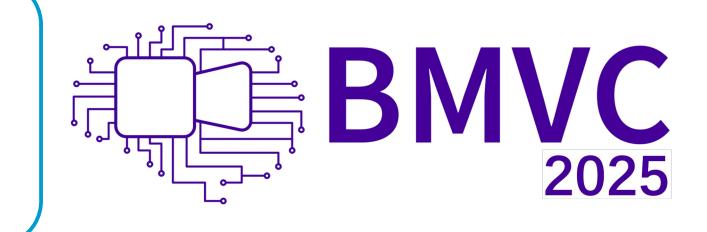


# Emre Gülsoylu<sup>1</sup>, André Kelm<sup>1</sup>, Lennart Bengtson<sup>2</sup>, Matthias Hirsch<sup>1</sup>, Christian Wilms<sup>1</sup>, Tim Rolff<sup>3</sup>, Janick Edinger<sup>2</sup> and Simone Frintrop<sup>1</sup>

<sup>1</sup> University of Hamburg, CV Group <sup>2</sup> University of Hamburg, DOS Group <sup>3</sup> University of Hamburg, HCI Group



# TRUDI and TITUS: A Multi-Perspective Dataset and A Three-Stage Recognition System for Transportation Unit Identification

## Motivation

**Task:** Detect and recognise ID codes on transportation units (TU), such as containers and trailers.

#### **Problem:**

- Lack of publicly available datasets for TU identification and benchmarking.
- Reliance on fixed camera settings and dynamic, real-world conditions.

## Solution:

- TRansportation Unit Detection and Identification (TRUDI) dataset contains aerial and ground imagery from multiple terminals, taken under various weather conditions.
- Three-stage Identification of Transportation UnitS (TITUS) pipeline was developed for robust, multi-perspective TU identification.

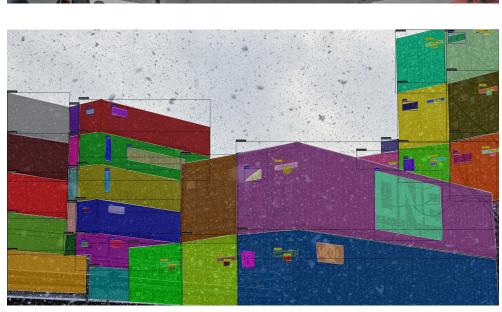
## TRUDI Dataset

35,034 annotated instances (containers, trailers, tank containers, ID text, logos) from ground and aerial perspectives.











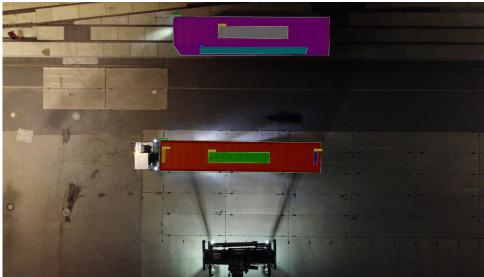


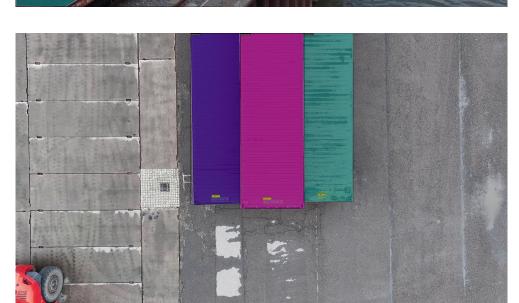












- Suitable for various CV tasks such as object detection, instance segmentation, logo detection, text detection, text recognition and text spotting.
- Varying lighting, weather, perspectives (ground/aerial), and image quality.
- 18-month collection across multiple countries, times of day, and seasons.

Category	# Instances
Container	11109
Tank Container	808
Trailer	2780
ID Text	14009
Logo	6328

