

Towards Safe AI for Automated Driving

Fabian Hüger, Volkswagen & CARIAD
CSCS 2021 (online), November 30, 2021



The results, opinions and conclusions expressed in this publication are not necessarily those of Volkswagen Aktiengesellschaft.

Agenda

1.

Introduction – CARIAD

2.

DNNs and Safety in Automated Driving

3.

KI-Absicherung Project & Approach

4.

Consequences

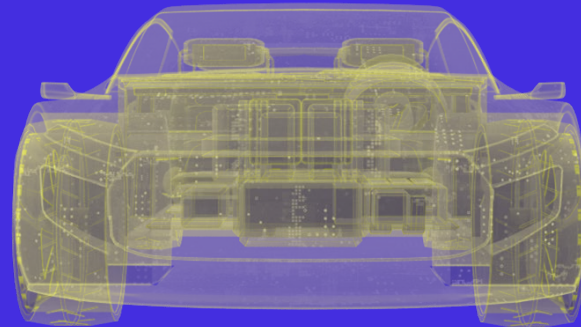
The car needs to be rethought

Connectivity

Making cars and mobility part of our customers' digital life.

Software driven

Rethinking the car from a software perspective, turning it into an intelligent companion.



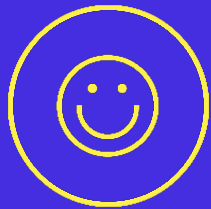
Electrification

Making mobility more sustainable.

Autonomous driving

Making cars safer and more comfortable for everyone.

CARIAD is here to make automotive mobility safer, more sustainable, and more comfortable.



Comfort

From enjoying the ride to enjoying digital life in your car – everything will become easier, more convenient, and more fun to use.



Safety

Automated and assisted driving will be much safer than any human at the steering wheel.



Sustainability

Continuous software updates keep our cars fresh for many years. Our smart navigation features save kilometers and resources, while reducing congestion.

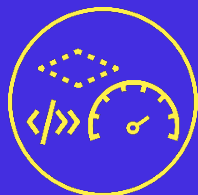
Our software platform delivers it all.

One software platform. Lots of benefits.



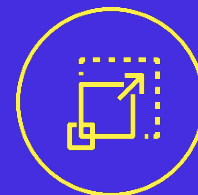
Updatability

Constant and efficient updatability enables attractive vehicles and the best, always fresh customer experiences.



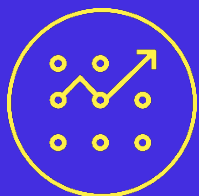
Speed

The seamless software platform and intelligent data analysis speed up development and time to market.



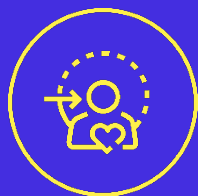
Scalability

The digital platform suits any car model – from entry-level to top-end. Applications can easily be customized.



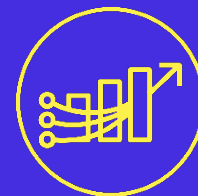
Simplicity

One unified platform reduces complexity – and less hardware reduces costs and weight.



Customer orientation

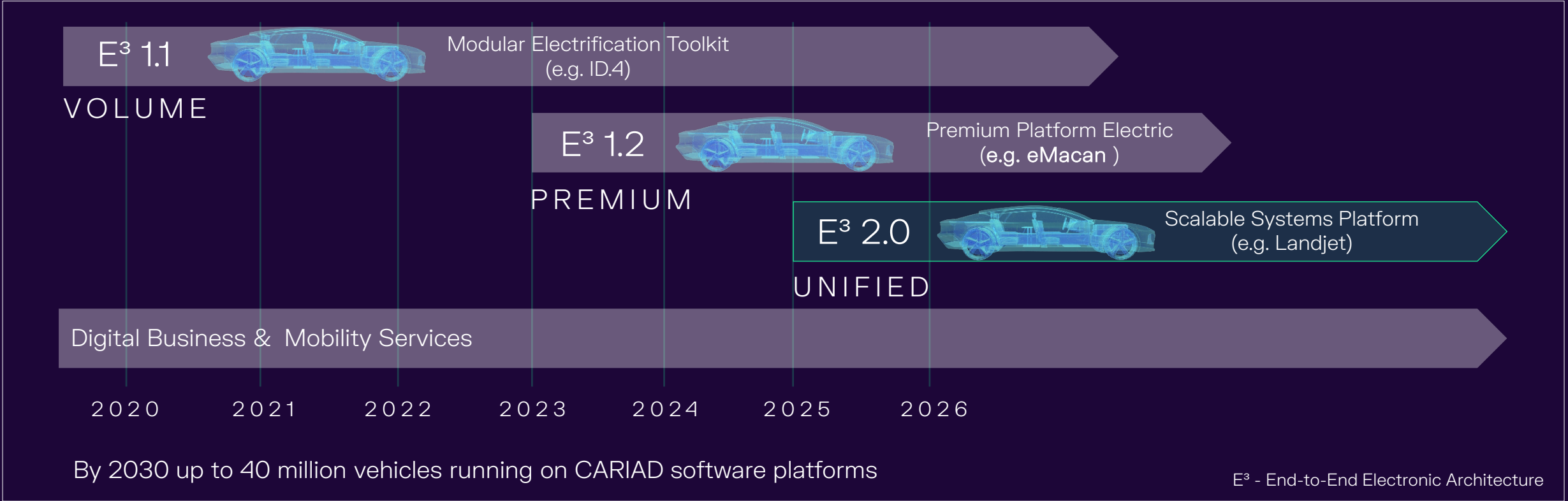
Data-oriented development helps us to learn from and react to customers' needs and desires.



New revenue streams

Car brands can generate new digital business models – from after sales to monetizing data or third-party apps.

Our platforms E³ 1.1 and E³ 1.2 are technological front runners, while E³ 2.0 will be the one platform in the Group starting 2025.



Agenda

1. Introduction – CARIAD
2. DNNs and Safety in Automated Driving
3. KI-Absicherung Project & Approach
4. Consequences

Automated Driving and AI

Processing chain of autonomous driving & the use of AI along



- Near real time (20Hz)
- Multi sensor
- Redundancy

UTILIZATION OF AI



Arguing Safety in Automated Driving Systems

AI goes safety critical

CENTRAL CHALLENGE

SAFETY

(FuSa + SOTIF)

Central Challenge in bringing highly automated driving on the road.

Argument on safe functioning needed to allow for acceptance & road permission

COMPLEXITY DRIVERS



Mere driving will not suffice to plausibilize safety – particularly challenging with respect to software updates over time. “Black-Box” approach seems impracticable



Handling complexity of the driving environment – open world, unknown unknowns, etc.



Need for continual safety monitoring & assurance – continuous monitoring

Agenda

1. Introduction – CARIAD
2. DNNs and Safety in Automated Driving
- 3. KI-Absicherung Project & Approach**
4. Consequences

Acknowledgement: The research leading to these results is funded by the German Federal Ministry for Economic Affairs and Energy within the project "Methoden und Maßnahmen zur Absicherung von KI basierten Wahrnehmungsfunktionen für das automatisierte Fahren (KI-Absicherung)". The authors would like to thank the consortium for the successful cooperation.



KI-Absicherung Project & Approach

The results, opinions and conclusions expressed in this publication are not necessarily those of Volkswagen Aktiengesellschaft.

Gefördert durch:



Bundesministerium für Wirtschaft und Energie

aufgrund eines Beschlusses des Deutschen Bundestages



Making the safety of AI-based
function modules for highly
automated driving verifiable

KI ABSICHERUNG

Safe AI for Automated Driving

Pedestrian detection

Challenge

AI Land



Promising new technology with unimagined possibilities

Established safety processes cannot be applied



Safety Land



Safe, trustworthy driving function



Industry consensus (Safe AI): Methodology for joint safety argumentation

Our Team: Experts from AI, Safety and Virtual Reality



OEMs



Tiers



Technology Provider



Research Institutes



University



External Partners



Consortium
Lead:
Volkswagen AG

Co - Lead:
Fraunhofer IAIS

Budget:
41 Mil. €

BMW Funding:
19.2 Mil. €

24 Partners

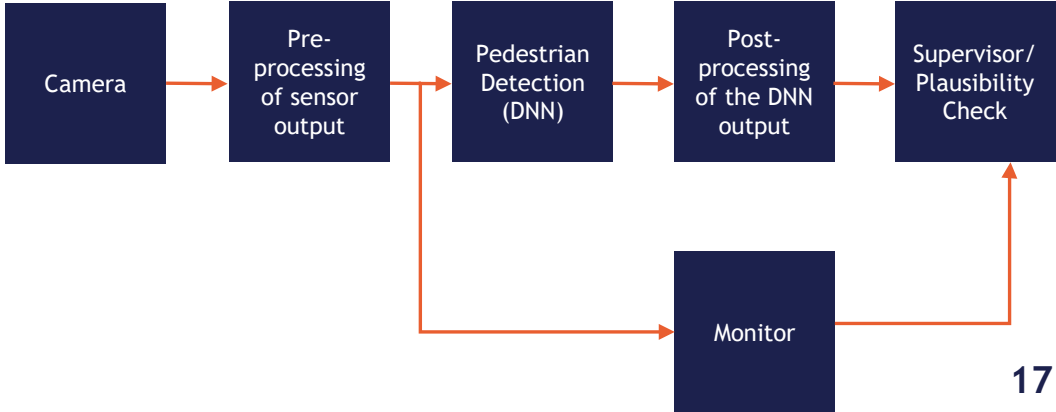
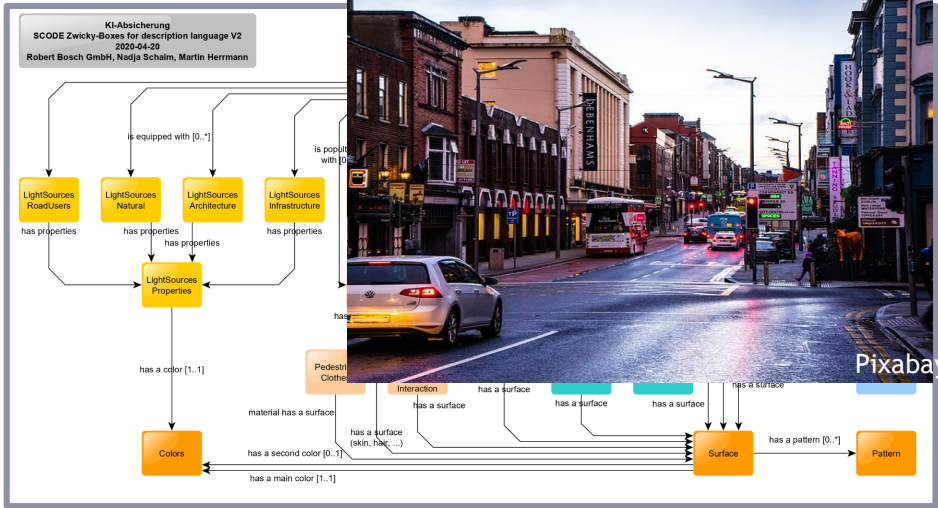
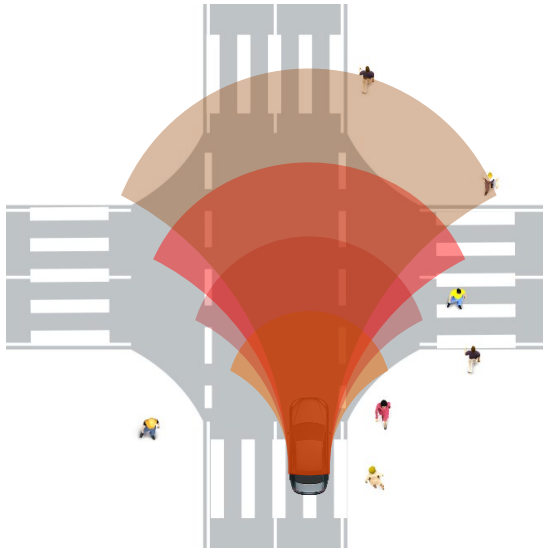
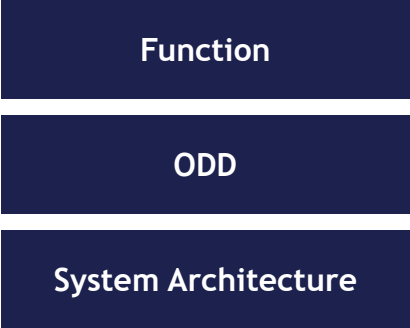
Duration: 36
month

