

# ***DENSO***

Crafting the Core

## **Are Transformers More Robust?**

## **Towards Exact Robustness**

## **Verification for Transformers**

SafeComp, 20.09.2023

### **Brian Hsuan-Cheng Liao**

Systems Engineering R&D

Corporate R&D

DENSO AUTOMOTIVE Deutschland GmbH (DNDE)

In collaboration with

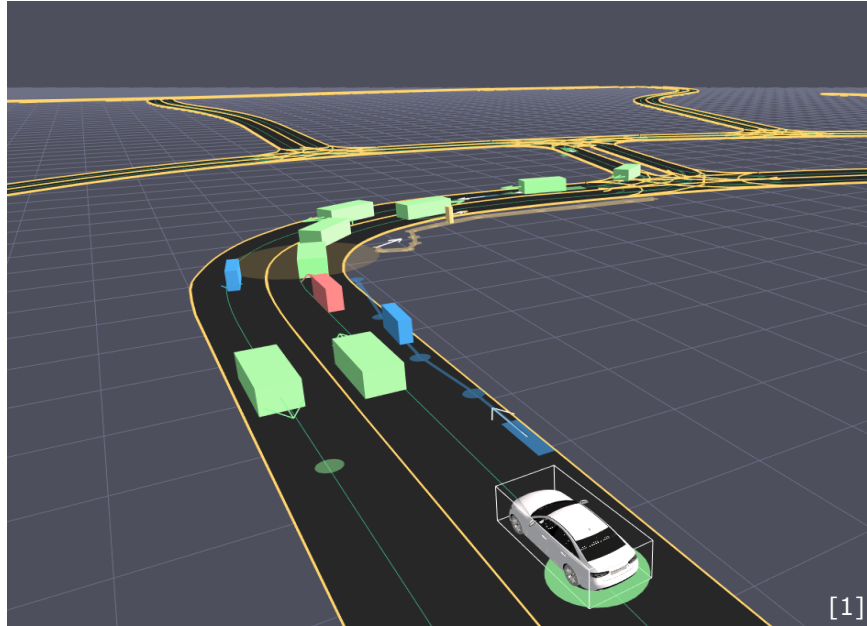
Dr. Chih-Hong Cheng (Fraunhofer IKS)

Dr. Hasan Esen (DNDE)

Prof. Alois Knoll (Technical University of Munich)



# Neural Networks Are Everywhere ...



Autonomous Driving  
Aircraft Autopiloting



Medical Diagnosis  
Surgical Robots

[1] Source: [Bloomberg](#).

[2] Source: [CeramTec](#).

# Especially Transformers ...

**Attention Is All You Need**

Ashish Vaswani  
Google Brain  
avaswani@google.com

Llion Jones  
Google Research  
llion@google.com

## AN IMAGE IS WORTH 16x16 WORDS: TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

Alexey Dosovitskiy<sup>1,†</sup>, Lucas Beyer<sup>\*</sup>, Alexander Kolesnikov<sup>\*</sup>, Dirk Weissenborn<sup>\*</sup>,

Xiaohua Zhang<sup>1</sup>,  
Georgios

## End-to-End Object Detection with Transformers

Nicolas Carion<sup>1,2</sup>[0000-0002-2308-9680], Francisco Massa<sup>2</sup>[000-0003-0697-6664],  
Gabriel Synnaeve<sup>2</sup>[0000-0003-1715-3356], Nicolas Usunier<sup>2</sup>[0000-0002-9324-1457],  
Alexander Kirillov<sup>2</sup>[0000-0003-3169-3199], and Sergey  
Zagoruyko<sup>2</sup>[0000-0001-9684-5240]

<sup>1</sup> Paris Dauphine University  
<sup>2</sup> Facebook AI

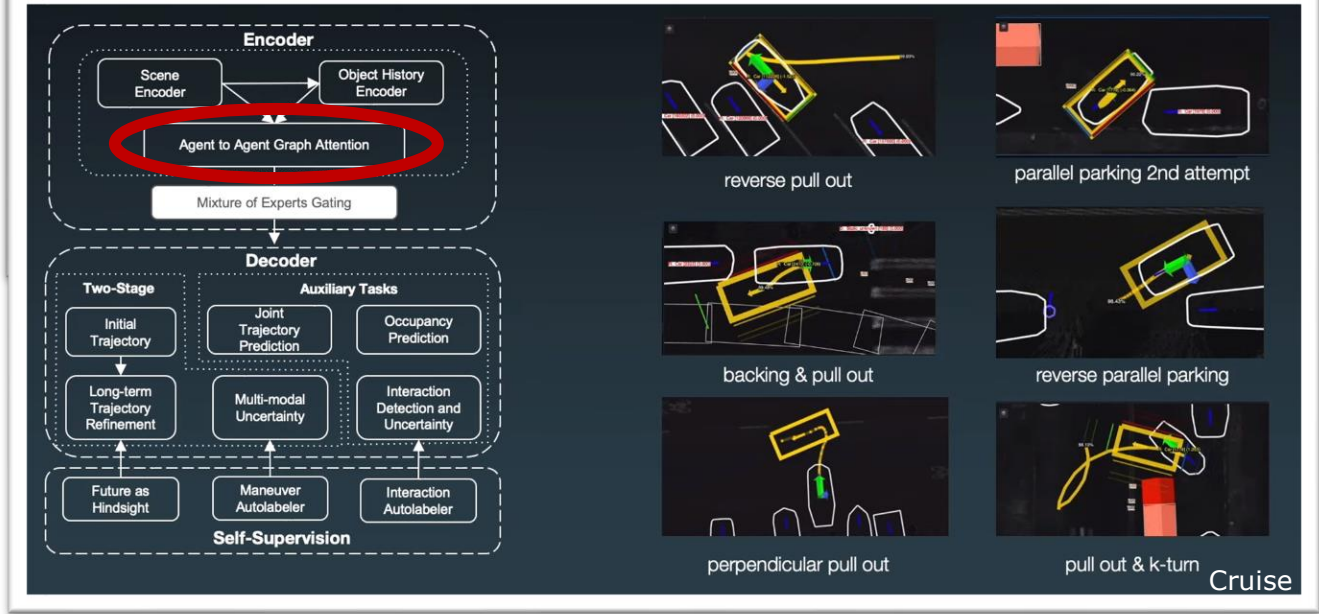
{alcinos, fmassa, gab, usunier, akirillov, szagoruyko}@fb.com

**Abstract.** We present a new method that views object detection as a direct set prediction problem. Our approach streamlines the detection pipeline, effectively removing the need for many hand-designed components like a non-maximum suppression procedure or anchor generation that explicitly encode our prior knowledge about the task. The main ingredients of the new framework, called DETR (DEtection TRansformer) or DETR, are a set-based global loss that forces unique predictions via bipartite matching, and a transformer encoder-decoder architecture. Given a fixed small set of learned object queries, DETR reasons about the relations of the objects and the global image context to directly output the final set of predictions in parallel. The new model is conceptually simple and does not require a specialized library, unlike many other modern detectors. DETR demonstrates accuracy and run-time performance on par with the well-established and highly-optimized Faster R-CNN baseline on the challenging COCO object detection dataset. Moreover, DETR can be easily generalized to produce panoptic segmentation in a unified manner. We show that it significantly outperforms competitive baselines. Training code and pretrained models are available at <https://github.com/facebookresearch/detr>.

