural evolution of malicious threats, resulting in updates to machine learning models for threat detection [2]. These IoCs are collected from online platforms like VirusTotal, in which attackers can also submit IoCs of their choice. In personalized medicine, it is envisioned that patient treatment is adjusted in real-time by analyzing information crowdsourced from multiple participants [14]. By controlling a few devices, attackers can submit fake information (e.g., sensor measurements), which is then used for training models applied to a large set of patients. Defending against such poisoning attacks is challenging with current techniques. Methods from robust statistics (e.g. [18], [27]) are resilient against noise but perform poorly on adversarially-poisoned data, and methods for sanitization of training data operate under restrictive adversarial models [13].

I. INTRODUCTION
As more applications with large societal impact rely on machine learning for automated decisions, several concerns have emerged about potential vulnerabilities introduced by machine learning algorithms. Sophisticated attackers have strong incentives to manipulate the results and models generated by machine learning algorithms to achieve their objectives. For instance, attackers can deliberately influence the training dataset to manipulate the results of a predictive model (in poisoning attacks [5], [40]–[42], [45], [48], [56]), cause misclassification of new data in the testing phase (in evasive attacks [3], [9], [21], [43], [44], [51], [52]) or infer private information on training data (in privacy attacks [19], [20], [60]). Several surveys from academics and industry highlighted

One fundamental class of supervised learning is linear regression. Regression is widely used for prediction in many settings (e.g., insurance or loan risk estimation, personalized medicine, market analysis). In a regression task a numerical response variable is predicted using a number of predictor

Jagielski et al. *Manipulating Machine Learning: Poisoning attacks and countermeasures for regression learning* - <https://arxiv.org/pdf/1804.00308.pdf> - 2021

2. Thematische Einführung – Security bei KI – neue Herausforderungen

Ein Beispiel: Sichere Algorithmen

Angriff:

Membership inference attack (durch gezielte Anfragen an das trainierte ML Modell Rückschlüsse auf die Trainingsdaten ziehen)

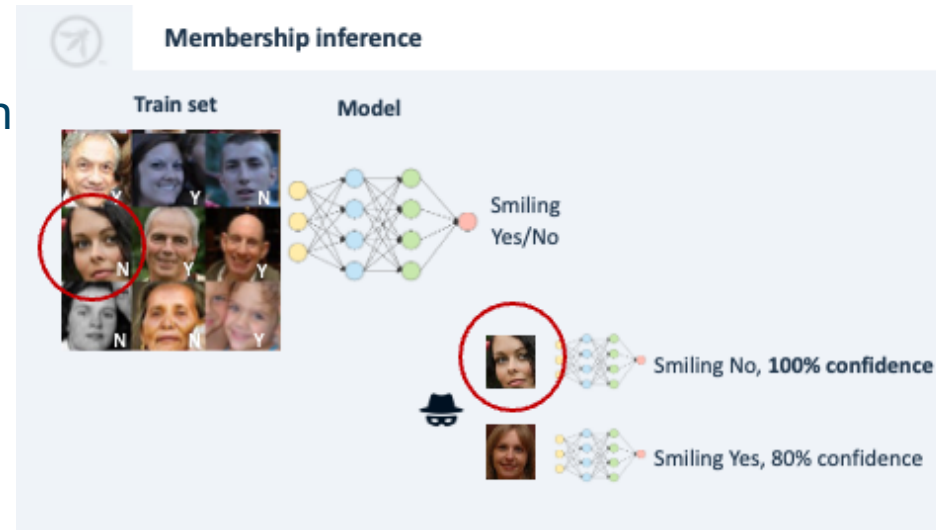
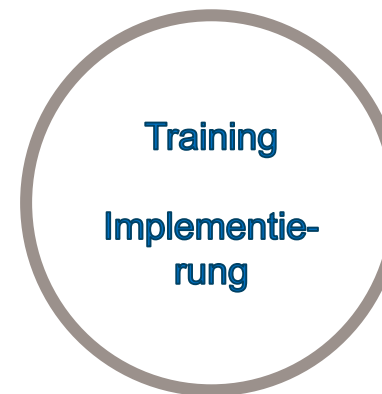
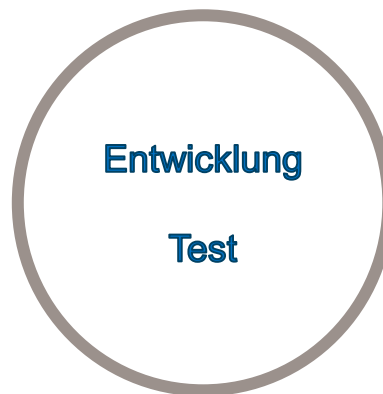
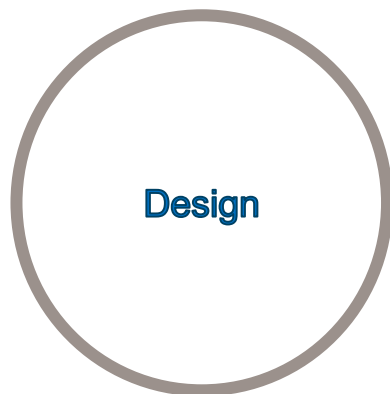