01.07.2019 -
20.06.2022

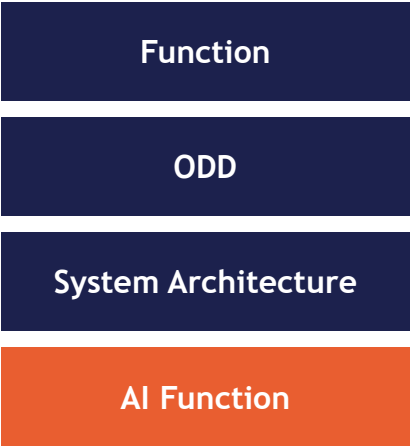
Gefördert durch:
 Bundesministerium
für Wirtschaft
und Energie

aufgrund eines Beschlusses
des Deutschen Bundestages

Our Approach: Specification



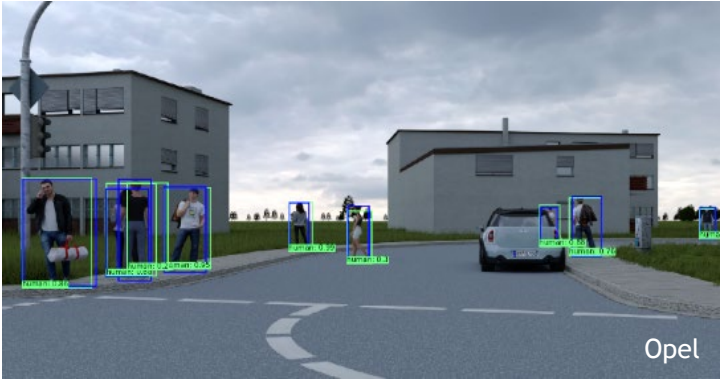
Our Approach: AI Function Pedestrian detection



Semantic Segmentation



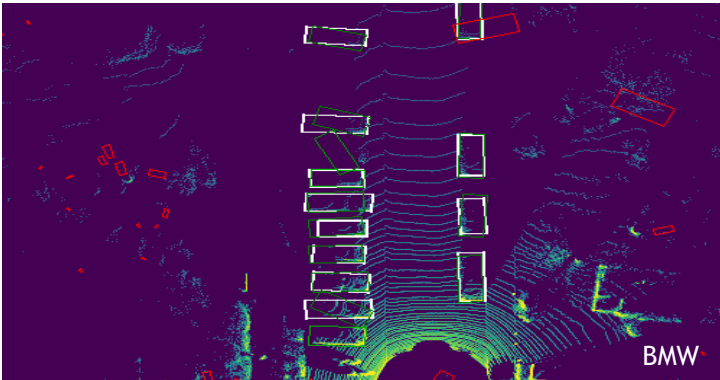
2D Bounding Box Detection



Instance Segmentation



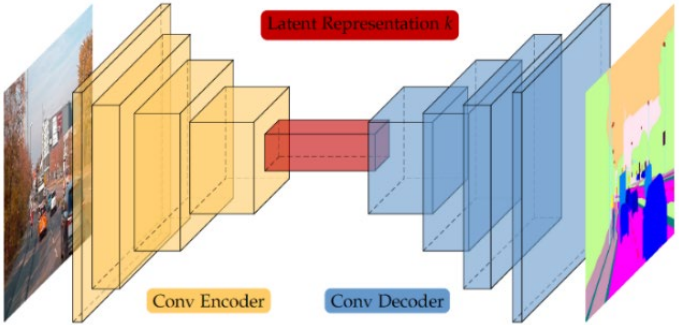
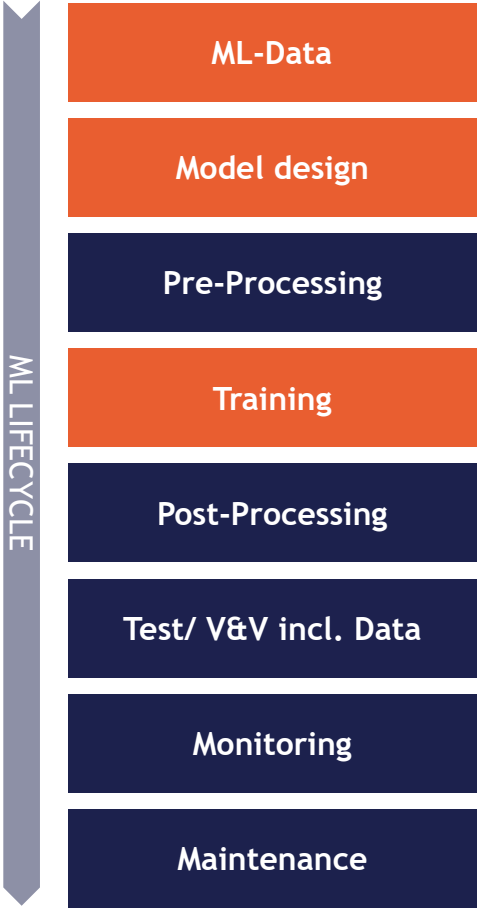
3D Bounding Box Detection



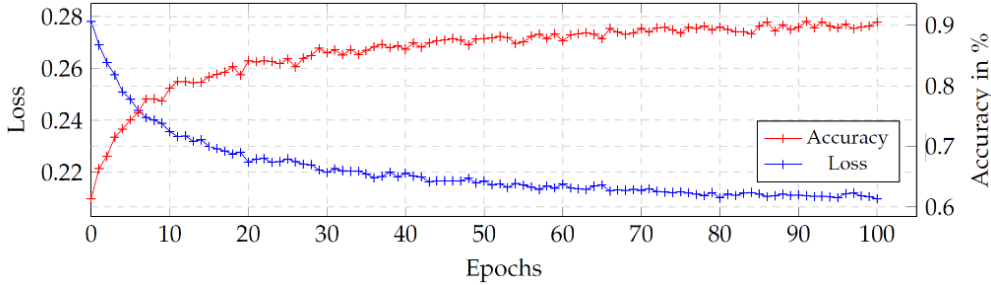
3D Pose estimation



Our Approach: Synthetic Data and ML-Lifecycle

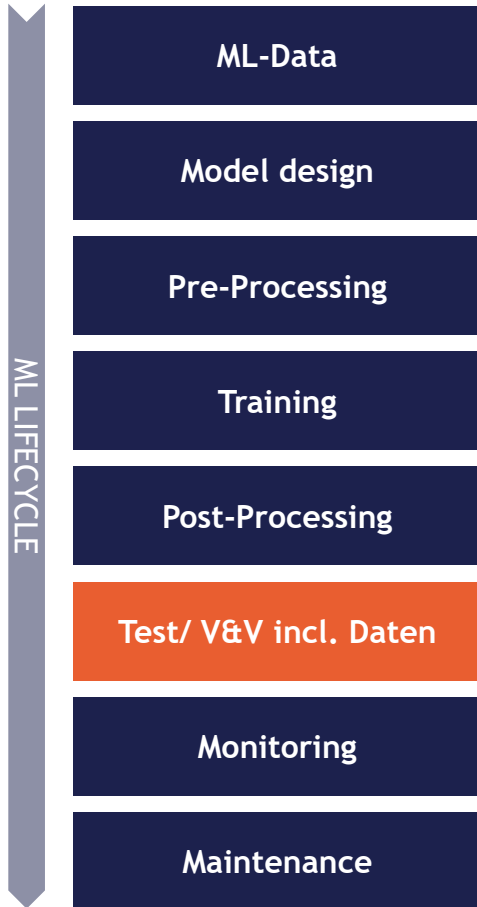


Volkswagen AG



Volkswagen AG

Our Approach: ML-Lifecycle-Validation data

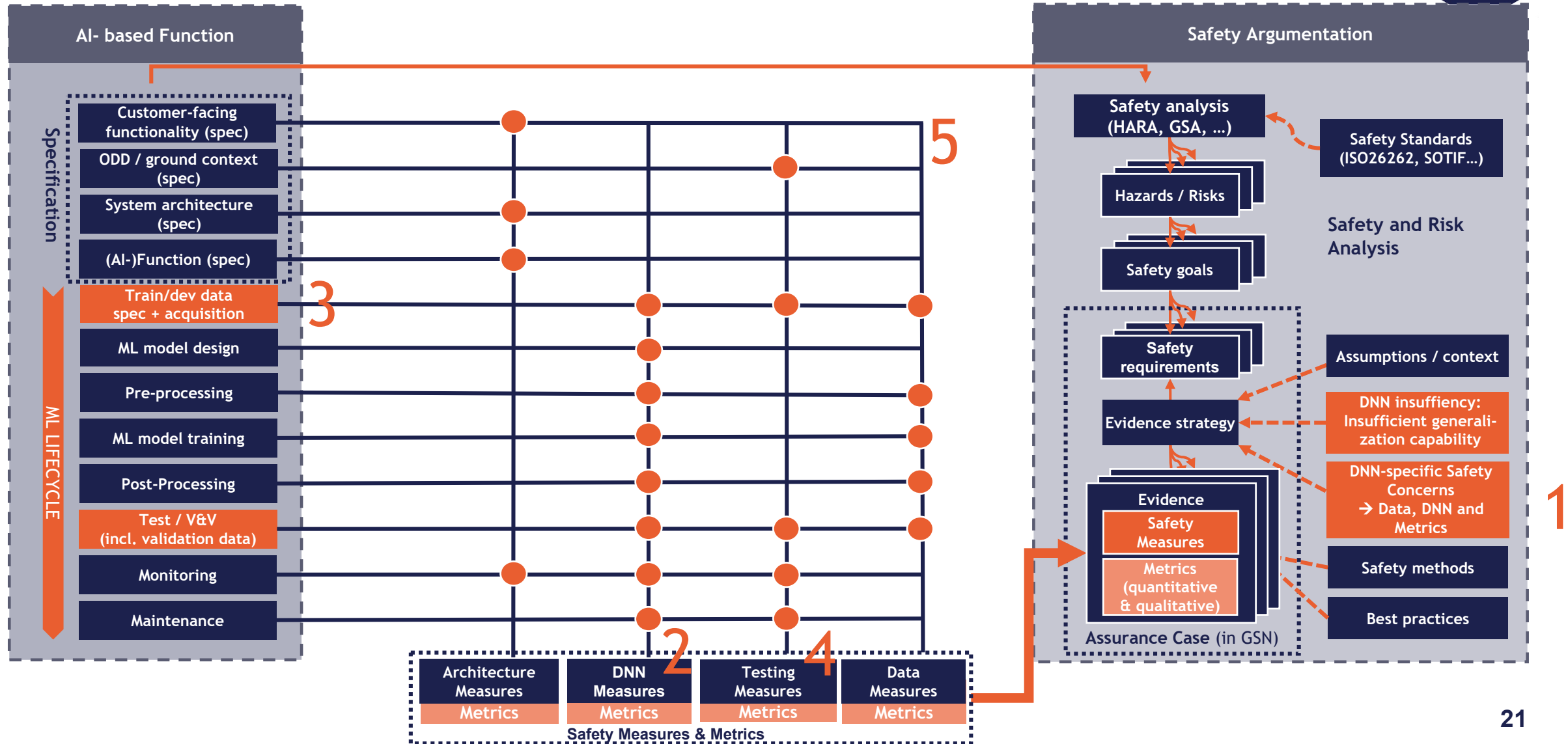


Continuous process for identification, specification and generation of synthetic data