While the Transformer language processing, attention is used to replace overall structure and a pure transformer very well on image data and transfer (ImageNet, COCO) results compared to traditional models.



[1]



[2]

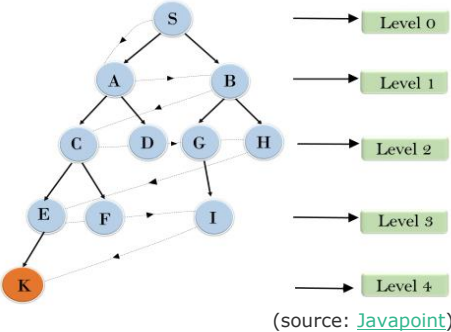
[1] Source: [Tesla AI Day 2021](#)

[2] Source: [Cruise Under the Hood 2021](#)

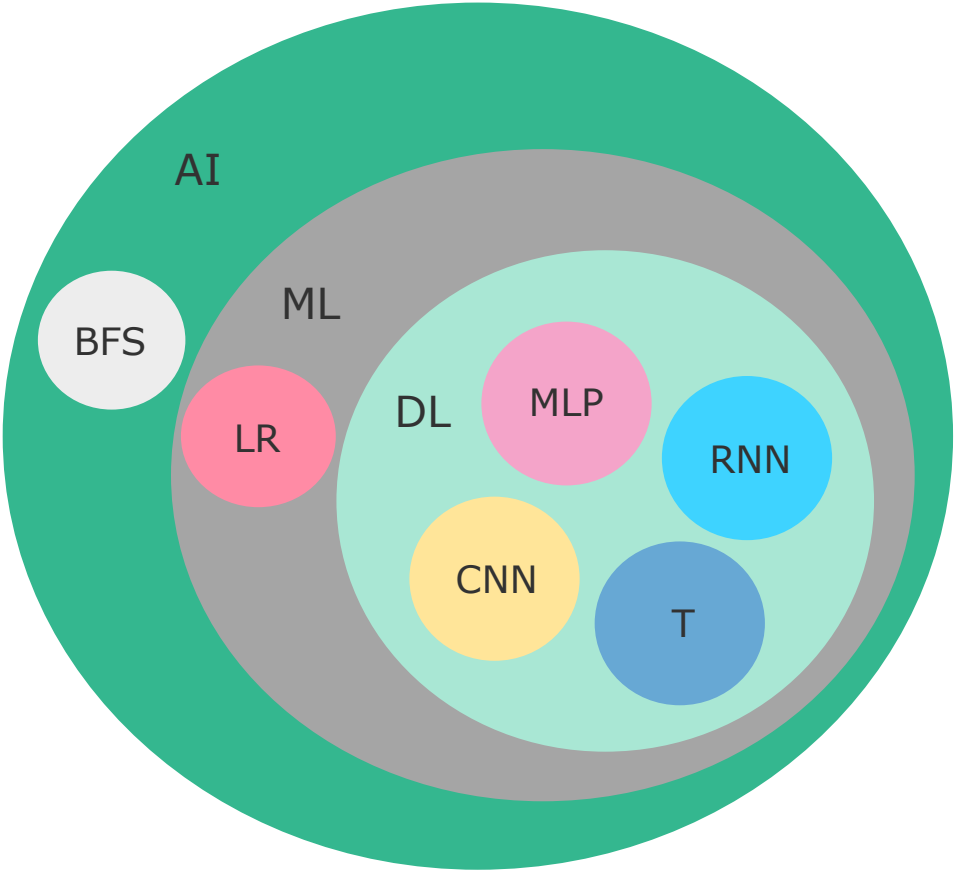
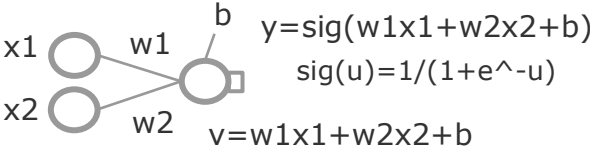
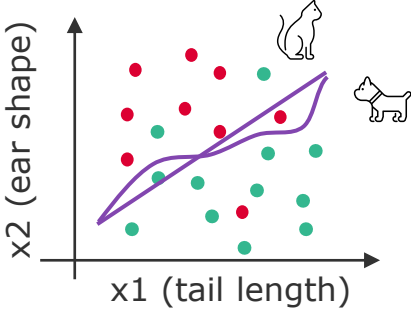


# But Why? A Look into these "AI" Models

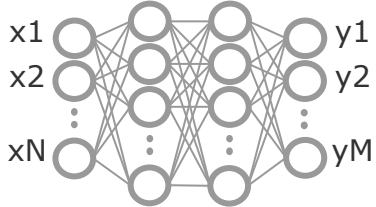
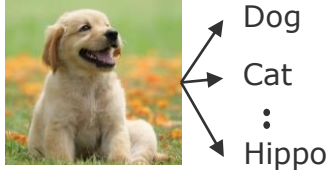
Breadth-first search



