

Explainable and Trustworthy Traffic Sign Detection for Safe Autonomous Driving: An Inductive Logic Programming Approach

Zahra Chaghazardi

Department of Computer Science, University of Surrey
United Kingdom

`z.chaghazardi@surrey.ac.uk`

Saber Fallah

Connected and Autonomous Vehicles Lab, University of Surrey
United Kingdom

`s.fallah@surrey.ac.uk`

Alireza Tamaddoni-Nezhad

Department of Computer Science, University of Surrey
United Kingdom

`a.tamaddoni-nezhad@surrey.ac.uk`

Traffic sign detection is a critical task in the operation of Autonomous Vehicles (AV), as it ensures the safety of all road users. Current DNN-based sign classification systems rely on pixel-level features to detect traffic signs and can be susceptible to adversarial attacks. These attacks involve small, imperceptible changes to a sign that can cause traditional classifiers to misidentify the sign. We propose an Inductive Logic Programming (ILP) based approach for stop sign detection in AVs to address this issue. This method utilises high-level features of a sign, such as its shape, colour, and text, to detect categories of traffic signs. This approach is more robust against adversarial attacks, as it mimics human-like perception and is less susceptible to the limitations of current DNN classifiers. We consider two adversarial attacking methods to evaluate our approach: Robust Physical Perturbation (PR2) and Adversarial Camouflage (AdvCam). These attacks are able to deceive DNN classifiers, causing them to misidentify stop signs as other signs with high confidence. The results show that the proposed ILP-based technique is able to correctly identify all targeted stop signs, even in the presence of PR2 and AdvCam attacks. The proposed learning method is also efficient as it requires minimal training data. Moreover, it is fully explainable, making it possible to debug AVs.

1 Introduction

The popularity of AVs is rising rapidly because of their potential to reduce human error on the road, leading to safer transportation. AVs are believed to make more accurate perceptions and react faster than humans. Deep Neural Networks (DNNs) play a significant role in developing perception systems for AVs. However, DNNs face significant challenges that must be addressed before AVs can be deployed safely [7]. The major challenges facing DNN-based vision systems in autonomous driving are discussed below.

DNN-based systems are often considered "black boxes" because their logic is not transparent. Since it is difficult to explain how the system makes the prediction, it is challenging to debug them when they make a wrong decision. For example, misclassifying objects, such as mistaking shadows for pedestrians,

is a common problem in AVs and making decisions based on these misclassifications can lead to fatal accidents. Considering the fatal Uber accident [23], given that the AV's DNN-based decision-making is opaque, there is no way to debug the system and ensure such mistakes do not happen again. Moreover, using algorithms with ambiguous logic makes it impossible to evaluate and trust them. This means that regulatory approval is not applicable to stochastic-based AV vehicles.

Furthermore, DNNs face significant challenges when it comes to learning from small data and achieving out-of-distribution generalizability and transferability to new domains. In real-world scenarios, particularly in security domains, there is often a lack of large, annotated, and carefully curated data sets to train these systems. This can make it difficult for DNNs to acquire knowledge from a few examples and transfer it to new domains, unlike humans, who can do so with ease. Anomaly detection tasks, in particular, are affected by this challenge due to the rarity of anomalous data. Anomalies can be caused by errors, faults, or adversarial attacks, which can lead to security and safety hazards. Adversarial examples provide evidence of a network's weakness in achieving high generalisation performance [31]. Improving generalizability is crucial for adapting models to new domains when there is insufficient data. Given the lack of generalizability, current DNNs are not able to incrementally learn and improve when deployed in real-life situations and transfer knowledge from one domain to another (multi-domain) [28].

In the real world, DNNs are vulnerable to adversarial attacks and can be deceived easily. In adversarial cases, minor perturbations will lead to misclassifications with high confidence. Adversarial attacks have been investigated for different vision tasks, such as image classification, object detection, and semantic segmentation. For example, it is possible to change the red traffic light to green for AV [35], make people invisible to AI [32] using small crafted adversarial patches held in front of the body or make the AV to misinterpret a stop sign as a speed limit sign [16].

Researchers have suggested a few solutions, such as transfer learning for transferring knowledge to another domain, to address challenges associated with DNN classifiers. However, the proposed solutions partially solve the problems and have many limitations. For example, the transfer learning approach faces a significant challenge regarding data sharing and several legal issues such as privacy and property law [21].

To strengthen the safety of autonomous driving, this paper proposes an explainable ILP-based solution focusing on traffic sign detection. The proposed method mimics human perception to recognise traffic signs by detecting high-level features, including signs' geometric shapes, colours and contents, that differentiate them from other signs. While DNNs only use low-level (pixel-level) features that can be easily misled [16] and need a large amount of data, this traffic sign detector only needs a handful of training images and is fully robust against adversarial attacks.