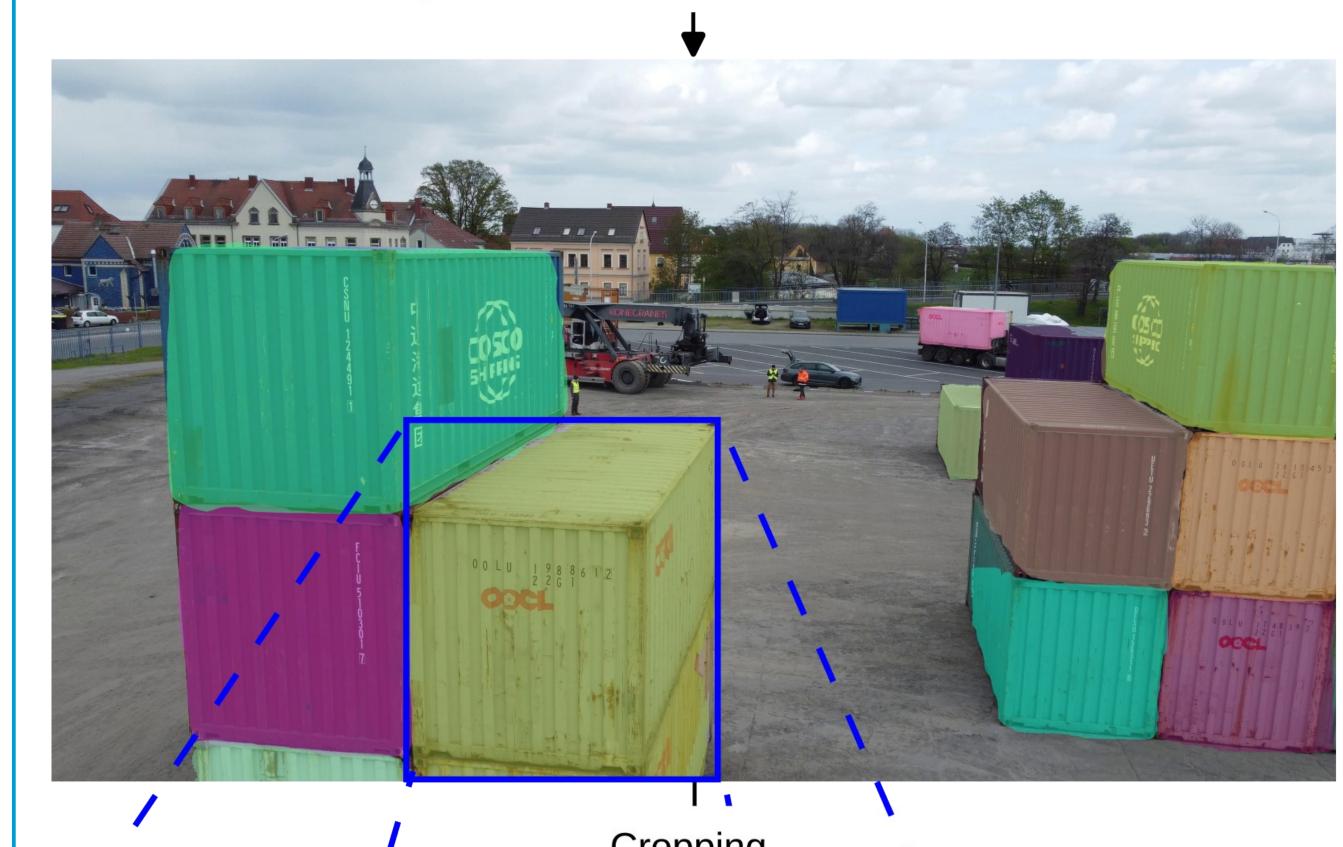
# Evaluation

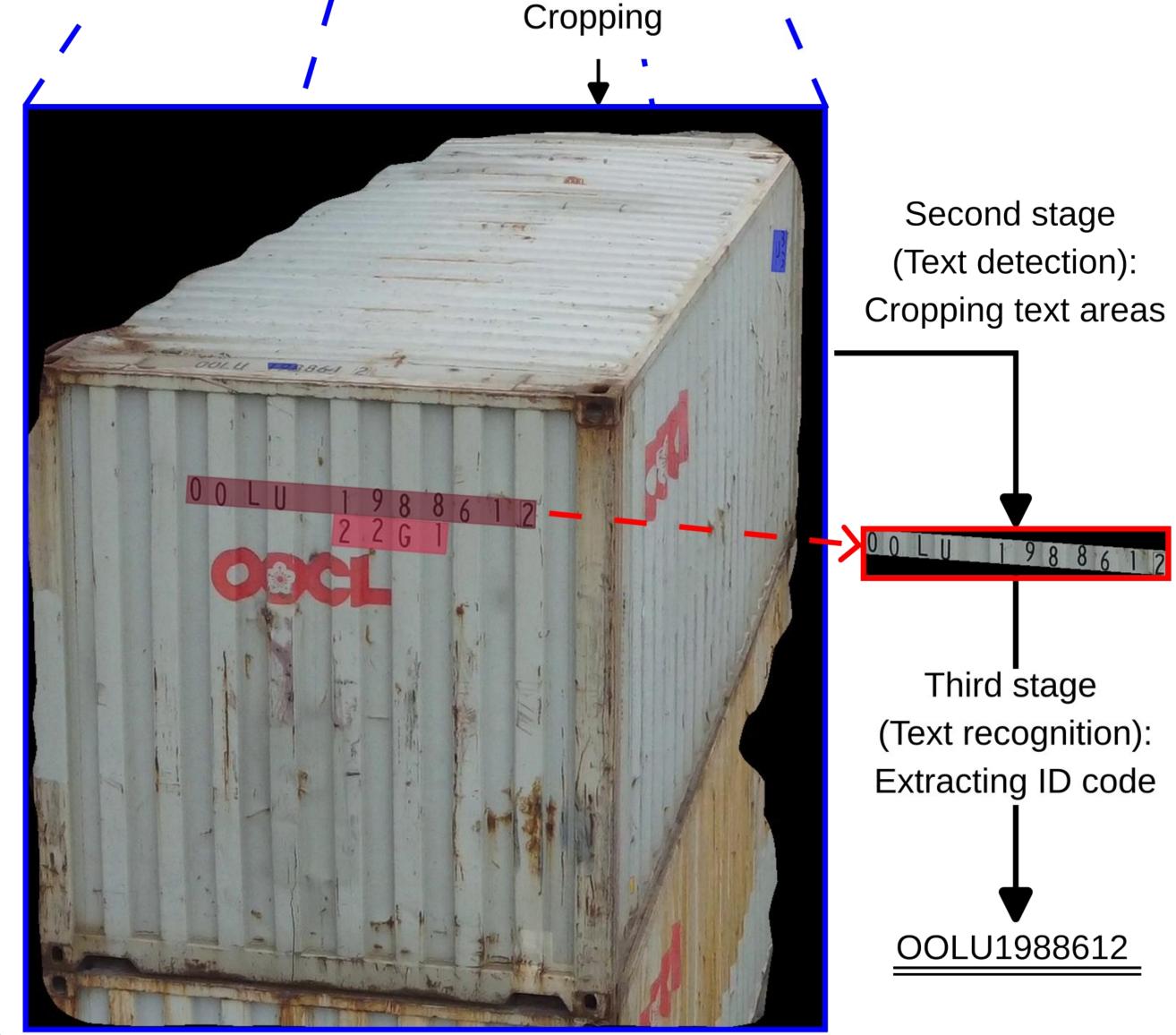
Stage	Perspective	Model	AP@0.50:0.95	AP@0.50	AR@0.50:0.95
Instance Seg.	Aerial Ground Combined	Mask R-CNN-Swin(1C)[1,2] Mask R-CNN-Swin(3C)-C[1,2] Mask R-CNN-Swin(3C)[1,2]	0.43 <b>0.44</b> <b>0.44</b>	0.62 <b>0.67</b> 0.65	0.48 <b>0.50</b> <b>0.50</b>
Stage	Perspective	Model	Recall	Precision	Hmean
Text Det.	Aerial Ground Combined	DBNet++-C[3] DBNet++-C[3] DBNet++[3]	0.79 0.79 0.79	0.68 <b>0.69</b> <b>0.69</b>	0.73 <b>0.74</b> <b>0.74</b>
Stage	Perspective	Model	Word acc.	Char. recall	Char. precision
Text Rec.	Aerial Ground Combined	RobustScanner-C[4] RobustScanner-C[4] RobustScanner[4]	<b>0.68</b> 0.54 0.63	<b>0.88</b> 0.73 0.83	<b>0.88</b> 0.73 0.84
Stage	Perspective	Precision	Recall	F1 Score	Accuracy
End-to-end	Aerial Ground Combined	<b>0.45</b> 0.25 0.39	<b>0.30</b> 0.19 0.27	<b>0.36</b> 0.22 0.32	<b>0.22</b> 0.12 0.19

# TITUS Pipeline



First stage (Instance segmentation): Segmentation masks for TU instances





### Paper + Dataset + Weights



### References

[1] Kaiming He, et al. Mask R-CNN. In Proc. of the IEEE International Conference on Computer Vision, 2961–2969, 2017.
[2] Ze Liu, et al. Swin transformer: Hierarchical vision transformer using shifted windows. In Proc. of the IEEE/CVF International Conference on Computer Vision, 10012–10022, 2021.
[3] Minghui Liao, et al. Real-time scene text detection with differentiable binarization and adaptive scale fusion. IEEE Trans. on Pattern Analysis and Machine Intelligence, 45(1):919–931, 2022.
[4] Xiaoyu Yue, et al. RobustScanner: Dynamically enhancing positional clues for robust text recognition. In Proc. European Conference on Computer Vision, 135–151. Springer, 2020.