Our Approach: Big picture





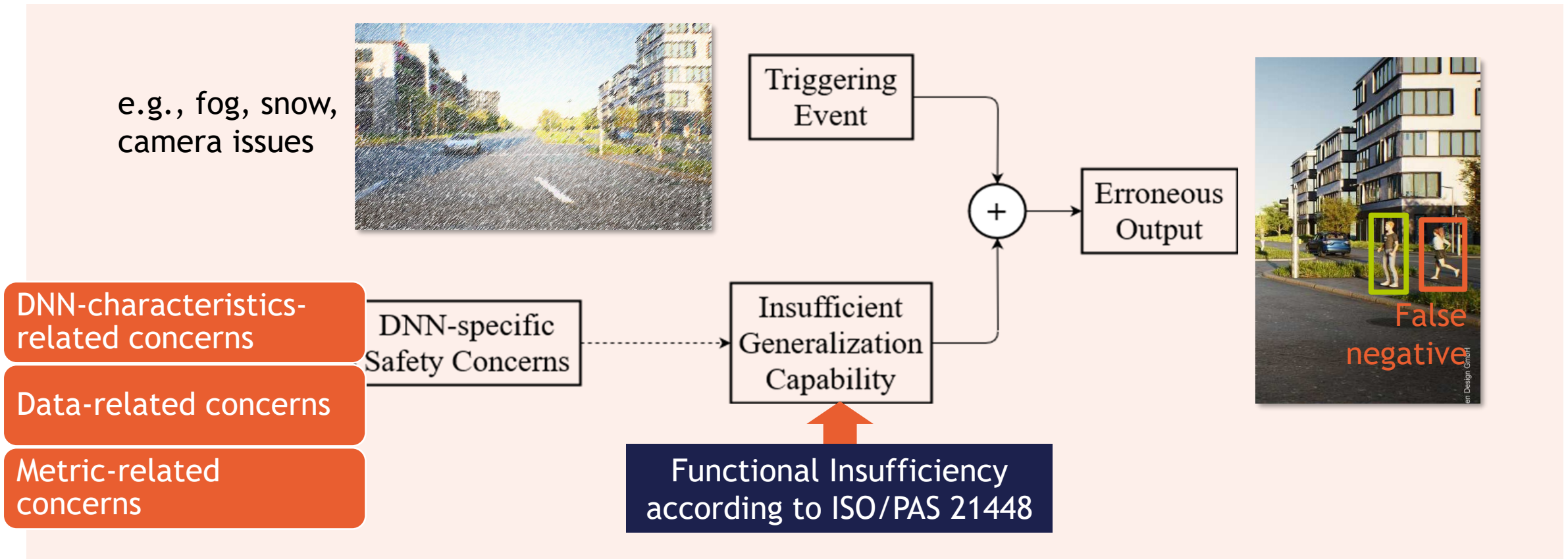
1

DNN-specific safety concerns



Our Approach: DNN-specific Safety Concerns (1/2)

We define **DNN-specific Safety Concerns (SCs)** as underlying issues of DNN-based perception which may negatively affect the safety of a system.





Based on:

O. Willers, S. Sudholt, S. Raafatnia, S. Abrecht: Safety Concerns and Mitigation Approaches Regarding the Use of Deep Learning in Safety-Critical Perception Tasks

T. Sämann, P. Schlicht, F. Hüger: Strategy to Increase the Safety of a DNN-based Perception for HAD Systems

G. Schwalbe, B. Knie, T. Sämann, T. Dobberphul, L. Gauerhof, S., V. Rocco: Structuring the Safety Argumentation for Deep Neural Network Based Perception in Automotive Applications

Functional Insufficiencies

DNN-characteristics-related concerns

Data-related concerns

Metric-related concerns

FI-1 INSUFFICIENT GENERALIZATION CAPABILITY

Wrong outputs by an AI-based function that was trained on a limited database. Erroneous input to output mapping or wrong approximation.

SC-1.1 UNRELIABLE CONFIDENCE INFORMATION

DNNs tend to be overconfident in their predictions under certain conditions or in general outputting unreliable confidence information.

SC-1.2 BRITTLINESS OF DNNs

Non-robustness against common perturbations such as noise or certain weather conditions as well as targeted perturbations known as adversarial examples

SC-1.2.1 LACK OF TEMPORAL STABILITY

Detection results rapidly changing in time whereas little change occurs in the ground truth

SC-1.3 INCOMPREHENSIBLE BEHAVIOUR

Inability to explain exactly how DNNs come to a decision.

SC-1.4 INSUFFICIENT PLAUSIBILITY

AI based functions usually lack basic plausibility checks, which are intended to identify detections of the perception function that violate physical laws.

SC-2.1 DATA DISTRIBUTION IS NOT A GOOD APPROXIMATION OF REAL WORLD

The distribution of data used in the development should be a valid approximation of the ODD in the real world.

SC-2.2 INADEQUATE SEPARATION OF TEST AND TRAINING DATA

Test data might be correlated to training data which might induce overfitting on test data.

Technologies Assessment

SC-2.3 DEPENDENCE ON LABELLING QUALITY

Labelling quality can directly affect the resulting model performance. Moreover, due to missing labelling quality, evaluation results might be misleading.

SC-2.3.1 MISSING LABEL DETAILS OR META-LABELS

Missing meta-labels or label details possibly leads to improper data selection or insufficient training objectives.

SC-2.4 SPECIFICATION OF THE ODD

An incomplete or incorrect ODD specification leads to incomplete data records for training and testing.

SC-2.5 DISTRIBUTIONAL SHIFT OVER TIME

A DNN is trained and tested at a certain point in time. Changes will occur naturally and therefore can potentially harm the performance of DNNs.

SC-2.6 UNKNOWN BEHAVIOUR IN RARE CRITICAL SITUATIONS

The long tail problem describes the fact that there exists an enormous amount of possibly safety-critical street scenes that have a low occurrence probability.

SC-3.1 SAFETY-AWARE METRICS

Some state-of-the-art metrics only evaluate the average performance of DNNs. Safety-aware metrics are required to sophisticatedly evaluate the performance of DNNs.

DNN-specific Safety Concerns



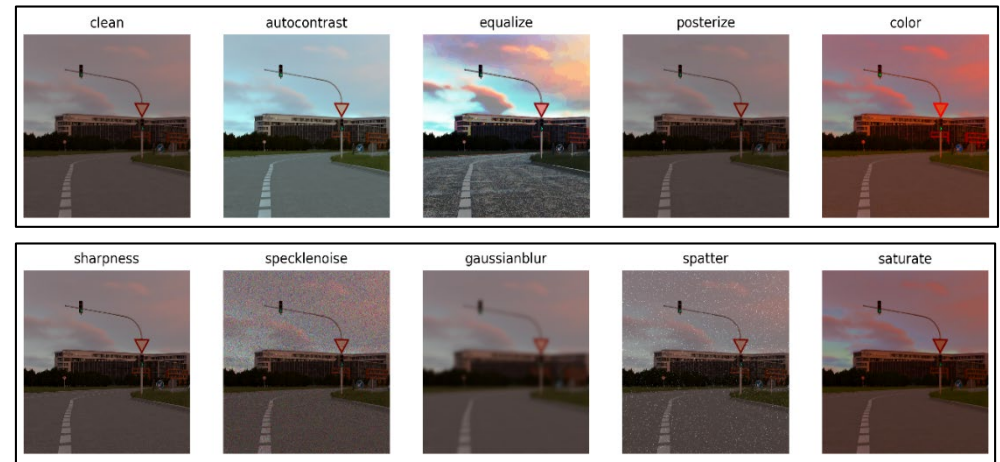
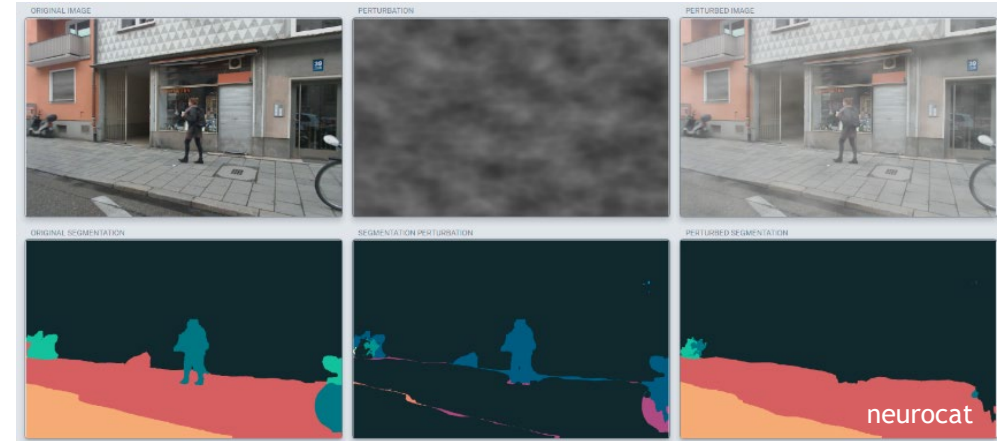
2

Exemplary Measures

Our Approach: Identify, Measure and Counteract „DNN-specific Safety Concerns”

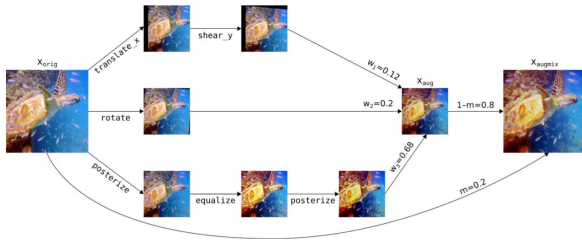
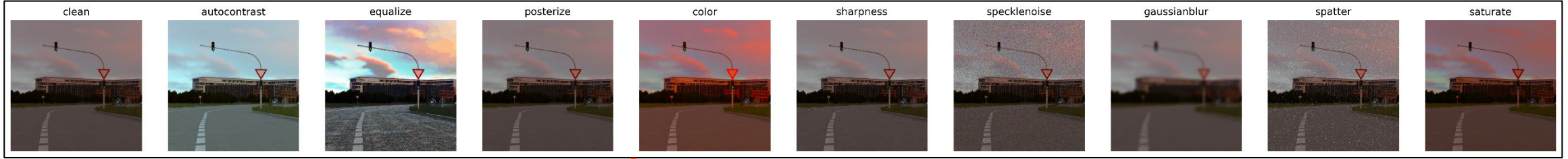
Addressed Safety Concern:
Brittleness of DNNs