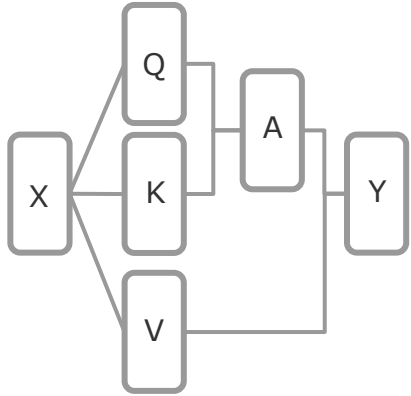
Logistic regression



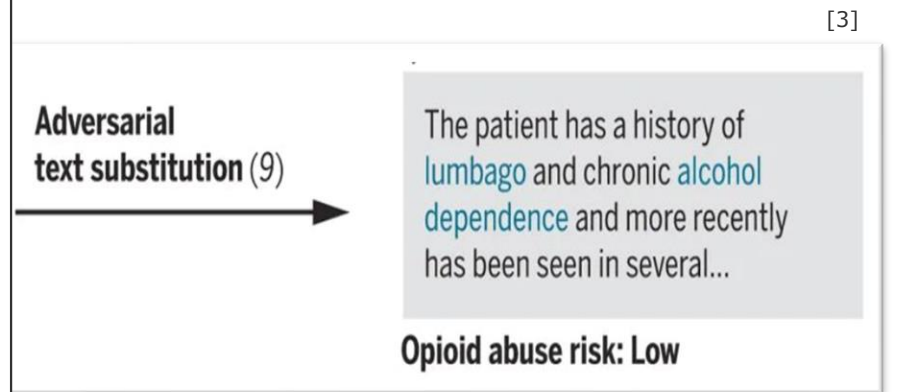
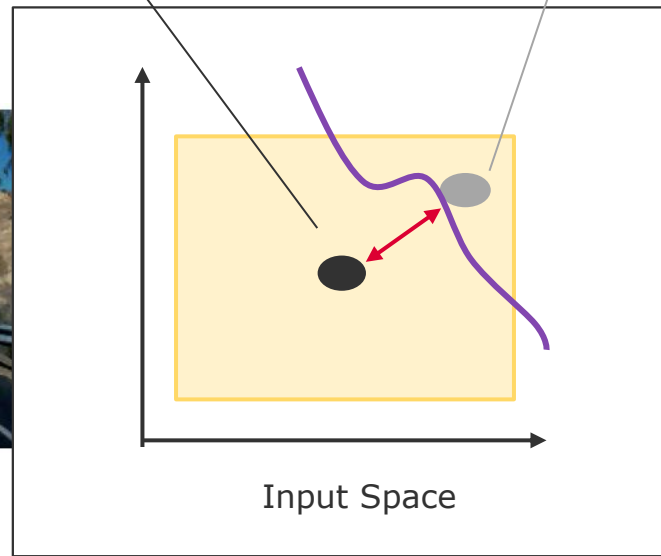
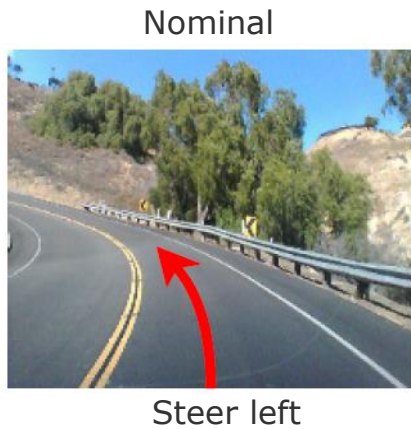
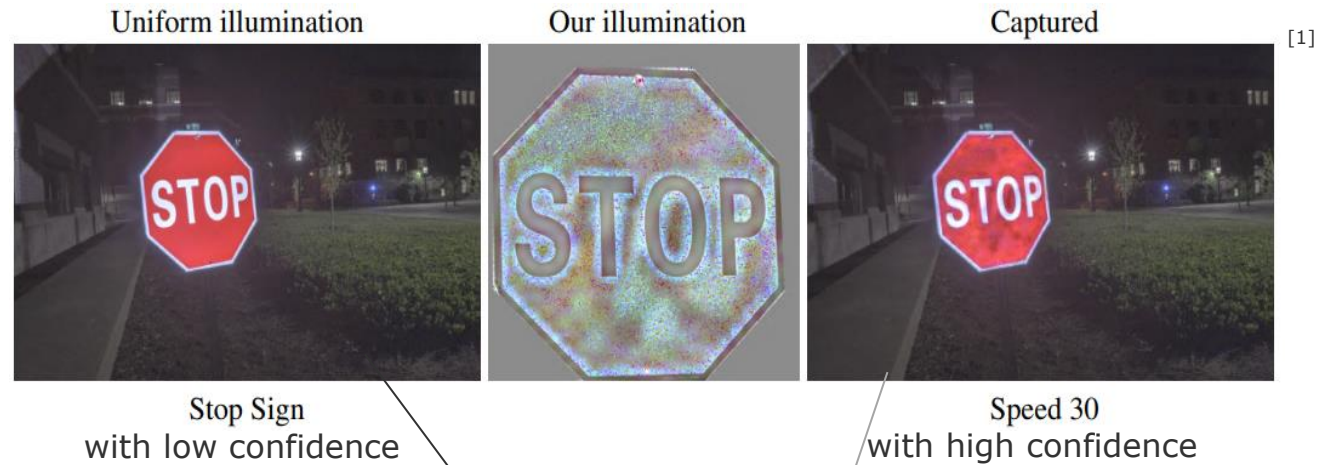
Multi-layer perceptron



Transformer (self-attention)