Bild:
<https://owasp.org/www-project-ai-security-and-privacy-guide/>



2. Thematische Einführung – Security bei KI – neue Herausforderungen

Ein Beispiel: Sichere Algorithmen

Angriff:

Membership inference attack (durch gezielte Anfragen an das trainierte ML Modell Rückschlüsse auf die Trainingsdaten ziehen)

Angriff auf die Vertraulichkeit



2. Vorstellung der Normungs- und Standardisierungsbedarfe

Ein Beispiel: Sichere Algorithmen

Angriff:

Membership inference attack (durch gezielte Anfragen an das trainierte ML Modell Rückschlüsse auf die Trainingsdaten ziehen)

Verteidigung:

z.B. Vermeidung der Auswahl ungeeigneter ML Verfahren

Angriff auf die Vertraulichkeit

z.B. Verringern ‚by-design‘



2. Vorstellung der Normungs- und Standardisierungsbedarfe

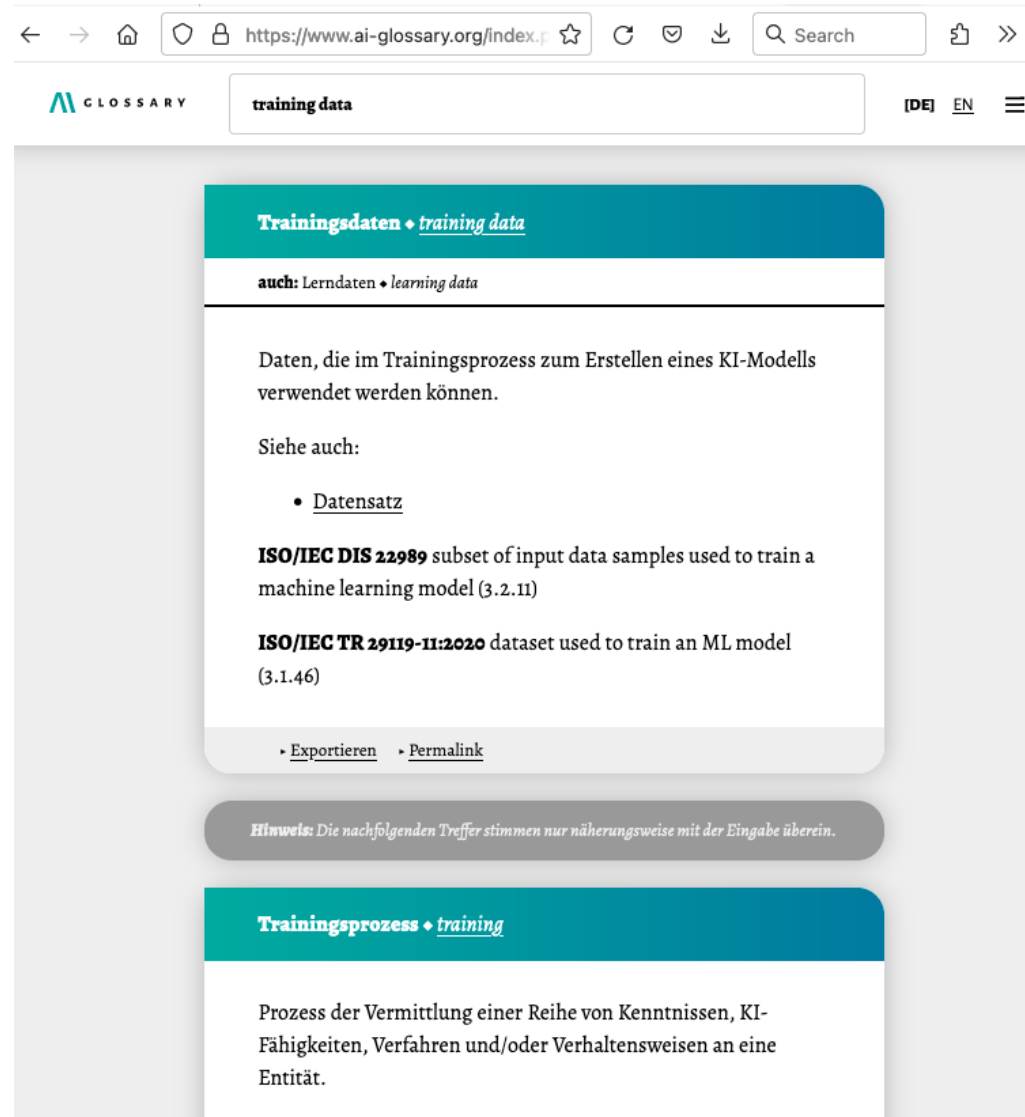
Begriffe für interdisziplinäre Kommunikation