- Addressing “Brittleness of DNNs” (Example)
 - **Requirement:** Robustness = Performance even under reasonable perturbations (gained from ODD definition, data analysis and sensor specs)
 - **Metric:** Performance under corruption
 - **Methods (e.g.)**
 - Augmentation Training (**AugMix**)
 - From a Fourier-Domain Perspective on Adversarial Examples to a **Wiener Filter** Defense for Semantic Segmentation
 - **Evidence:** Effectiveness of measure via metric



Our Approach: Identify, Measure and Counteract „DNN-specific Safety Concerns” via AugMix

Addressed Safety Concern:
 Brittleness of DNNs
 Corruption Robustness



Combined using AugMix

- + Improved robustness
- + Improved generalization
- + Data efficient augmentation strategy

AUGMIX: A SIMPLE DATA PROCESSING METHOD TO IMPROVE ROBUSTNESS AND UNCERTAINTY

Dan Hendrycks*
 DeepMind
 hendrycks@berkeley.edu

Norman Mu*
 Google
 normanmu@google.com

Ekin D. Cubuk
 Google
 cubuk@google.com

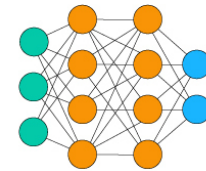
Barret Zoph
 Google
 barretzoph@google.com

Justin Gilmer
 Google
 gilmer@google.com

Balaji Lakshminarayanan†
 DeepMind
 balajiln@google.com



Training



Evaluation on 14 *unseen* „real-world“ corruptions

Our Approach: Identify, Measure and Counteract „DNN-specific Safety Concerns” via **AugMix**

Adressed Safety Concern:
Brittleness of DNNs
Corruption Robustness

Augmented Image

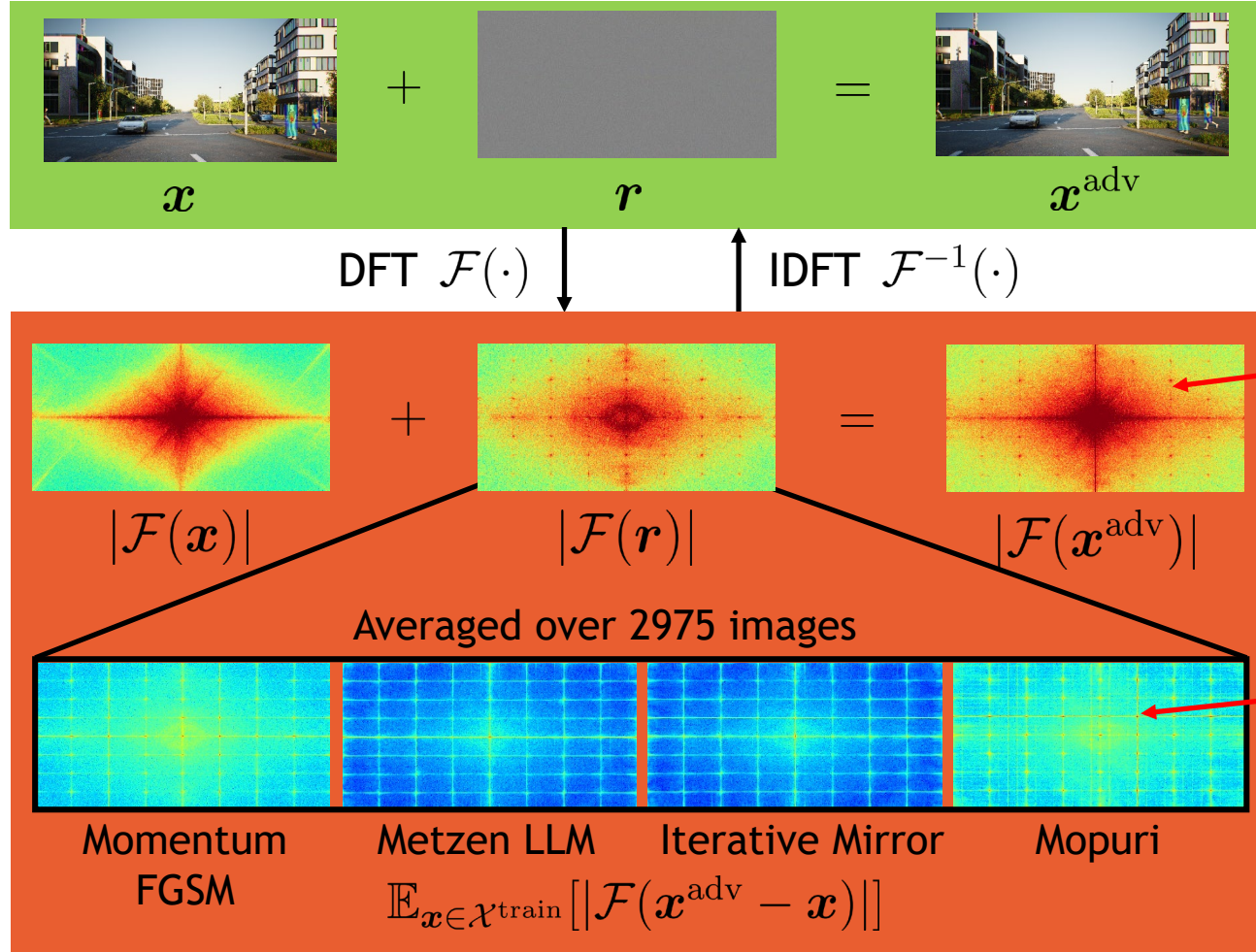
Baseline Segmentation

Defended Segmentation



Our Approach: Identify, Measure and Counteract „DNN-specific Safety Concerns”

Addressed Safety Concern:
Brittleness of DNNs
Adversarial Attacks



Adversarial examples are imperceptible in the spatial domain

Strong visible artifacts in the frequency domain

These artifacts are image-type and attack-type independent

- Spatial domain
- Frequency domain

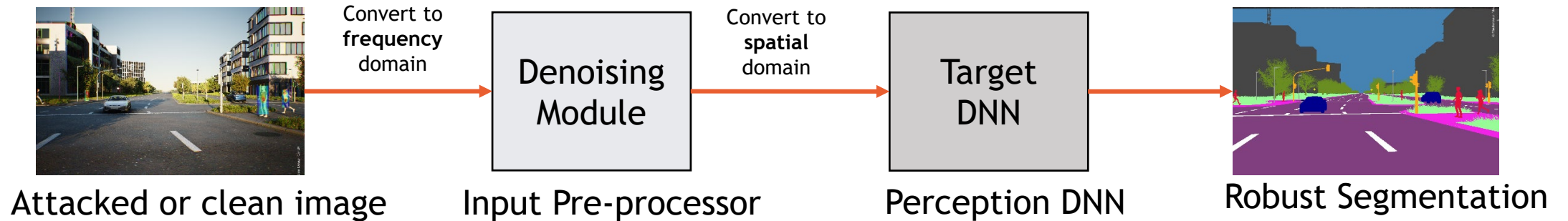
Our Approach: Identify, Measure and Counteract „DNN-specific Safety Concerns” via **Wiener Filters**

Addressed Safety Concern:
Brittleness of DNNs
Adversarial Attacks

Wiener Filters (WF) as an online denoising module

Steps:

1. Convert input image to DFT domain.
2. Apply pre-computed WF as a multiplicative filter.
3. Convert to spatial domain using IDFT.
4. Feed image to target DNN.



Our Approach: Explore Mechanisms!



- Heatmap-based Attention Consistency Validation
- Mixture of Experts
- Domain Randomization in Optimized Dataset Selection
- MC Dropout
- Uncertainties For Anomaly Detection
- Hybrid Learning using Concept Enforcement
- Active Learning
- Adversarial Training
- Hybrid and robustness-focussed Compression
- ...

Approx 80
Mechanisms are
developed and
evaluated

Inspect, Understand, Overcome: A Survey of Practical Methods for AI Safety

Sebastian Houben¹, Stephanie Abrecht², Maram Akila¹, Andreas Bär¹⁵, Felix Brockherde¹⁰, Patrick Feifel⁸, Tim Fingscheidt¹⁵, Sujan Sai Gannamaneni¹, Seyed Eghbal Ghobadi⁸, Ahmed Hammam⁸, Anselm Haselhoff⁹, Felix Hauser¹¹, Christian Heinzemann², Marco Hoffmann¹⁶, Nikhil Kapoor⁷, Falk Kappel¹², Marvin Klingner¹⁵, Jan Kronenberger⁹, Fabian Küppers⁹, Jonas Löhdefink¹⁵, Michael Mlynarski¹⁶, Michael Mock¹, Firas Mualla¹³, Svetlana Pavlitskaya¹⁴, Maximilian Poretschkin¹, Alexander Pohl¹⁶, Varun Ravi-Kumar⁴, Julia Rosenzweig¹, Matthias Rottmann⁵, Stefan Rüping¹, Timo Sämann⁴, Jan David Schneider⁷, Elena Schulz¹, Gesina Schwalbe³, Joachim Sicking¹, Toshika Srivastava¹², Serin Varghese⁷, Michael Weber¹⁴, Sebastian Wirkert⁶, Tim Wirtz¹, and Matthias Woehrle²

¹Fraunhofer Institute for Intelligent Analysis and Information Systems

²Robert Bosch GmbH

³Continental AG

⁴Valco S.A.

⁵University of Wuppertal

⁶Bayerische Motorenwerke AG

⁷Volkswagen AG

⁸Opel Automobile GmbH

⁹Hochschule Ruhr West

¹⁰umlaut AG

¹¹Karlsruhe Institute of Technology

¹²Audi AG

¹³ZF Friedrichshafen AG

¹⁴FZI Research Center for Information Technology

¹⁵Technische Universität Braunschweig

¹⁶QualityMinds GmbH

Survey: available at
<https://www.ki-absicherung-projekt.de/>



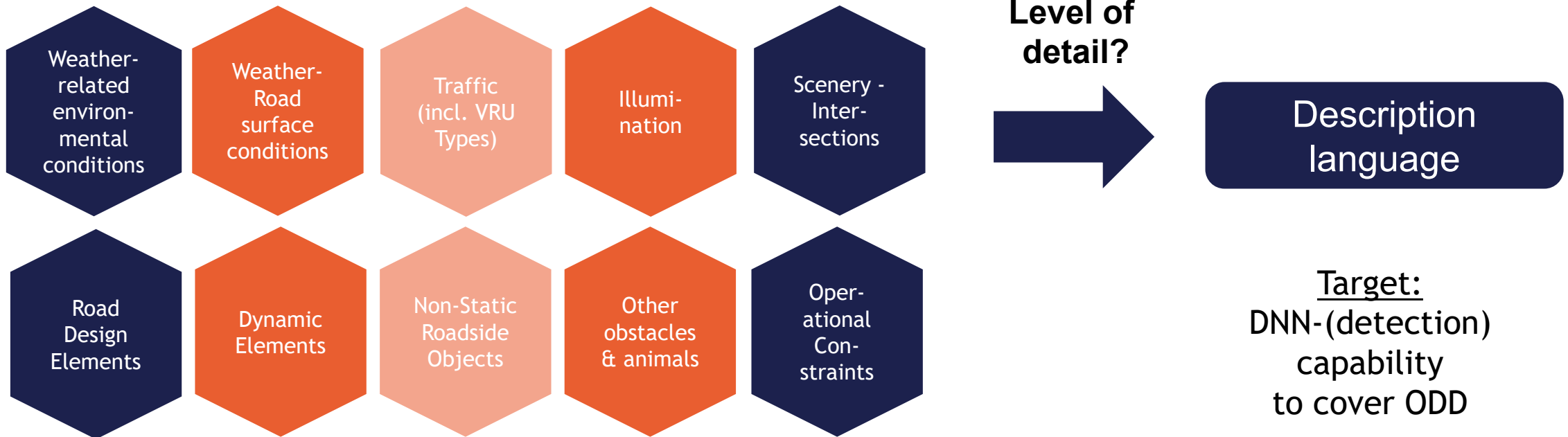
3

Our systematic approach for data

Operational design domain (ODD)



- An ODD describes / specifies operating conditions under which a given driving automation system or feature is specifically designed to function [...]
- Taxonomy and Definitions for Terms Related to Driving Automation Systems (examples)



A description language & data input space modeling is needed to...