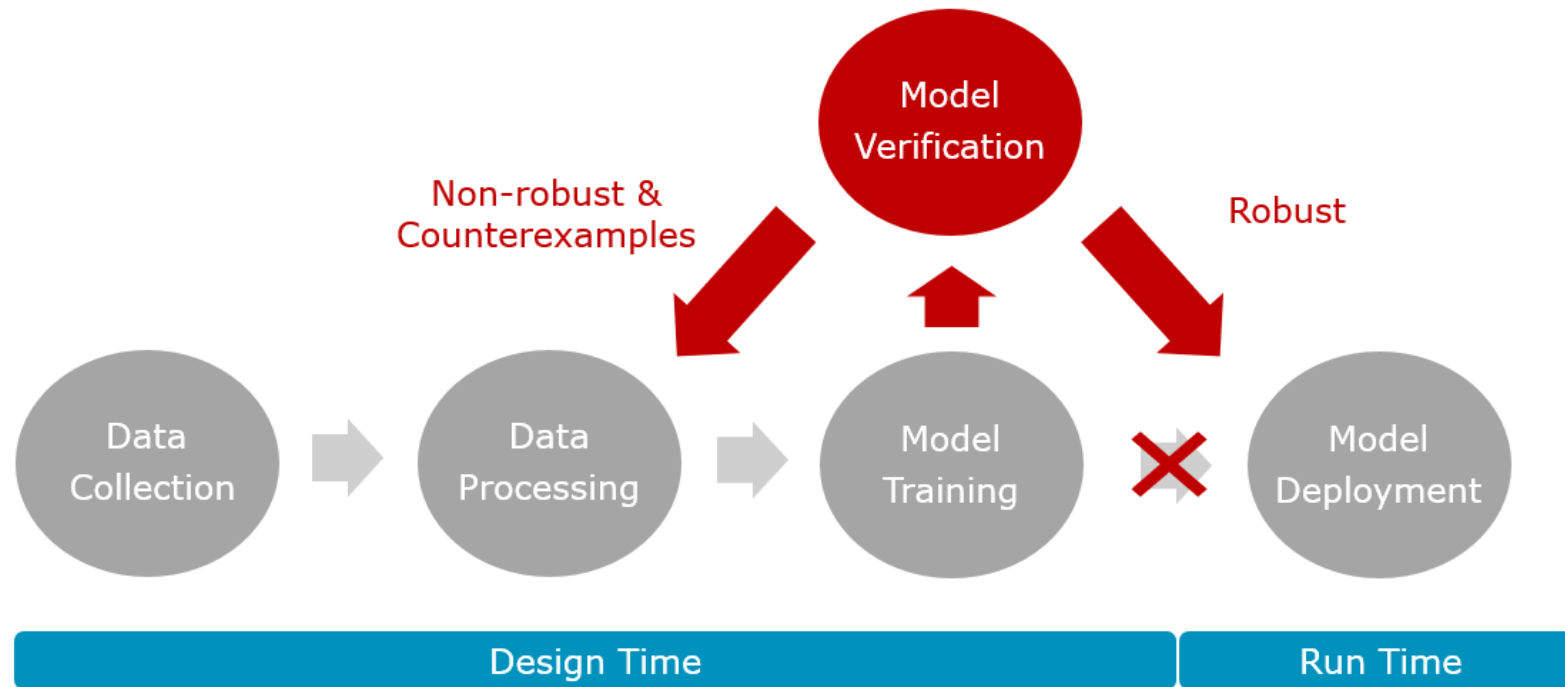


# The Robustness Problem of NNs

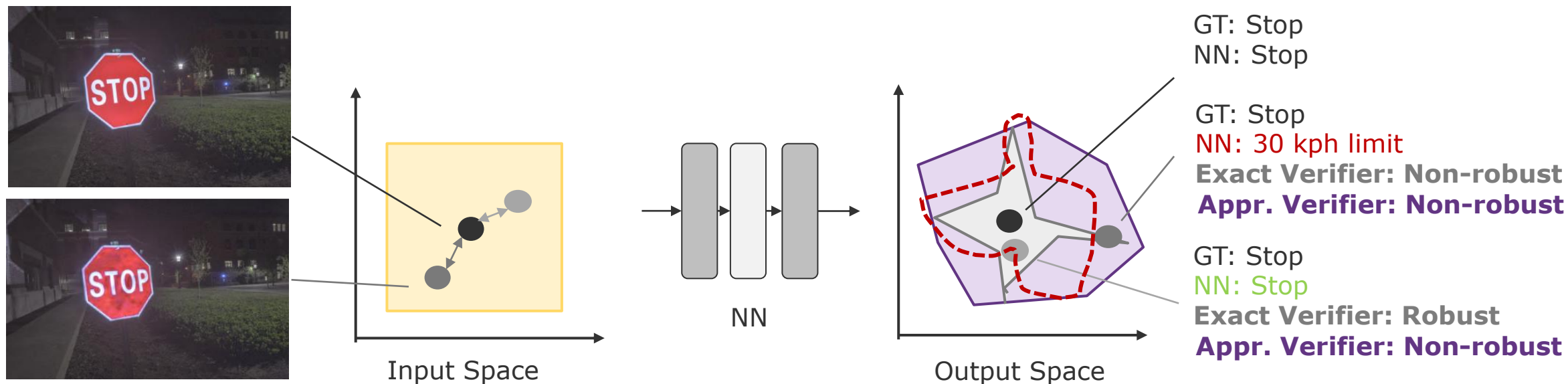


[1] "Optical adversarial attacks," Gnanasambandam *et al.*, ICCV, 2021.  
 [2] "DeepXplore: Automated whitebox testing of deep learning systems," Pei *et al.*, SOSP, 2017.  
 [3] Cary *et al.*, "Adversarial attacks on medical machine learning," in Science, 2020.

# Robustness Verification as a Potential Remedy

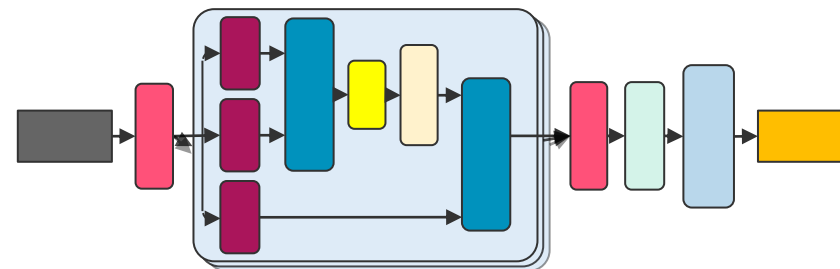
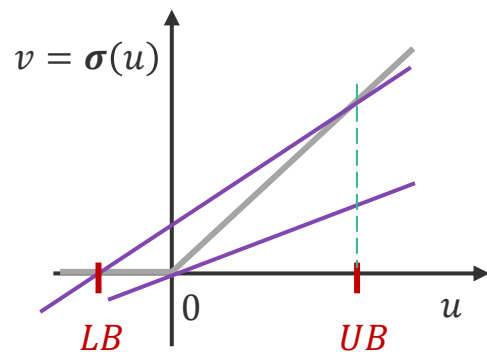


# Main Approaches and Common Challenges

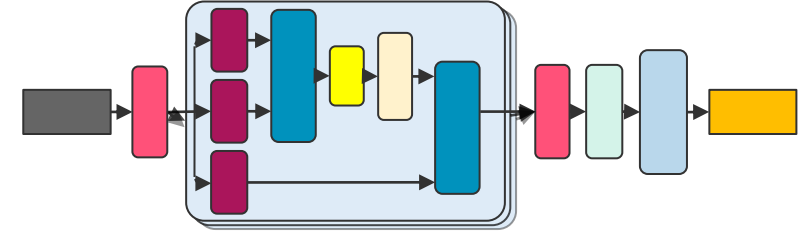
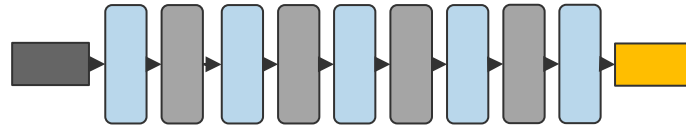


■ Admissible perturbation region   
 ⬠ Exact adversarial polytope   
 - - - NN decision boundary   
 ⬠ Over-approximated adversarial polytope

- Challenge 1: The non-convex activations, e.g., ReLU
- Challenge 2: Emerging architectures and operations



# State of the Art in NN Robustness Verification



	MLPs (or CNN/RNN variants)	Transformers
Approximate Verifier	$\sim 10^6$ neurons <sup>[1]</sup> Image Classification (CIFAR-10)	$\sim 10^4$ neurons <sup>[3]</sup> Sentiment Analysis (Yelp)
Exact Verifier	$\sim 10^5$ neurons <sup>[2]</sup> Image Classification (CIFAR-10)	Lacking, but a network of size $\sim 10^3$ should be verifiable