Several studies have investigated the application of Inductive Logic Programming (ILP) in image recognition tasks. ILP has been employed in Logical Vision [12, 11], incorporating the abductive perception technique [29] to extract high-level interpretation of objects such as classical 2D shapes by utilising low-level primitives, such as high contrast points. ILP has also been used for 3D scene analysis [17] with 3D point cloud data. However, to our knowledge, a traffic sign detection based on the ILP has not been proposed previously for traffic sign classification. Therefore, our approach is a novel contribution to this context.

The paper is structured as follows. Section 2 surveys some successful adversarial examples in AVs. Section 3 describes the framework for robust traffic sign detection using ILP. Section 4 details experiments. In this section, the Aleph-based approach is compared with the Metagol-based approach. Metagol can learn hypotheses with only one positive and one negative example, while Aleph needs at least eight positive and negative examples to have the same accuracy as Metagol. Also, the ILP-based system is compared with the DNN-based classifier on adversarial examples. The results show that the ILP-based

approach is considerably more resilient to adversarial attacks. Finally, Section 5 summarises the outcomes and discusses further work.

2 Adversarial Attacks on AVs' Perception

In this section, we survey a sample of successful adversarial attacks in autonomous driving that easily deceived DNN-based vision classifiers. An adversarial attack aims to generate adversarial examples as the input for machine learning systems. However, adversarial examples are only negligibly modified from the real examples; they lead to misclassification [19].

When the fragility of deep neural networks to specific input perturbations was discovered for the first time, it was shown that an adversarial attack could turn a bus into an ostrich for an AI system [31]. Another algorithm named Show-and-Fool [8] was introduced to evaluate the robustness of an image captioning system. This method attained a 95.8% attack success rate for adversarial examples via applying a minor perturbation on image pixels which are invisible to humans, turning a stop sign into a teddy bear for the AI system.

The authors of [20] devised a method whereby semantic image segmentation could be attacked using adversarial perturbation to blend out (vanish) a desired target. They showed the existence of universal noise, which removes a target class (e.g. all pedestrians) from the segmentation while leaving it mostly unchanged otherwise. The robustness of the popular DNN-based semantic segmentation models evaluated against adversarial attacks on urban scene segmentation [2]. The results showed that the segmentation performances of all models seriously dropped after the attacks.

Later it was shown that adversarial examples could be misclassified by deep learning systems in real life [22]. Previous works have threatened the model by feeding machine learning classifiers directly, which is not always possible in the real world.

Another paper [16] proposed the Robust Physical Perturbations (RP2) technique to fool a Convolutional Neural Network (CNN) based road sign classifier in the physical world under various distances and viewpoints using different robust visual adversarial perturbations. This approach caused targeted misclassification, which changed a stop sign into a speed limit sign for the AI system. They also proposed a disappearance attack, causing a stop sign hidden from state-of-art object detectors like Mask R-CNN and YOLO [15]. An Adversarial Camouflage (AdvCam) approach [14] generated adversarial photos to fool a DNN classifier at various detecting angles and distances. With a few stains invisible to humans, this technique can cause the classifier to misclassify the objects, such as misidentifying a stop sign as a "barber shop" with .82% confidence.

Fig. 1 illustrates targeted stop signs with successful physical-world attacking approaches named RP2 and AdvCam, misleading the state-of-the-art DNN classifiers.

An Adaptive Square Attack (ASA) method [24] has been suggested that can attack the black box by generating invisible perturbation for traffic sign images, successfully leading to sign misclassification. Five adversarial attacks and four defence methods have been investigated on three driving models adopted in modern AVs [13]. They demonstrated that while these defence methods can effectively defend against a variety of attacks, none can provide adequate protection against all five attacks.

One recent work proposed three sticker application methods, namely RSA, SSA and MCSA, that can deceive the traffic sign recognition DNNs with realistic-looking stickers [4]. Another attack included painting the road, which targeted deep neural network models for end-to-end autonomous driving control [5]. Another work demonstrated a successful physical adversarial attack on a commercial classification system to deceive an AV's sign classifier[25].



Figure 1: Targeted physical perturbation by a) AdvCam and b) RP_2 misleading DNN classifiers, SL45 is speed limit 45 sign.

BadNets algorithm [18] was implemented to deceive a complex traffic sign detection system leading to maliciously misclassifying stop signs as speed-limit signs on real-world images.

These adversarial attacks on the deep-learning models pose a significant security threat to autonomous driving.

3 Robust Traffic Sign Detection Using ILP

Inductive Logic Programming (ILP) is a machine learning method which uses logic-based representation and inference. Depending on the type of logical inference and the search algorithm, there are different ILP systems, such as Aleph [3] and Metagol [10], that used in this paper.

Due to a logic-based representation and inference, ILP has the potential for human-like abstraction and reasoning. These logic-based AI approaches have the ability to learn unknown complex tasks with only a few examples. It complements deep learning because logic programs are interpretable and data-efficient, leading them towards a strong generalisation. Moreover, these rule-based approaches, which are explicitly symbolic, are sometimes considered safer than neural approaches [1].