Complexity of language



Be able to describe / **specify operating conditions** (and edges of ODD*) as of PAS 1883:2020 and others



Systematically capture important knowledge and describe the (expected) **key input space dimensions** and their **possible variations** having an influence on the functional performance of a DNN-based function (→ Zwicky Boxes & Ontology)



Perform training and assurance **data coverage estimations** for data driven AI-based systems



Describe **Corner cases / rare critical situations** to be considered in training / test data sets



For synthetic perception data production & meta-data: describe data dimensions that should be varied & **incrementally generate new data** by analyzing coverage and generating missing combinations

DNN-specific Safety Concerns (examples)



Data distribution is not a good approximation to target domain



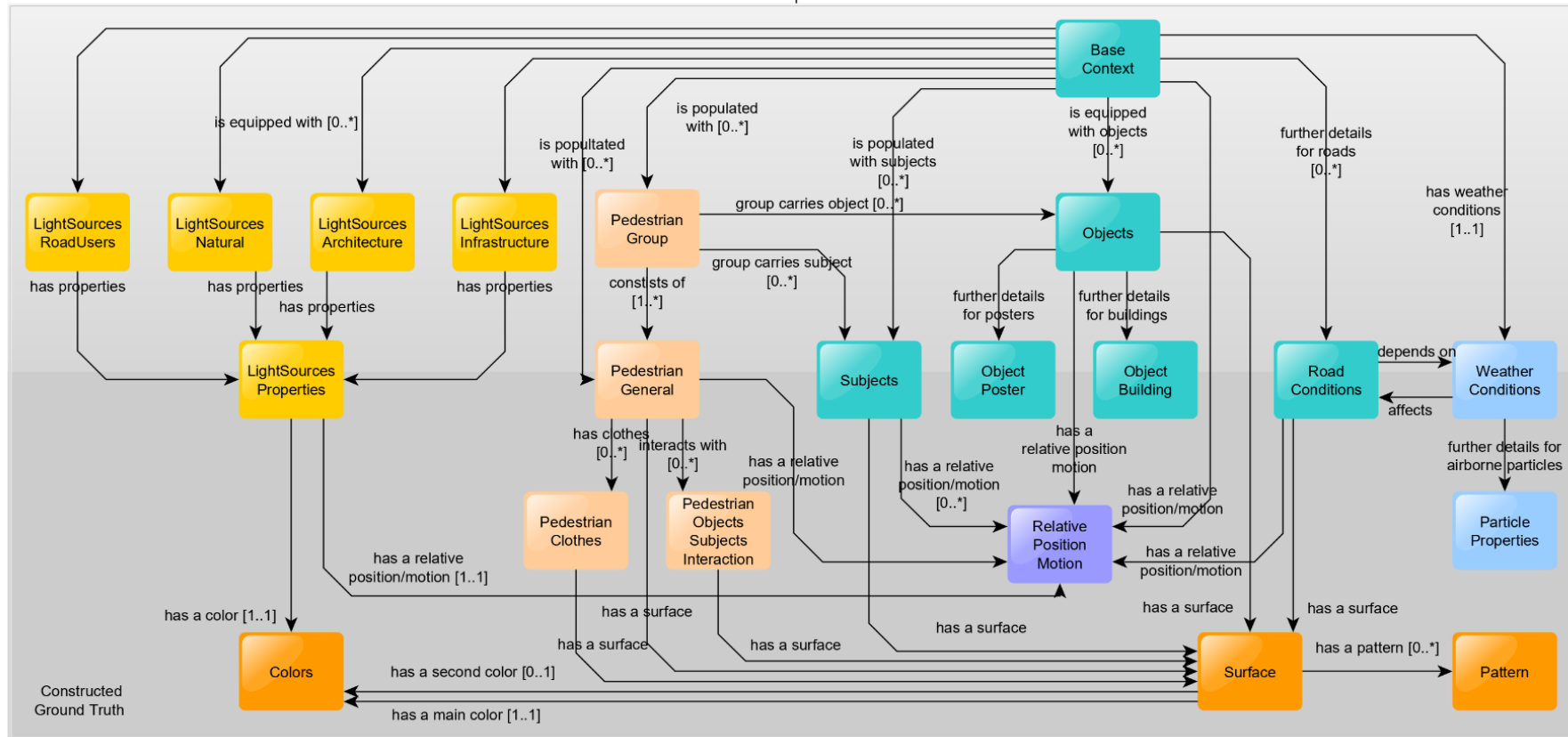
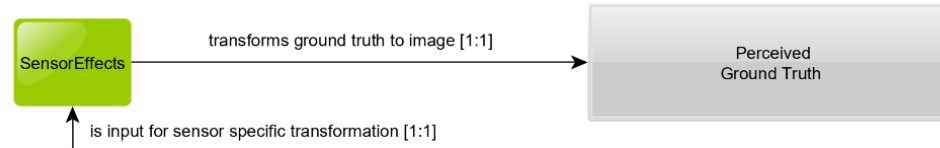
Unknown behavior in rare safety-critical situations

Based on: O. Willers, S. Sudholt, S. Raafatnia, S. Abrecht: Safety Concerns and Mitigation Approaches Regarding the Use of Deep Learning in Safety-Critical Perception Tasks

High Level view of Ontology / Domain model derived from SCODE Zwicky-Boxes



KI-Absicherung
SCODE Zwicky-Boxes for description language V5b
 2020-06-23
 Robert Bosch GmbH, Nadja Schalm, Martin Herrmann



Total

- ~250 dimensions
- ~1000 alternatives
- Several Sub-domains

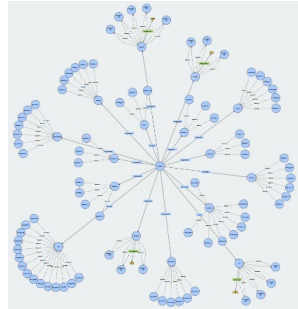
Approach

- Review of public data sources / existing standards
- Brainstorming with experts
- Expert interviews
- Iterative refinement
- Needs to be challenged / extended by identified corner cases



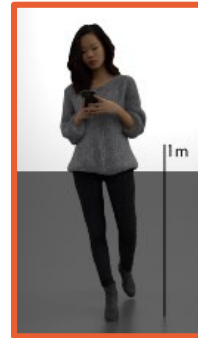
Data representations of the data input space aligned to ontology

Ontology Graph (Relations)



Visualization of KI Absicherung pedestrian sub ontology

Asset & Object descriptions for data analytics



Pedestrian:Age "adult"
 Pedestrian:BodyHeight "160cm-200cm"
 Pedestrian:BodyShape "thin"
 Pedestrian:BodyType "hourglass"
 Pedestrian:FaceShape "oval"
 Pedestrian:Gender "female"
 Pedestrian:HairColor "black"
 Pedestrian:HairLength "long"
 Pedestrian:HairStyle "other"
 Pedestrian:Pigmentation "medium"
 Pedestrian:Pose "walking"
 Pedestrian:SkinModification "no"
 Pedestrian:SpecialHandicap "no"

Source: BIT-TS

Systematic Combination of variations

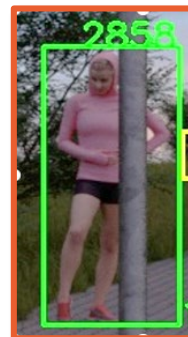
Dimension	Person 1	Person 2	Person 3	...
Age	Child	Teenager	Adult	
Gender	Male	Female	Male	
Body height	80-120 cm	120-160 cm	160-200 cm	
Pose	Running	Lying	Walking	
Pedestrian Location	Middle of street	Left side walk	Right side walk	
...	

Representations of variations

DAYTIME	morning	day	evening	night	
HAZE/FOG	no		yes		
STREET CONDITION	dry	wet	icy	snow	broken
SKY	cloudy		no	clear	
RAIN	no		yes		
REFLECTION ON ROAD	no		yes		
SHADOW ON ROAD	no		yes		
VRU TYPE	adult		child		
VRU POSE	pedestrian	jogger	cyclist		
VRU CONTRAST TO BG	low		high		

Zwicky Box - Discretized variations of important dimensions (Bosch)

Object GT Annotations for DNN-Training & Testing



Height = 55 px
 Width = 10 px
 Occlusion_level: 80%
 Occluded_body_part: arm
 Occluder: lamp
 Within_breaking_distance_30kph: true

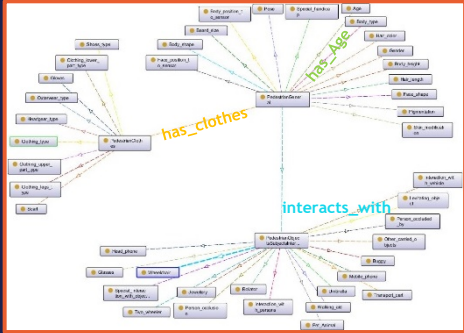
Source: BIT-TS

Systematically identify and describe the (known / expected) **key input space dimensions** and their **possible variations & combinations** having an influence on the functional performance of a DNN-based function

Structured Incremental dataset generation to boost data coverage (Vision)



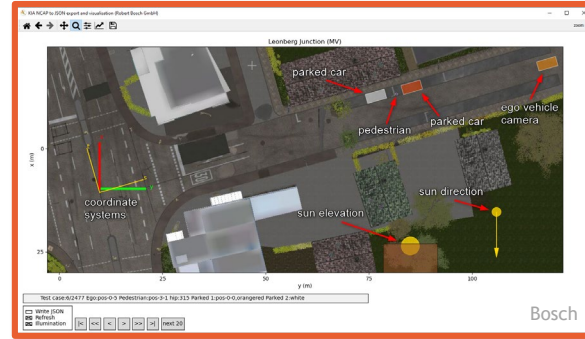
Domain model / data input space



Possible Variations based on Ontology



Data Request tool



JSON



Synth. Data production



Annotations & Meta data



Data Coverage Analysis

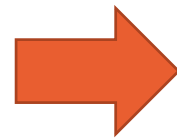
Dimension	Person1	Person2	Person3
Age	Child	Teenager	Adult
Gender	Male	Female	Male
Body height	80-120cm	120-160cm	160-200cm
Pose	Running	Lying	Walking
Pedestrian Location	Middle of street	Left side walk	Right side walk
...

