# KI ABSICHERUNG
### Safe AI for Automated Driving

# Visual Exploration and Semantic Analysis of DNN Weaknesses with ScrutinAI

## Elena Haedecke, Michael Mock

### Investigating DNN predictions

The performance of a DNN usually depends on a combination of different dimensions and cannot be explained by looking at a single dimension of the input. Rather, several dimensions have to be investigated as potential influencing factors. In order to efficiently arrive at relevant insights within the large amount of data and dimensions, tool support is needed that guides the human in terms of visual analytics (VA). The goal is to allow the analyst and/or the auditor to use the insights gained during the analysis process to find evidence for the origins or the absence of specific insufficiencies, ultimately contributing towards an overall safety argumentation. Thus, the workflow of the tool focuses on efficiently guiding the human analyst on utilizing domain knowledge and concentrating on semantic aspects of the model and data (see Figure 1 (a)).
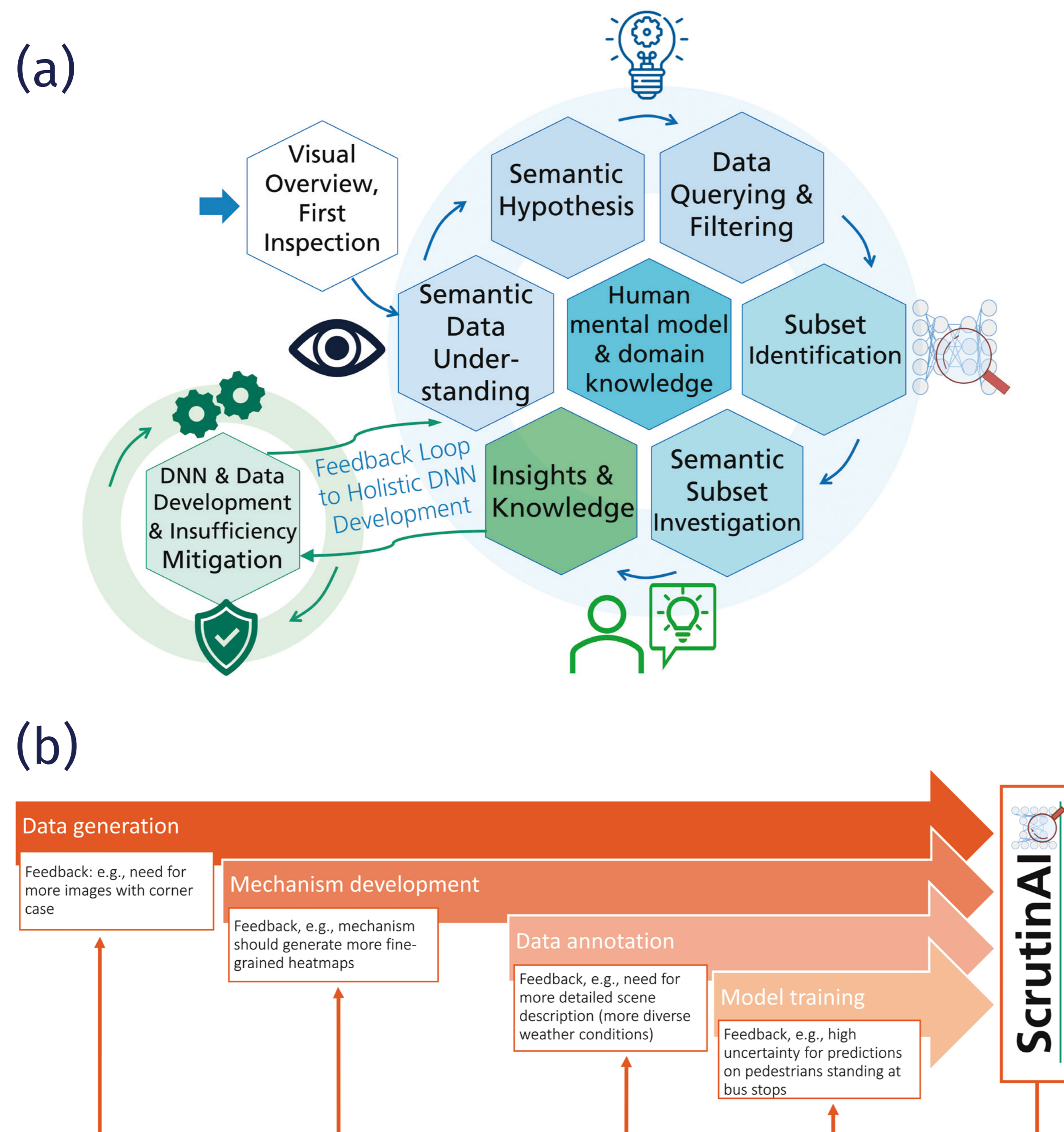
### Safety Hypothesis:

The method addresses the safety concern Incomprehensible behavior. It enables the analyst and/or the auditor to identify presence and causes of incorrect predictions by visually exploring the DNN decisions from different perspectives with the interactive tool ScrutinAI (see Figure 2).