www.ai-glossary.org

Entstanden aus der Glossararbeitsgruppe während der Entwicklung der NRM KI v2

Wird weiter gepflegt und regelmäßig um Begriffe erweitert werden

Versionierte Begriffsdefinitionen; Permanente Links; Angabe zu den Standards und Normen



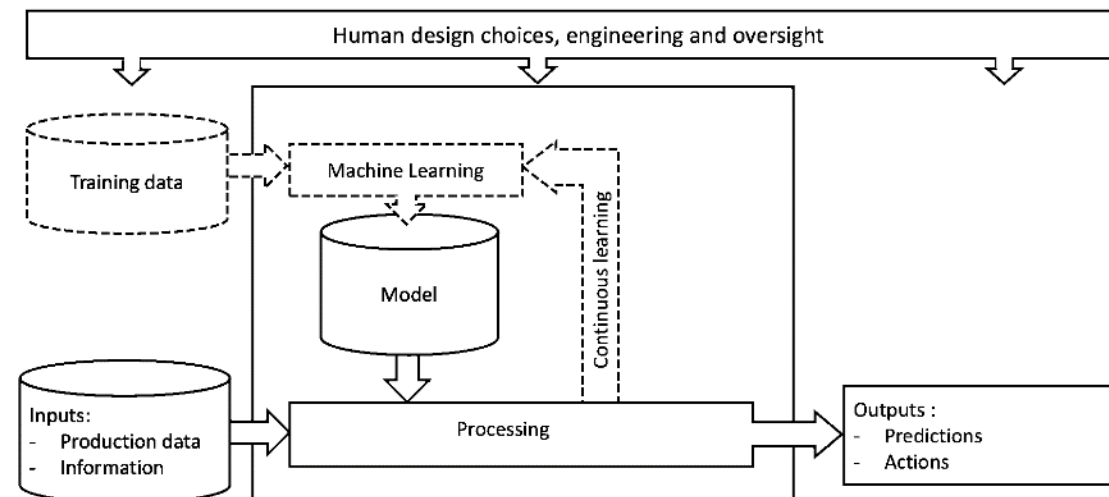
The screenshot shows a web browser window with the URL <https://www.ai-glossary.org/index.p>. The page title is "training data". The main content area is titled "Trainingsdaten • *training data*". It includes a synonym "auch: Lerndaten • *learning data*". The definition states: "Daten, die im Trainingsprozess zum Erstellen eines KI-Modells verwendet werden können." It also lists "Siehe auch:" with a bullet point for "Datensatz". Two ISO/IEC standards are cited: "ISO/IEC DIS 22989 subset of input data samples used to train a machine learning model (3.2.11)" and "ISO/IEC TR 29119-11:2020 dataset used to train an ML model (3.1.46)". At the bottom of the definition box, there are links for "Exportieren" and "Permalink". A warning message below reads: "Hinweis: Die nachfolgenden Treffer stimmen nur näherungsweise mit der Eingabe überein." Below this, the next entry is titled "Trainingsprozess • *training*".

2. Vorstellung der Normungs- und Standardisierungsbedarfe

Bedarf 02-05

Abstrakte Zerlegung der KI-Komponente in Daten und Prozesse

- ISO/IEC 22989:2022 macht einen ersten Versuch, dieser ist aber auf sehr hoher Abstraktionsebene.