Known Performance Limitations

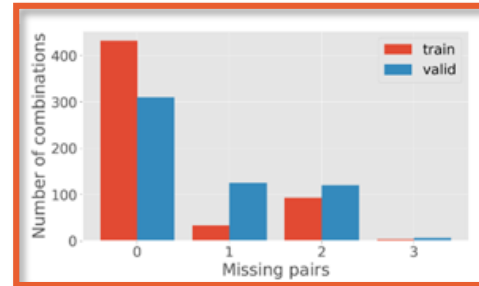
(expert know-how)

- Person occlusions
- Low contrast: similar color to background
- Uncommon person locations
- Uncommon poses

Constrained test-space



Missing Combinations



Optimized Combinatorial Testing (example)

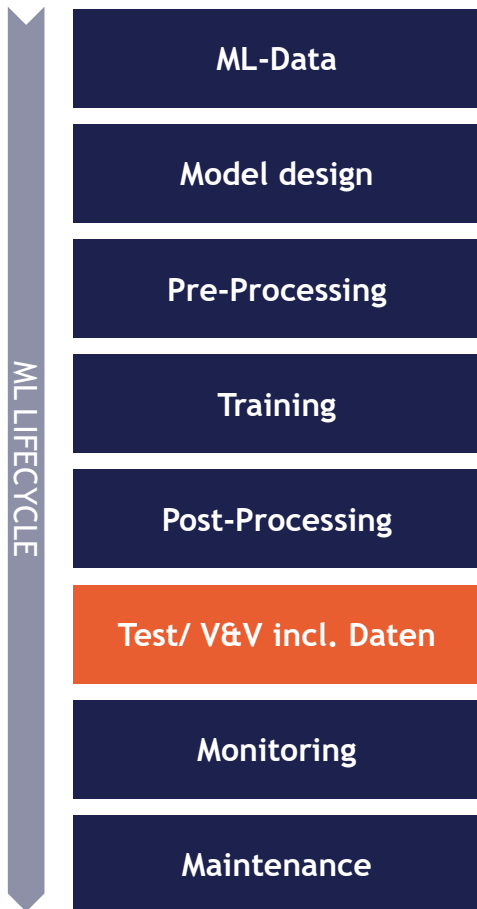




4

NCAP inspired test data production process

ML-Lifecycle-Validation data

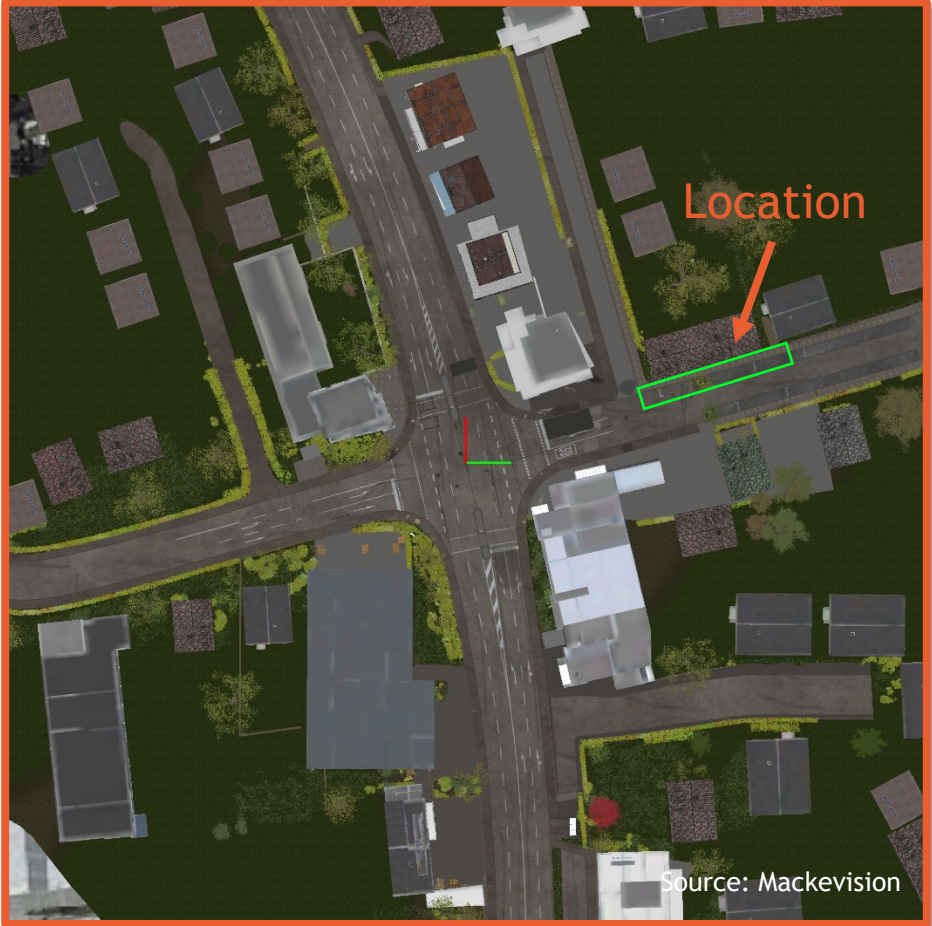
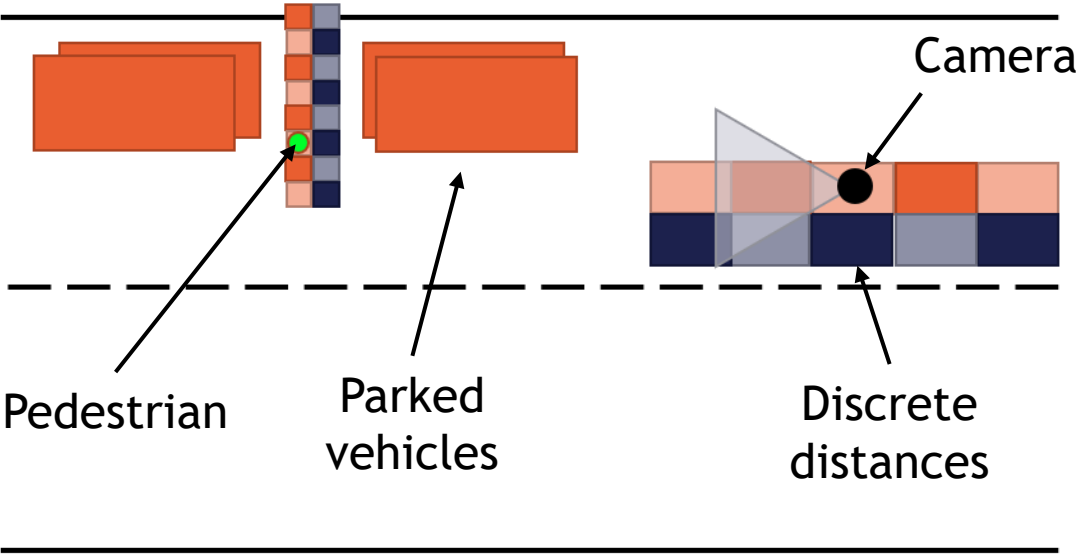


Definition of base scenario and location on base context



Story

A pedestrian is approaching the ego vehicle between two parking cars under different environment conditions



Discretization of dimensions in “Zwicky Boxes”

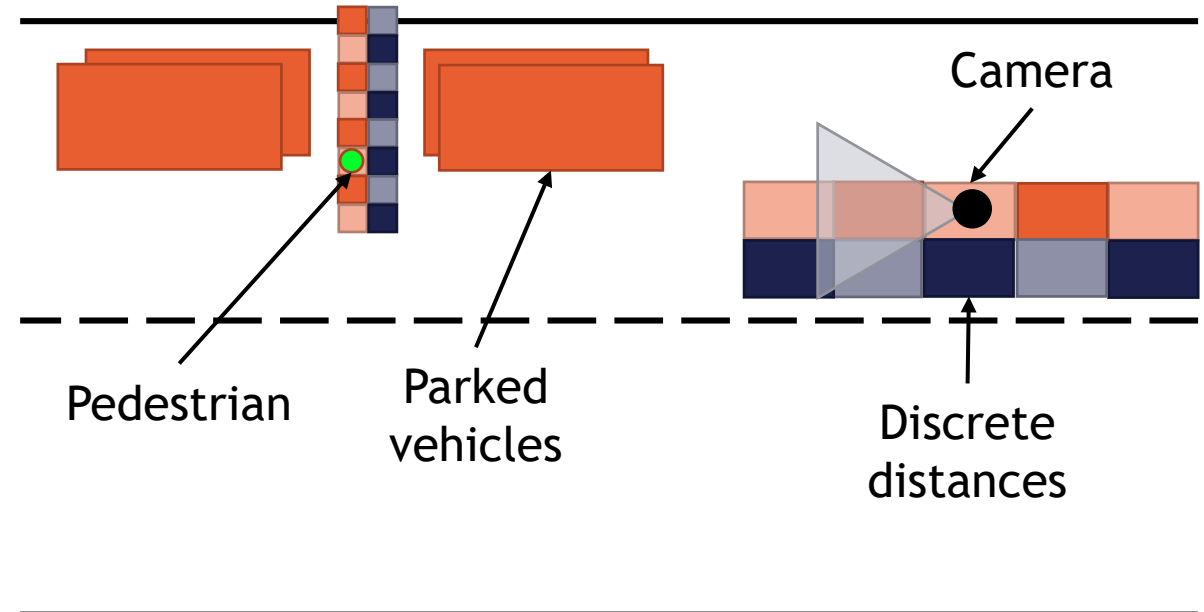
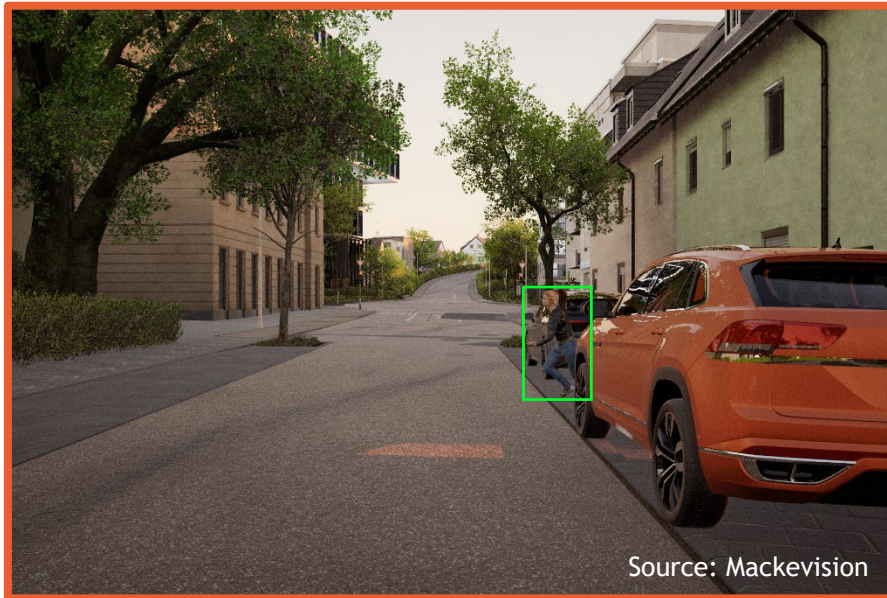


