

Rudolf Hoffmann / 13.03.2025





# Robustheitsanalyse für Vision-Modelle HTW Berlin

Rudolf Hoffmann / 13.03.2025



## Agenda

- 1 Projektüberblick
- **2** Robustheit von Vision Modellen
- **3** Deepbench: Vertrauen durch Robustheit
- **4** Use Cases und Corruptions
- 5 Domäne-Spezifische Analyse
- **6** Evaluation





#### **Arbeitspaket A7**

 Erstellung eines Frameworks zur Risiko- und Robustheitsanalyse von KI-Modellen

#### **Ziel des Frameworks**

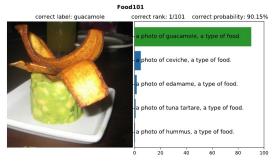
 Nutzerzentrierte Evaluierung existierender KI-Modelle zur Identifikation des optimalen Ansatzes für spezifische Use-Cases

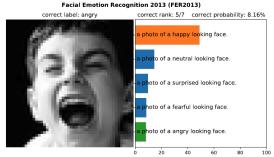
#### Relevanz für KI-Sicherheit

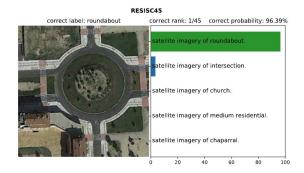
- Bewertung von KI-Modellen gegen Manipulation und Fehlfunktionen -Verlässlichkeit
- Minimierung von Sicherheitsrisiken durch domänenspezifische Evaluierung



#### **Robustheit von Vision Modellen**

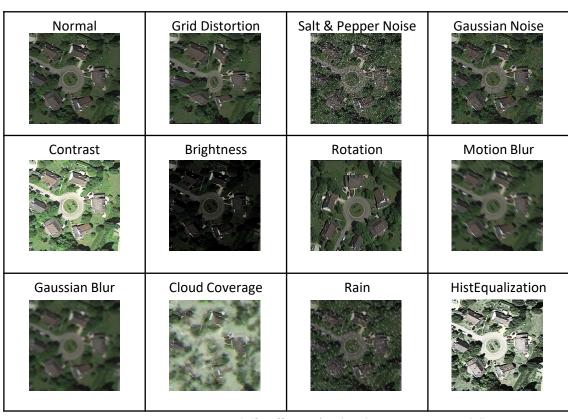




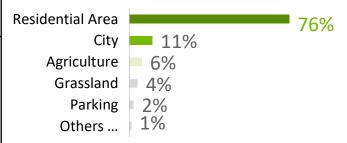




### **Robustheit von Vision Modellen**



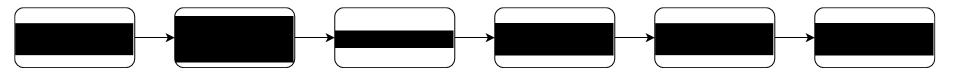
#### Beispielwerte





### Deepbench: Vertrauen durch Robustheit

- Flexibles und skalierbares Framework zur Robustheitsanalyse von KI-Modellen
- Fokus auf Vision und Vision-Language Modelle
- Use-Case spezifische Analyse





#### **Use Cases**



**Autonomes Fahren** 



**Manufacturing Quality** 



MedicalDiagnosis



**Handheld Device** 



**Security (Emotions)** 



Satellite



# **Corruptions**

#### Blur





















Noise



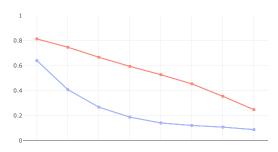
# **Domäne-Spezifische Analyse**

	Medical Diagno- sis	Auto- nomous Vehicles	Manu- factur- ing	People Recog- nitio	Aerial Imaging	Hand- held
Brightness	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>
Cloud Generator					<b>√</b>	
Contrast	✓	✓	✓	<b>√</b>	<b>√</b>	✓
Gaussian Blur	✓		✓			
Gaussian Noise	✓	✓	✓	✓	<b>√</b>	✓
Colour Shift		✓	✓			✓
Grid Distortion						
Grid Deformation	✓		✓	<b>√</b>	<b>√</b>	
Hist Equalization	✓			✓		
Horizontal Flip	✓		✓	✓	<b>√</b>	✓
Vertical Flip					<b>√</b>	
Rotation	✓	✓	✓	✓	✓	✓
Motion Blur		✓				✓
Perspective		✓	✓	<b>√</b>	<b>√</b>	✓
Transformation						
Rain		✓				
Salt and Pepper		✓		✓		✓
Shadow		✓	✓	<b>√</b>	<b>√</b>	✓

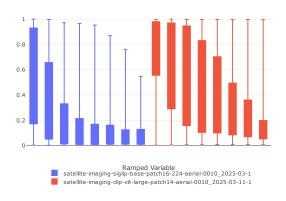


### **Evaluation: Beispiel Gaussian Noise**

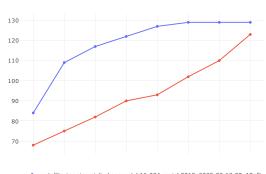




Boxplot of Ground Truth Score for GaussianNoise



Label Flips for GaussianNoise



satellite-imaging-siglip-base-patch16-224-aerial-0010\_2025-03-12-09\_18\_5
 satellite-imaging-clip-vit-large-patch14-aerial-0010\_2025-03-11-19\_02\_25 (L



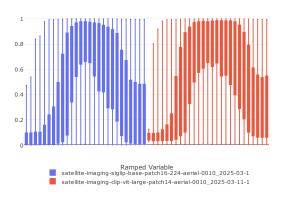
Gaussian Noise



### **Evaluation: Beispiel Brightness**



Boxplot of Ground Truth Score for Brightness







satellite-imaging-siglip-base-patch16-224-aerial-0010\_2025-03-12-09\_18\_5!
satellite-imaging-clip-vit-large-patch14-aerial-0010\_2025-03-11-19\_02\_25 (L



հեա

### **Abschluss**