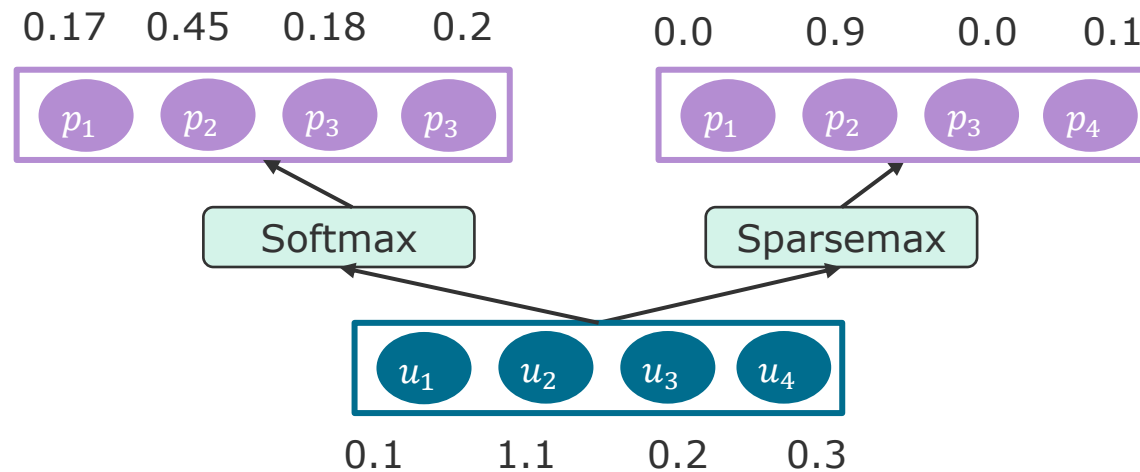
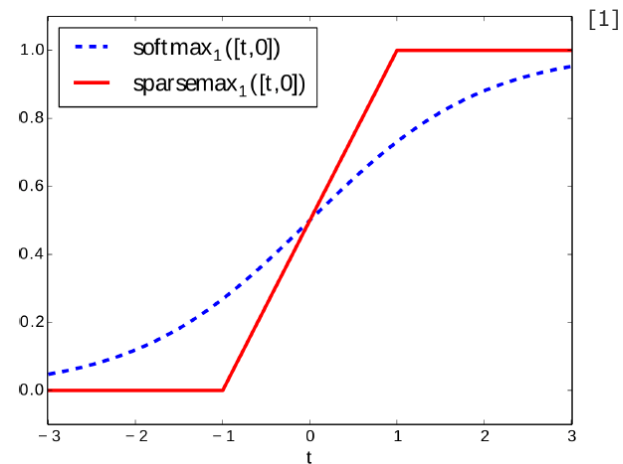
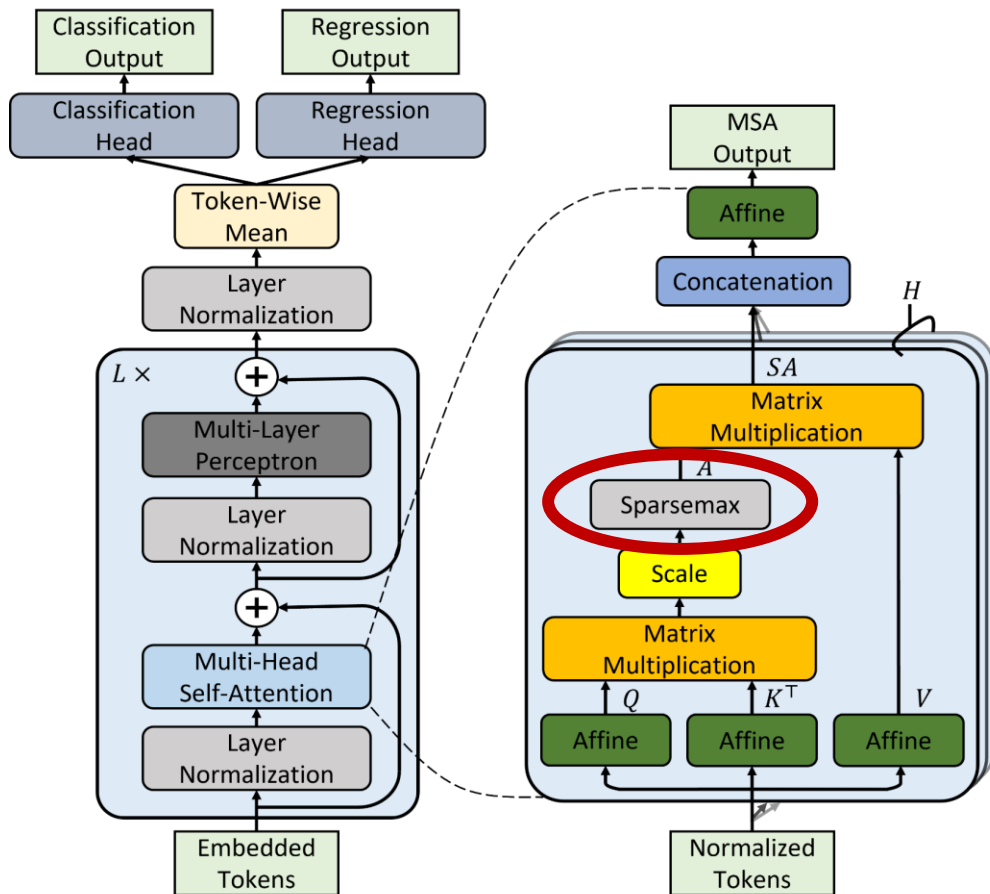
[1] Xu *et al.*, "Enabling complete NN verification with rapid and massively parallel incomplete verifiers," in *ICLR*, 2021.

[2] Tjeng *et al.*, "Evaluating robustness of NNs with MIP," in *ICLR*, 2019.

[3] Shi *et al.*, "Robustness verification for transformers," in *ICLR*, 2020.



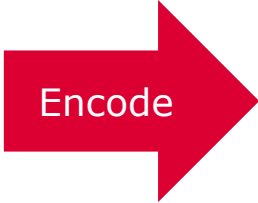
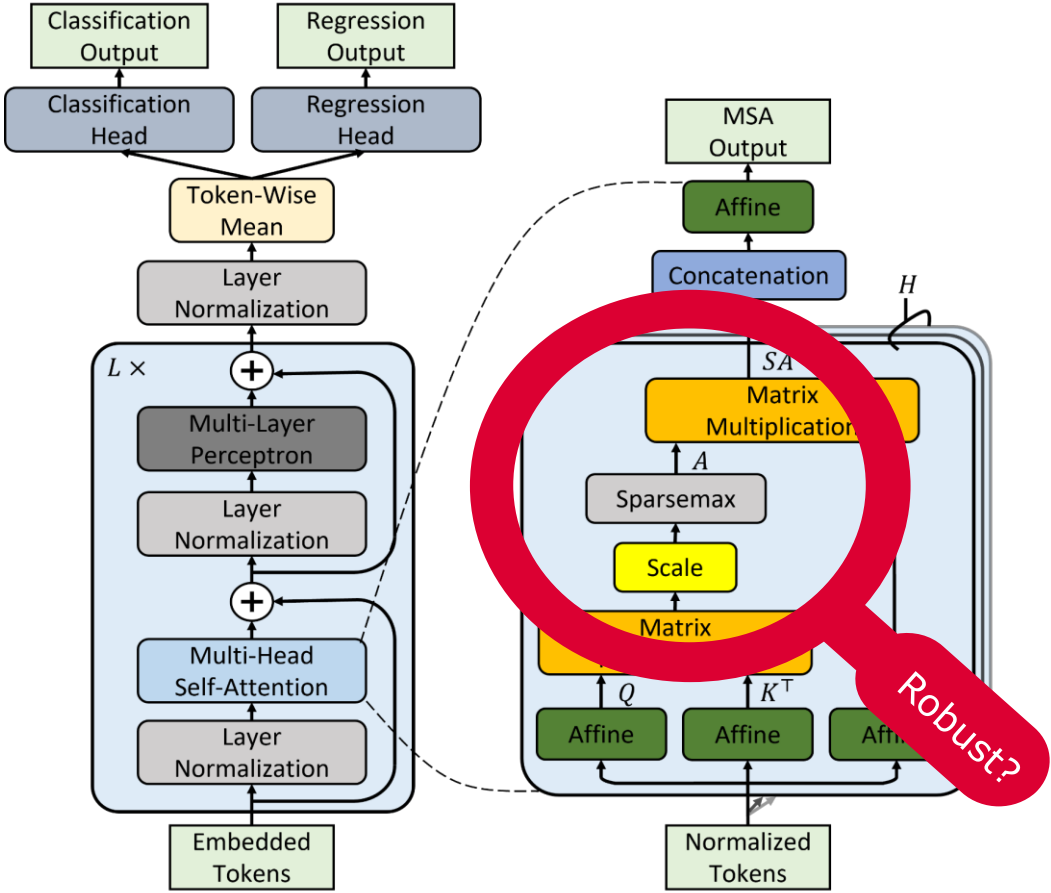
# Our Focus – Sparsemax Transformers



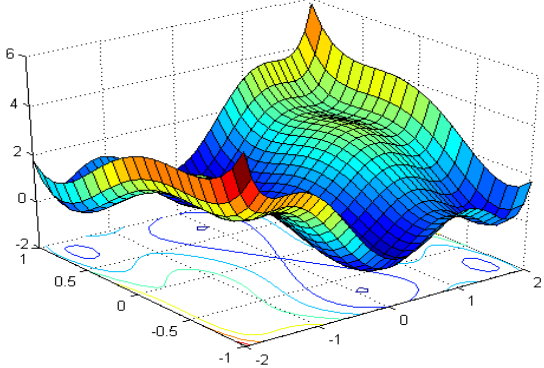
[1] Martins *et al.*, "From softmax to sparsemax," in *ICML*, 2016.