Ego XY position	pos-0-0		pos-0-1		pos-0-2		pos-0-3		pos-0-4		pos-0-5		pos-1-0		pos-1-1		pos-1-2		pos-1-3		pos-1-4		pos-1-5									
Pedestrian XY position	pos-0-0		pos-1-0		pos-2-0		pos-3-0		pos-4-0		pos-5-0		pos-6-0		pos-7-0		pos-0-1		pos-1-1		pos-2-1		pos-3-1		pos-4-1		pos-5-1		pos-6-1		pos-7-1	
Pedestrian pose	pose01					pose02					pose03					pose04					pose05											
Pedestrian asset	A1			A2			A3			A4			A5			A6			A7			A8			A9			A10				
Pedestrian hip direction	d0			d45			d90			d135			d180			d225			d270			d315										
Parked vehicle 1 type	BMW1					BMW2					BMW71					VW ID.3					VW Golf 8					VW Atlas						
Parked vehicle 1 XY position	pos-0-0			pos-0-1			pos-0-2			pos-1-0			pos-1-1			pos-1-2			pos-2-0			pos-2-1			pos-2-2							
Parked vehicle 1 color	BMW Black	BMW Cerium grey	BMW Melbourne red	BMW Mineral grey	BMW Misano blue	BMW Sao Paolo yellow	BMW Snapper Rocks blue	BMW Sunset orange	BMW White	VW Gletscher Weiss	VW Mangangrau	VW Mekana Turquoise	VW Mondsteingrau	VW Scale Silver	VW Stonewashed Blue	VW Energetic Orange	VW Deep Black	VW Delfingrau	VW Kings Red													
Parked vehicle 2 type	BMW1					BMW2					BMW71					VW ID.3					VW Golf 8					VW Atlas						
Parked vehicle 2 color	BMW Black	BMW Cerium grey	BMW Melbourne red	BMW Mineral grey	BMW Misano blue	BMW Sao Paolo yellow	BMW Snapper Rocks blue	BMW Sunset orange	BMW White	VW Gletscher Weiss	VW Mangangrau	VW Mekana Turquoise	VW Mondsteingrau	VW Scale Silver	VW Stonewashed Blue	VW Energetic Orange	VW Deep Black	VW Delfingrau	VW Kings Red													
Illumenation	direct sun										diffuse light																					
Sun direction	d0			d45			d90			d135			d180			d225			d270			d315										
Sun elevation	low					medium										day																
Road surface	A					B					C					D																

Source: Robert Bosch GmbH

- **Discretization:** The most critical dimensions are identified and discretized
- **Test coverage:** With pairwise testing it's possible to achieve a high error coverage in traditional software testing

Data production - Example data snapshot 1

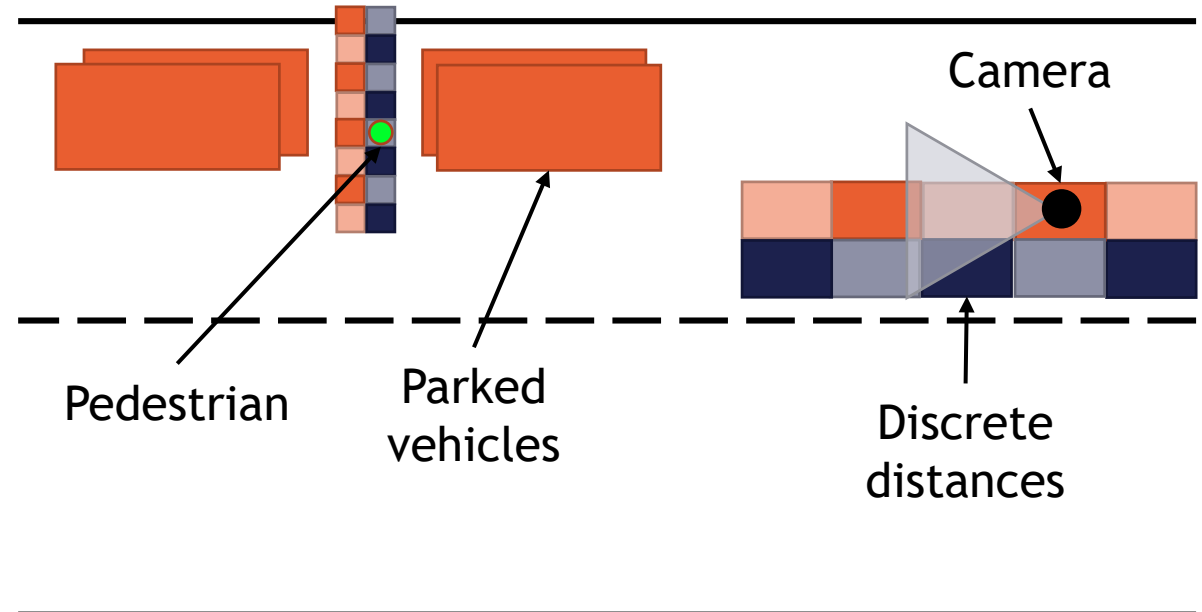


- **Safety critical:** Pedestrian has a running pose towards the camera
- The perception function shall be able to detect the pedestrian early enough without any image perturbations



- Those images are well suited as a reference for the analysis of brittleness in DNN's

Data production - Example data snapshot 1



- **Safety critical:** The legs are extended to the driving lane
 - **Uncommon pose:** Pedestrian lays between two vehicles and is difficult to see
- ➔
- In which combinations is the object detector **not** capable to perceive the pedestrian?

Examples for data post processing



original



fog



frost



brightness



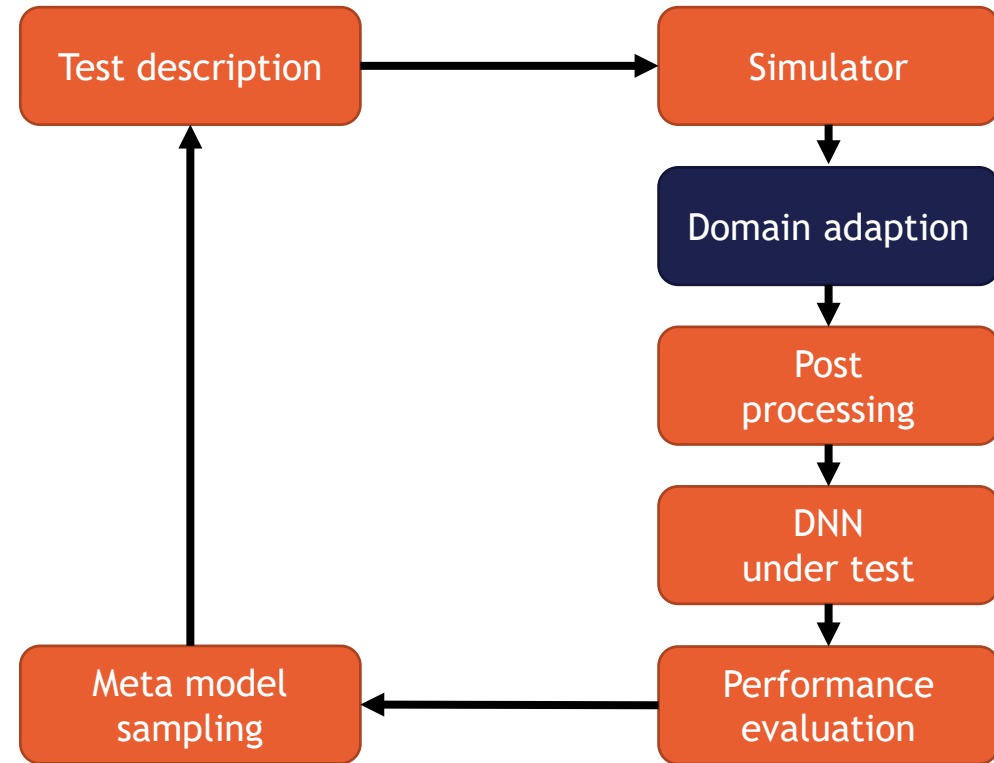
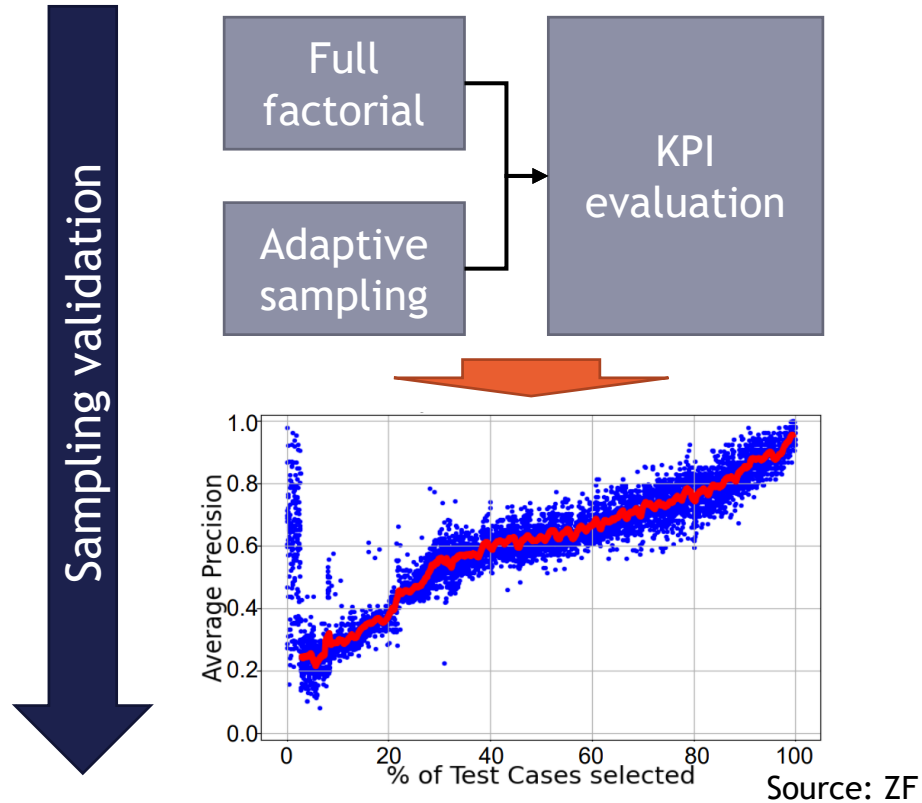
contrast



Motion blur



Test space exploration optimization



→ The most performance critical test cases are identified early in the test exploration
“Adaptive test case selection for DNN-based perception functions”
Paper release: <https://ieeexplore.ieee.org/document/9582499>