ISO 22989 – p.41 – Figure 5

- ISO/IEC 27090:draft versucht KI-spezifische Angriffe zu beschreiben, hier könnte ein detaillierteres Architekturdiagramm eingebracht werden, um die beschriebenen Angriffe/Risiken hinsichtlich zu lokalisieren

2. Vorstellung der Normungs- und Standardisierungsbedarfe

Bedarf 02-06

Existierende KI-Angriffe und Risiken mit existierenden zertifizierbaren IT-Sicherheitszielen abgleichen

- Abbildung von Angriffen auf KI-System als Angriffe auf IT-Sicherheits-Schutzziele entsprechend einer Beschreibung der schutzwürdigen KI-Komponenten und des Lebenszyklus

- Beispiel

Angriff: “Data Poisoning“

→ Verletztes Schutzziel: Violation integrity of trainingdata until and during „training phase“

→ Control: Integrity, Authenticity (bzw. Data Governance) für Trainingsdaten

Angriff: „Membership inference attack“

→ Verletztes Schutzziel: Violation of confidentiality of trainingdata during deployment phase

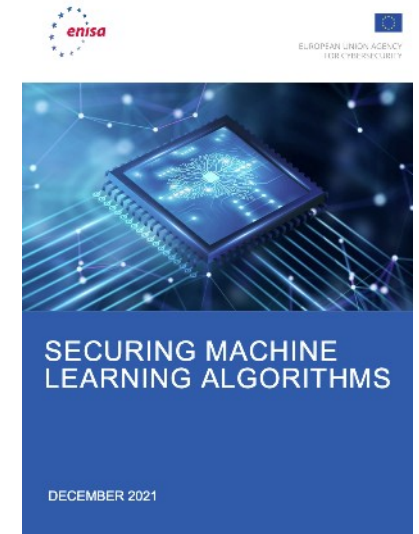
- ISO/IEC 27090:draft versucht KI-spezifische Angriffe zu beschreiben, hier könnte ein erstes Mapping versucht werden
- Ähnlich ISO/IEC 27091:draft für Privacy

2. Vorstellung der Normungs- und Standardisierungsbedarfe

Bedarf 02-07

Standardisierung von KI-Produkt- und Prozessprüfverfahren für Security und Privacy

- Wir brauchen für eine Zertifizierung klare Messverfahren und vor allem prüfbare Messwerte für die Security- und Privacy-Eigenschaften von KI-Systemen
- Prüfverfahren und Akkreditierungsverfahren (für die Prüfenden) sind essenziell, um die Qualität der Prüfung durch unabhängige Dritte sicherzustellen
- Passende generische Ansätze (ggf. ISMS Ansätze aus der IT-Sicherheit übernehmen)
- Zusätzlich spezielle Prüfungen abhängig von der KI Methode oder dem speziellen Algorithmus
- Es gibt erste Ansätze z.B. für Machinelles Lernen von der ENISA : <https://www.enisa.europa.eu/publications/securing-machine-learning-algorithms>



2. Vorstellung der Normungs- und Standardisierungsbedarfe

Bedarf 02-08

Ausarbeitung eines horizontalen Querschnittsstandards und vertikale Ausprägungen zu Security

Empfehlenswert wäre die Herausarbeitung von

horizontalen Themen zu Cybersecurity und Privacy für KI zur Prüfung und Zertifizierung, die alle Sektoren betreffen, sowie eine Schnittstelle zu sektorspezifischen Anforderungen.

Ein horizontales Thema wäre beispielsweise die Anforderung an eine geeignete Zugriffskontrolle. Als vertikale Ausprägung können wiederum spezielle Security-Anforderungen aus dem sektoralen Umfeld angesehen werden, wie u. a. für den Bereich der Medizinprodukte.

Themenschwerpunkt: KI-Sicherheit bei IT-Systemen

2. Vorstellung der Normungs- und Standardisierungsbedarfe

Bedarf 02-09

Entwicklung von Metriken und Controls gemäß den Standardisierungsanforderungen des geplanten EU AI Act

Entwicklung von Standardisierung zu

Cybersecurity-Anforderungen aus dem AI Act für Metriken und Controls zur Messung und Vermeidung von Cyberangriffen sowie Methoden für Prüfung, Auditierung und Zertifizierung inklusive Anforderungen an die Kriterien für die Prüfmaßnahmen und Prüfenden.