# Reducing Robustness Verification to an Optimization Problem

## Sparsemax Transformers



## MIQCP (Mixed Integer Quadratically Constrained Programming)



$$\min Dist_p(x, x') \quad \text{--- (1)}$$

$$s.t. \quad x' \in B_p(x), \quad \text{--- (2)}$$

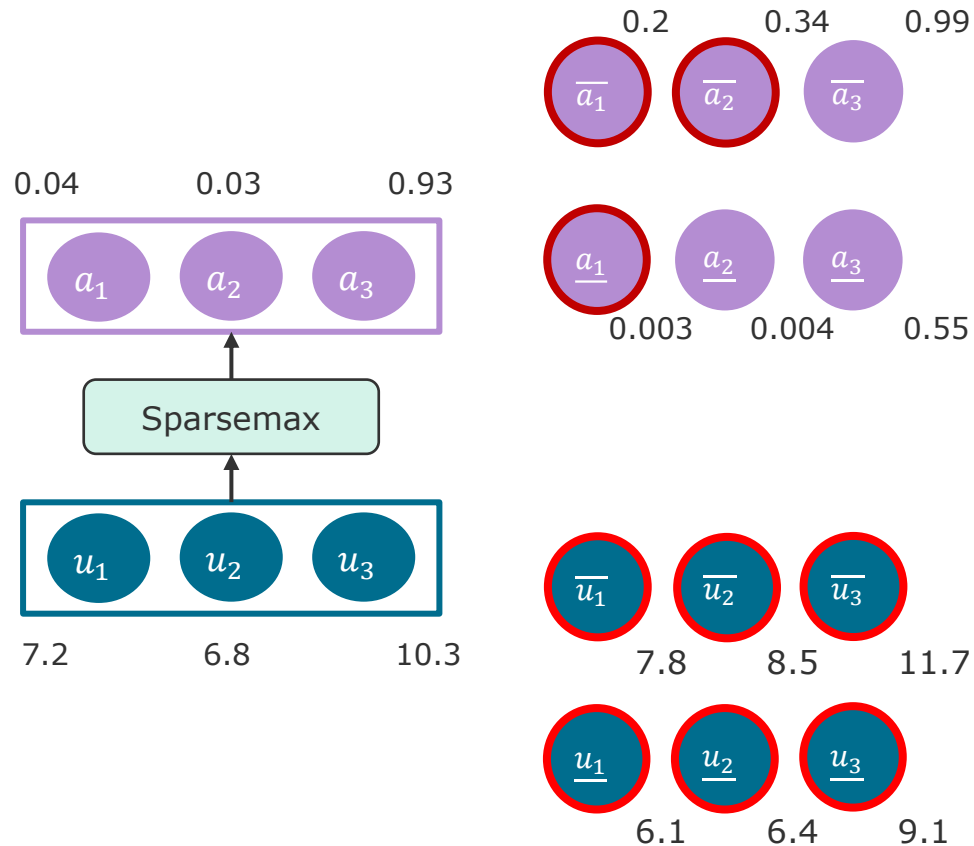
$$f(x) = gt(x) \wedge f(x') \neq gt(x). \quad \text{--- (3)}$$

- ① Find the min perturbation,
- ② within a budget, e.g.,  $\epsilon$
- ③ where the perturbed input causes a wrong prediction.

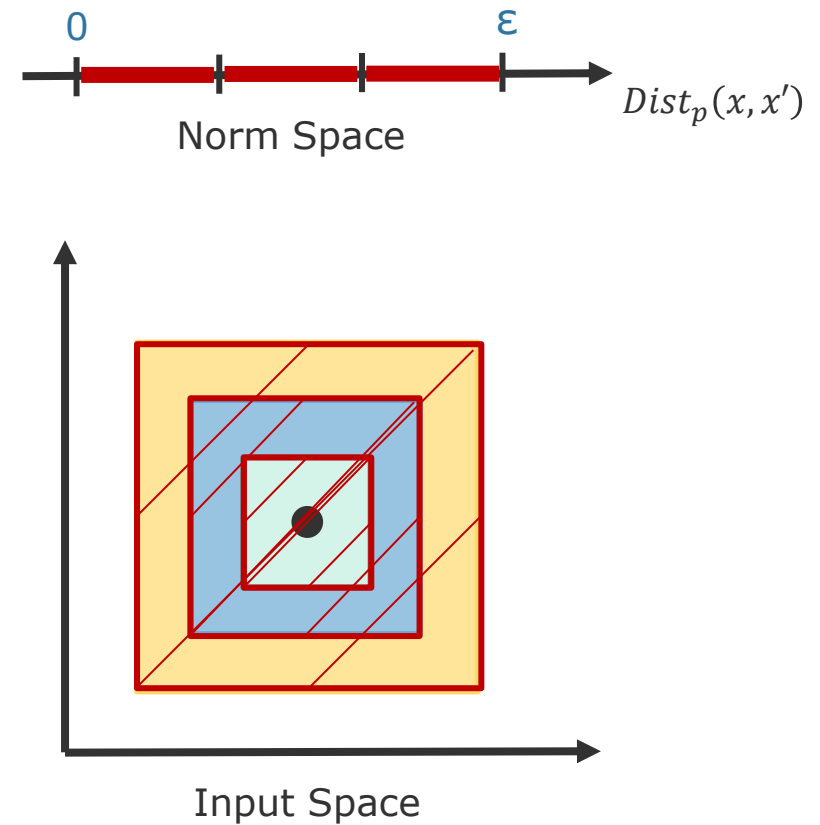
Robust?