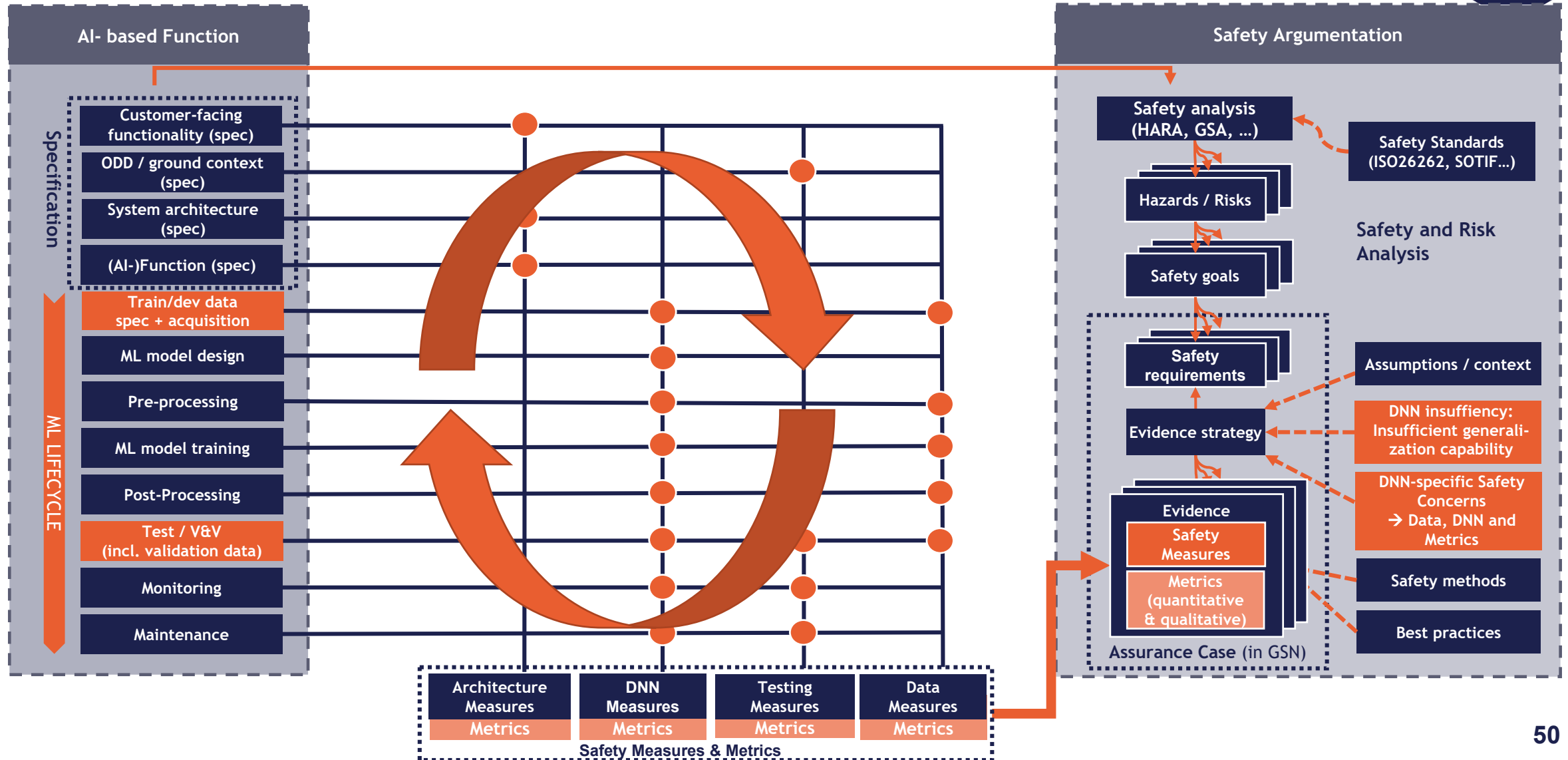


5

How do we work



Our Approach: Summary



Our Approach: Evidence Workstreams



Empowering experts from safety engineering and ML to produce measures and evidences

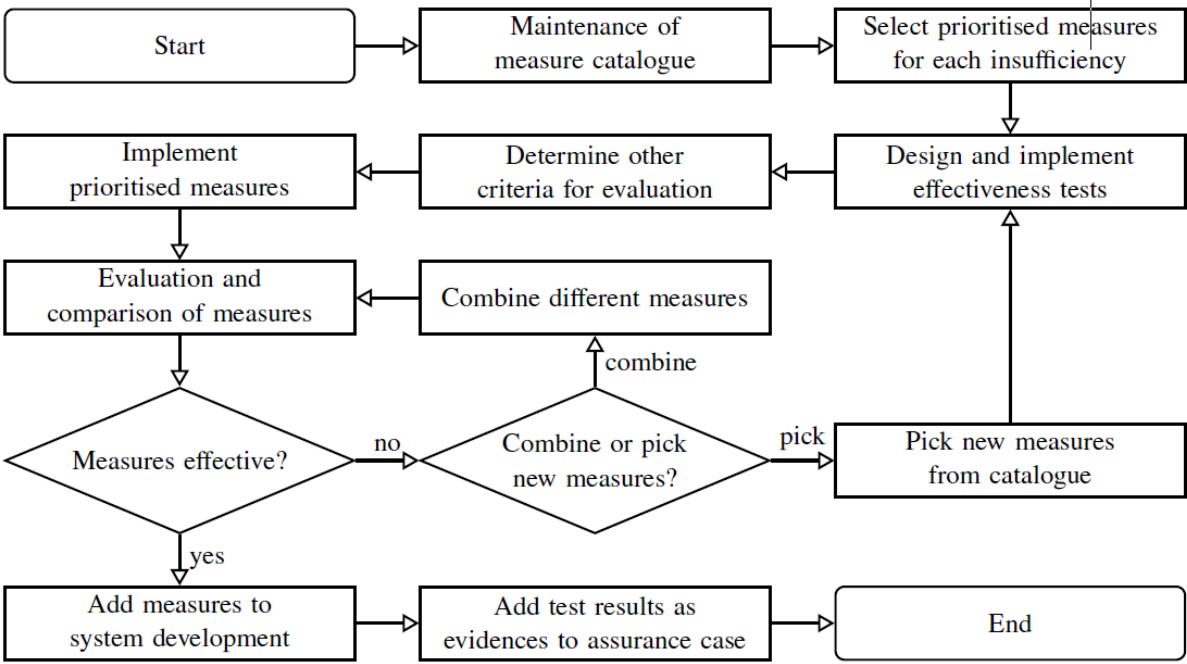
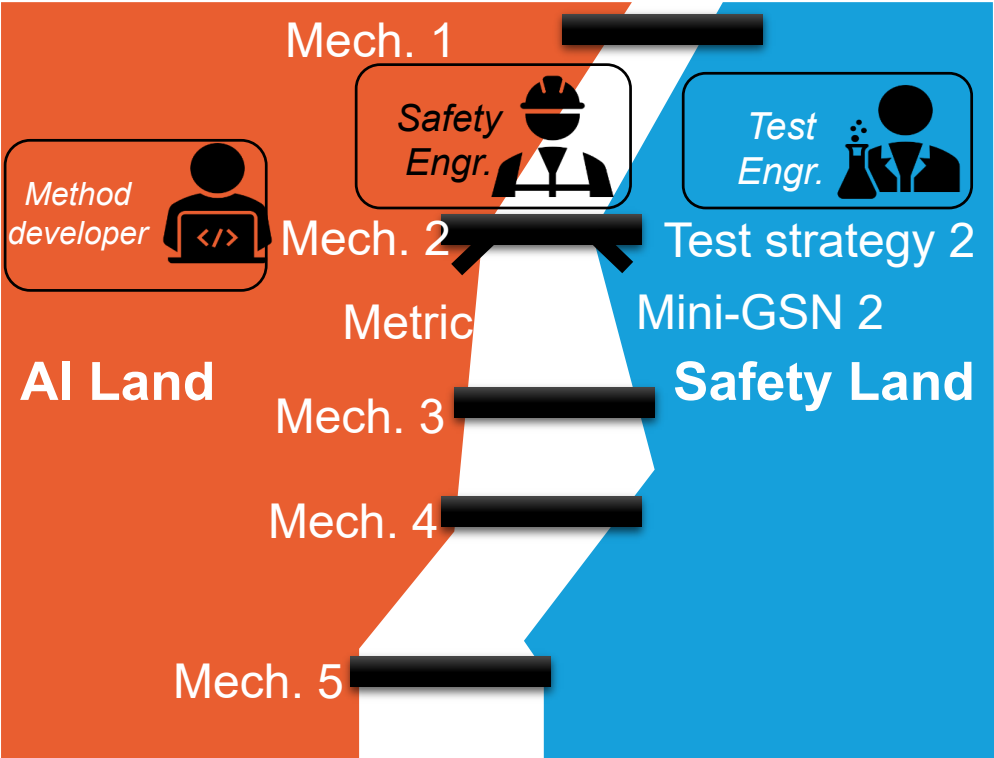


Fig. 3. Flowchart of the evidence workflow process

Agenda

1. Introduction – CARIAD
2. DNNs and Safety in Automated Driving
3. KI-Absicherung Project & Approach
4. Summary

Summary

Findings & Consequences

- Safe AI is a central challenge for highly automated driving
- KI-Absicherung provides an approach for Safe AI
- Approach may serve as template for the industry and beyond (see ISO PAS 8800)
- Deep integration of AI-specifics into development PMT is necessary (continuous assurance of AI)

Contact:

Fabian Hüger

Artificial Intelligence Safety
@Volkswagen CARIAD

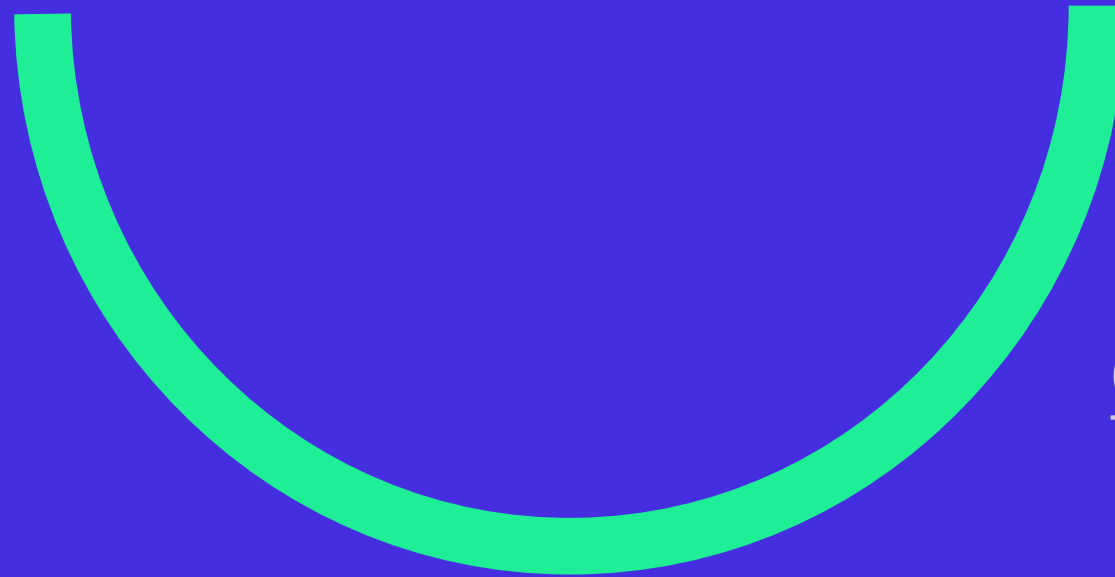
Contact: fabian.hueger@volkswagen.de



<https://scholar.google.de/citations?user=ISPOitUAAAAJ>

www.ki-absicherung-projekt.de  [@KI_Familie](https://twitter.com/KI_Familie)  [KI Familie](https://www.linkedin.com/company/ki-familie)

Thank you!



QUESTIONS?