### Deriving local and global explanations

To identify noticeable patterns and systematic weaknesses of the DNN and to investigate their causes, common performance metrics are linked to human-understandable semantic concepts defined by domain experts. This helps to establish semantic hypotheses, which are the basis for further, more in-depth investigation of subsets of interest, that potentially represent weaknesses of the model. By examining individual predictions, local explanations can be derived, and by considering larger subsets of data, inferences can be made about the global behavior of the model. During this investigation, the analyst either finds sufficient and strong evidence to verify or falsify the hypothesis, or proceeds with deeper analysis. Insights and knowledge gained during the investigation are fed back to the other stakeholders to mitigate the shortcomings of the DNN (see Figure 1 (b)).

### References:

Haedecke et al.: ScrutinAI: A Visual Analytics Approach for the Semantic Analysis of Deep Neural Network Predictions, 2022. EuroVis Workshop on Visual Analytics (2022)

(a)



(b)



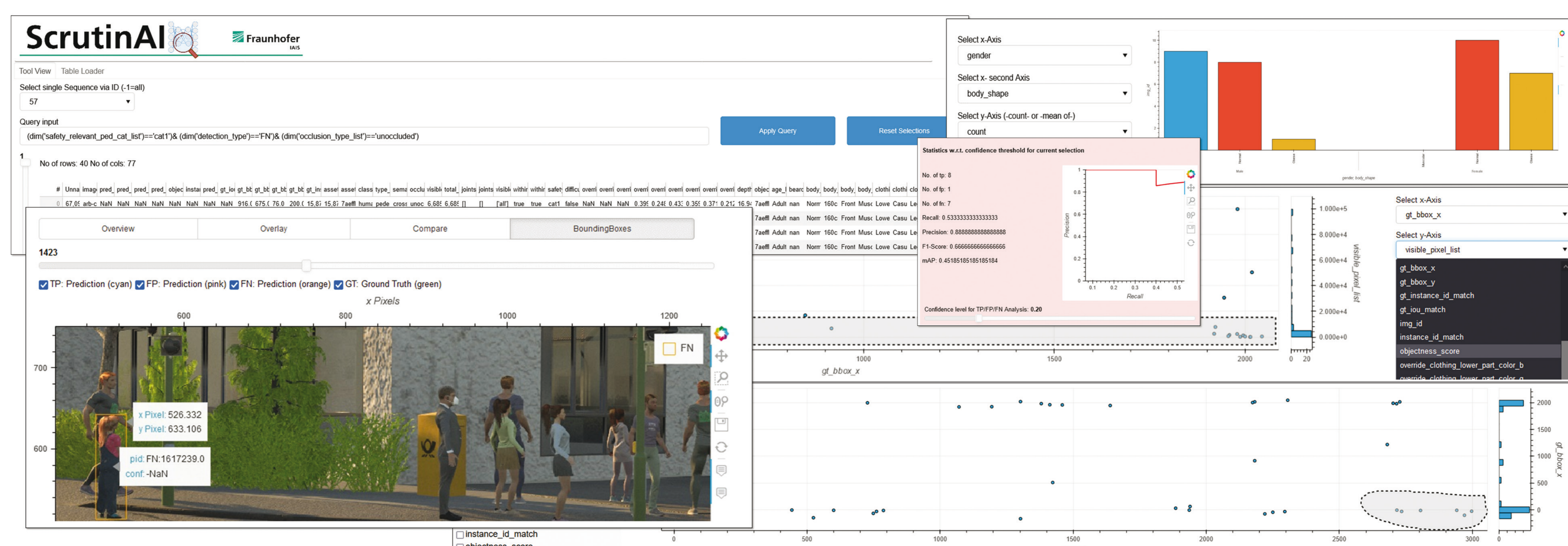*Figure 1: Workflow (a) and feedback loop (b).*



*Figure 2: ScrutinAI visually represents model inputs, outputs, and metadata through linked, interactive elements. The dynamic adjustment of parameters, textual queries or the graphical selection of interesting data subsets allow an investigation of the influence and interdependence of attributes.*

# Fraunhofer
## IAIS

**For more information contact:**
**elena.haedecke@iais.fraunhofer.de**
**michael.mock@iais.fraunhofer.de**

**https://www.ki-absicherung-projekt.de**       🐦 @KI_Familie       in KI Familie

**KI FAMILIE**

Image Sources:
Figure 1: Fraunhofer IAIS; Figure 2: Fraunhofer IAIS / MackeVision