# Accelerating Heuristics

- Interval analysis on Sparsemax activation<sup>[1]</sup>

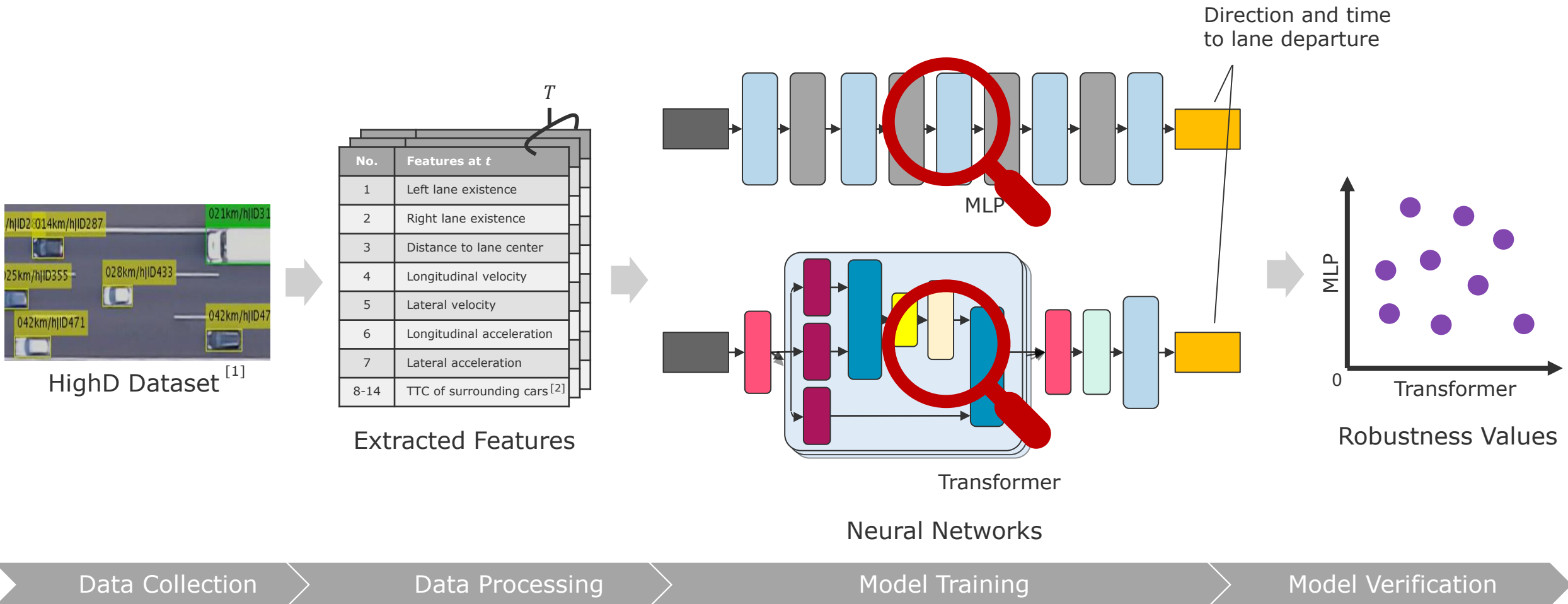


- Norm-space partitioning



[1] A recent improvement in Wei *et al.*, "Convex bounds on the softmax function with applications to robustness verification," in *AISTATS*, 2023.

# Experiment – Lane Departure Warning



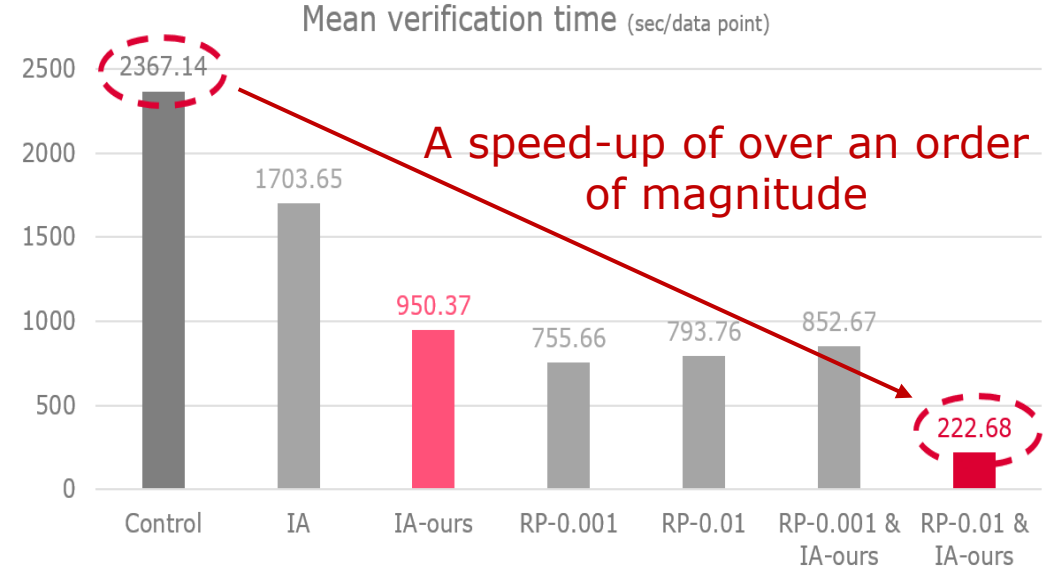
[1] <https://www.highd-dataset.com/>.

[2] TTC: Time to collision.