Dabei erscheint es wichtig, eine gemeinsame Arbeitsgruppe mit den Gremien der Cybersecurity und KI in den Standardisierungsorganisationen von Deutschland, der EU und eventuell auch international zu etablieren.

3. Interaktive Priorisierung der Bedarfe (Conceptboard)

Interessierte für die Bedarfsumsetzung

Während des Workshops wurde Conceptboard verwendet. Das Board ist nun geschlossen. Sie können sich aber noch per E-Mail an Jan.Roesler@din.de für die Umsetzung der Bedarfe melden – bitte Bedarfs-Code mit angeben.

3. Interaktive Priorisierung der Bedarfe (Mentimeter)

Priorisierung der Bedarfe für den Workshop – Ergebnis



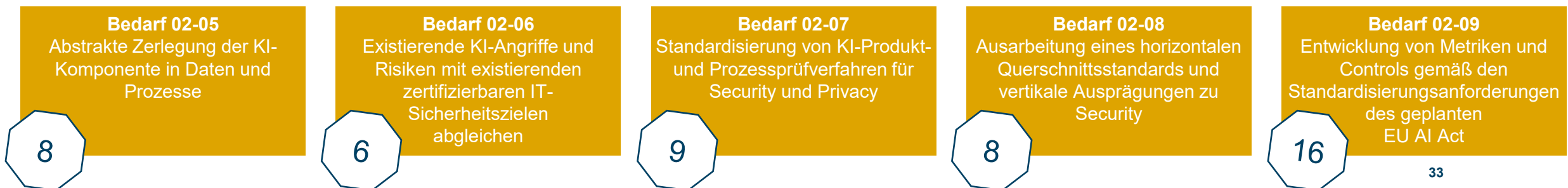
Für die Abstimmung konnten alle 5 Bedarfe von 1 bis 5 geordnet werden.

4. Diskussion der Bedarfe anhand der Priorisierung

Blitzlicht zur Diskussion der Bedarfe

- Die Bedarfe 02-09, 02-07 und 02-06 wurden diskutiert und ein gemeinsames Verständnis geschaffen. Allgemein gilt, dass die Bedarfe nicht losgelöst voneinander betrachtet werden können und über die Grenzen der IT-Sicherheit hinaus Schnittmengen mit anderen Bedarfen der Roadmap und anderen Normen/Standards haben. Daher soll verstärkt nach Überschneidungen (wie funktionale Sicherheit, Metriken) gesucht werden.
- Der Bedarf 02-09 erfordert die Entwicklung eines KI-Standards zur Robustness.
- Bedarf 02-07 ist als Ergänzung zu Bedarf 02-09 anzusehen und stellt ein Querschnittsthema dar. Folgende Fragestellungen sind zu beantworten:
 - Welche Standards gibt es hier schon und was ist jetzt (durch KI) anders und muss berücksichtigt werden?
 - Welche Änderungen sind nötig, damit Zertifizierungen inkl. KI möglich werden?
 - Welche Ansätze können aus *Information Security Management System* (ISMS) übernommen werden?
 - Welche Rolle spielt KI-Engineering und andere „KI-Methoden“ mit gezielten Eigenschaften → Cyber Resilliance Act
- Existierende Normen und Standards zu IT-Sicherheit ohne Bezug zu KI bleiben bestehen, müssen aber durch weitere KI-Standards (gerade in Bezug zu KI in IT-Sicherheit) ergänzt werden. Da keine aktueller Standard für IT-Sicherheit auch für KI-Systeme angewandt werden kann. Diese müssen speziell für die Anwendung auf KI-Systeme entwickelt werden.

Aktuelle Anmeldezahlen für die weitere Umsetzung (Nachmeldungen bitte an kuenstliche.intelligenz@din.de)



5. Nächste Schritte

Wie geht's weiter?

Aufbereitung der
WS-Ergebnisse

Rücksprache mit
Normenausschüssen

Gründung geeigneter
Arbeitsgruppen

Einladung zu den
nächstens Meetings

Bis bald und auf Wiedersehen!

Jan Rösler

Projektmanager

DIN Strategische Themenentwicklung KI

Jan.Roesler@DIN.de → Anmeldung für Bedarfsumsetzung

DIN

Deutsches Institut für Normung e. V.

Am DIN-Platz

Burggrafenstraße 6

10787 Berlin

www.din.de



DIN