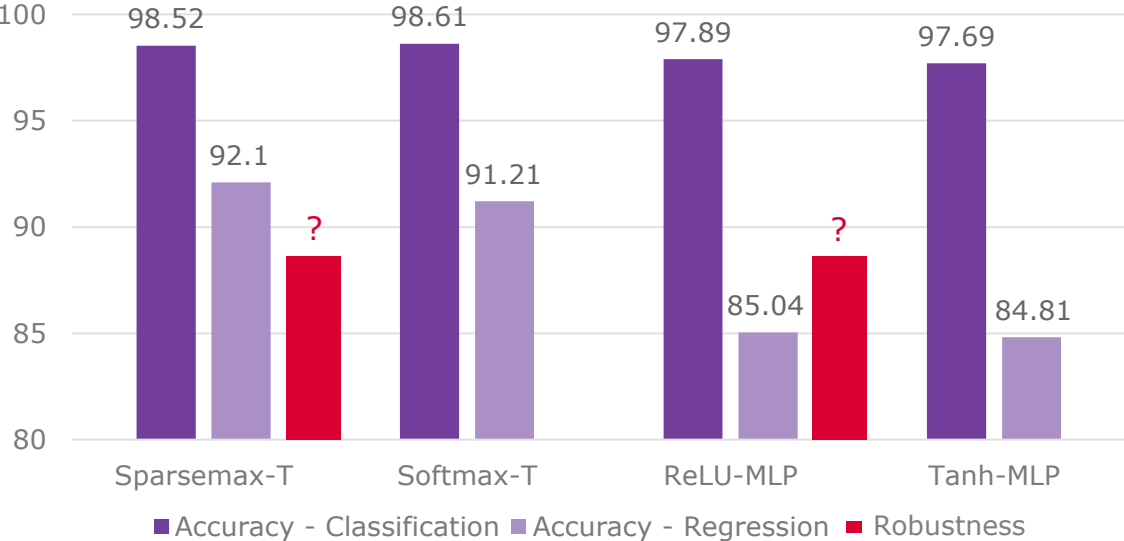


# Results – Ablation Study and Accuracy Assessment

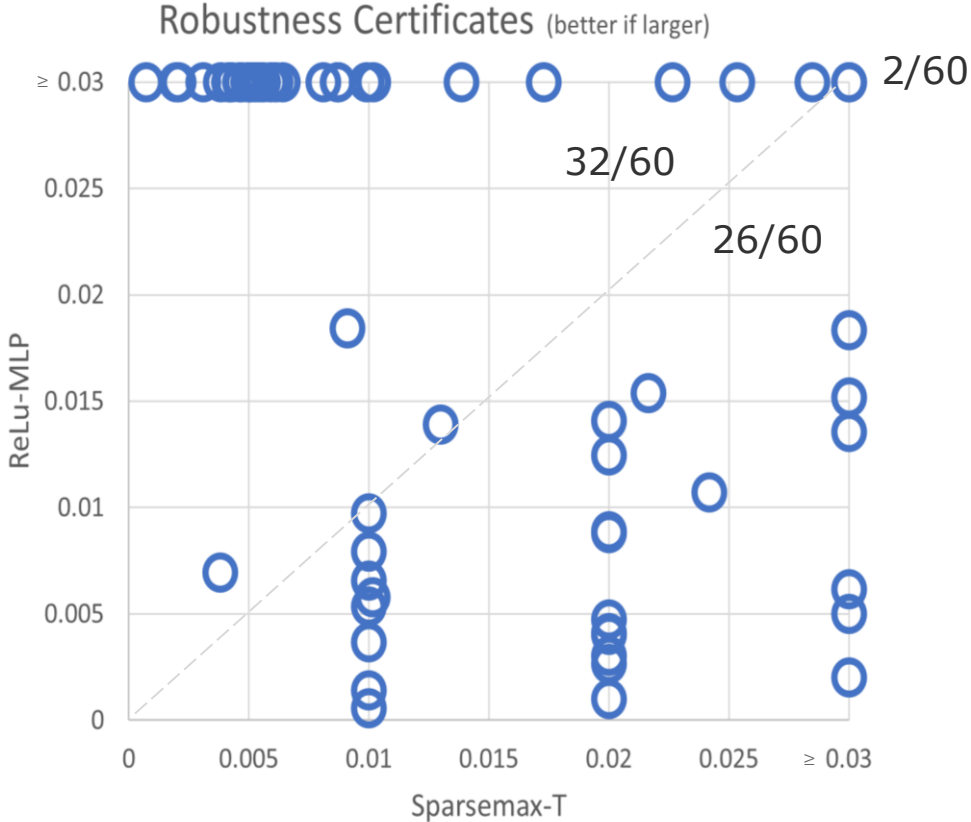
- On the proposed techniques



- Classification and regression accuracy



# Results – Robustness Benchmark



Statistics from verifying 5 ReLU-MLPs and 5 Sparsemax-Ts with 100 data points

	Mean	Var.	No. Max (0.01)
ReLU-MLPs	--	--	100
Sparsemax-Ts	0.0041	0.0006	69

- The Transformer seems less robust than the MLP.
- A recent empirical work in vision tasks concludes similarly.<sup>[1]</sup>
- Yet, some others disagree.<sup>[2, 3]</sup>

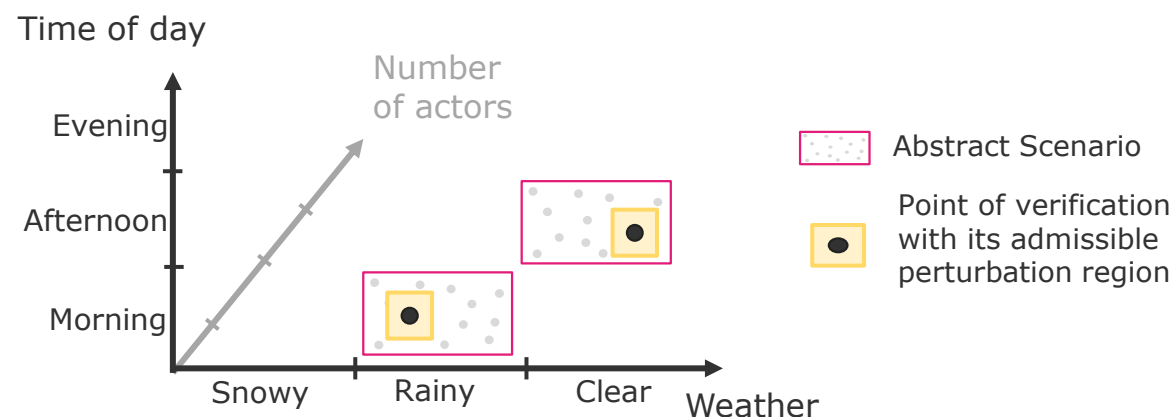
[1] Wang *et. al.*, "Can CNNs be more robust than Transformers," in *ICLR*, 2023.  
 [2] Bhojanapali *et. al.*, "Understanding robustness of transformers for image classification," in *ICCV*, 2021.  
 [3] Shao *et. al.*, "On the adversarial robustness of vision transformers," in *UCLR*, 2021.

# Summary

- Robustness is a safety-related concern in NNs.
- Our study focuses on exact robustness verification for a specific variant of Transformers.
- Robustness, as an application-oriented property, needs to be verified before NN deployment.

## Limitations and Open Directions

- Softmax not considered → Iterative bound tightening (e.g., with Branch-and-Bound)
- Small NNs and simple task → Real-world applications (e.g., via probabilistic verification)
- Point-wise analysis → Domain-covered assurance (e.g., with combinatorial testing)
- Design-time verification → Run-time verification (e.g., with different sensor modalities)

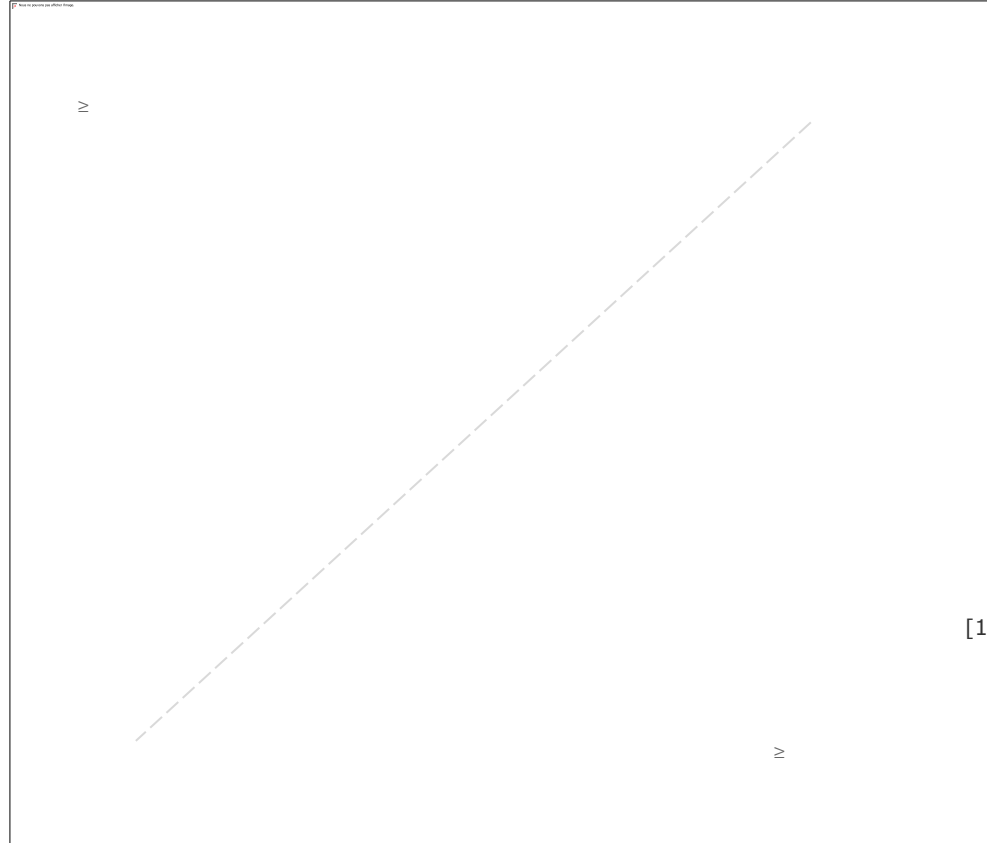


***DENSO***

Crafting the Core



# Incomplete Verification due to Timeouts



[1] The points marked with "Lower" give the lower bounds for the Transformer.