ILP aims to learn a hypothesis (rule) using a few positive and negative examples and Background Knowledge (BK); this induced rule, alongside BK, should cover as many positive and as few negative examples as possible [26]. For inducing the rules, BK should include all essential predicates to represent the relevant information.

One of the advantages of ILP is its ability to use BK, including facts and rules in the form of logical expressions, which could be related. In ILP, choosing appropriate BK based on well-selected features is essential to obtaining good results [9]. Moreover, using BK makes ILP incremental. For example,

suppose we want to learn animal signs in traffic sign detection, choosing sign "a" contains an animal(*contains(a, animal)*) as a BK, which holds when traffic sign "a" has an animal symbol. Then we provide BK with various different animal shaped symbols (deer, cow, ...). In that case, if we see a new animal sign that doesn't exist in our BK, we can add it to our BK without relearning, and there is no need to change the hypothesis. This feature makes it possible to have real-time interaction with drivers towards customised autonomous driving.

Our proposed ILP-based stop sign detection system is demonstrated in Fig. 2. The first step is pre-processing all the images, including training and test images, and turning them into a symbolic representation to provide BK. In the pre-processing phase, high-level features of traffic sign images, including colour, shape, text and digits, are extracted and represented as a set of logical facts for the next step. For feature extraction, computer vision tools such as OpenCV can extract high-level features using low-level features such as pixel colours or colour gradients.

In the next step, a set of positive and negative training examples (E) and a set of logical facts as BK extracted from the previous step will be provided to the ILP system. The system aims to learn a hypothesis H such that $B, H \models E$ where \models is logical entailment.

We use Aleph and Metagol as the ILP system to induce the rule for stop sign detection.

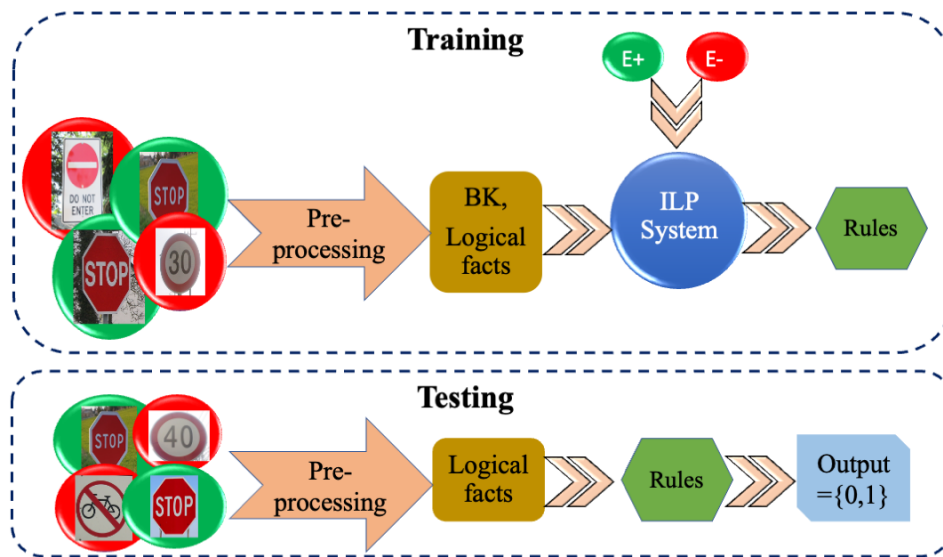


Figure 2: ILP- based traffic sign classifier

We used Aleph5 as the ILP system in one of our experiments; it is an old ILP system developed in Prolog and based on inverse entailment. Aleph's algorithm resolves the relationship between the determination predicate and the determining predicate to generate a general theory.

Metagol is employed in our other experiment. It is an ILP system based on Meta Interpretive Learning (MIL) [27] implemented in Prolog. By instantiating metarules, MIL learns logic programs from examples and BK. In addition, MIL not only learns the recursive definition and fetches higher-order meta-rules but also supports predicate invention.

4 Experimental Evaluation

This experiment aims to learn "*traffic_sign*", which is the target predicate. For simplicity, only the stop sign is investigated; other traffic signs can be included to have a complete traffic sign classifier. We provide Aleph and Metagol with the same BK. The Aleph mode declarations are illustrated in Table 1, and the Metagol-based system is supplied with the metarules demonstrated in Table 2, uppercase letters represent predicate symbols (second-order variables), and lowercase letters represent variables.

Table 1: Mode declarations for Aleph experiments.

:	$-modeh(1, traffic_sign(+sign, \#class)).$
:	$-modeb(*, colour(+sign, \#colour)).$
:	$-modeb(*, shape(+sign, \#shape)).$
:	$-modeb(*, word(+sign, -w)).$
:	$-modeb(*, closely_match(+w, \#word)).$
:	$-modeb(*, number(+sign, -n)).$
:	$-modeb(*, digits(+n, \#int)).$

In Aleph mode declaration, "*modeh*" indicates that the predicate should appear in the head of the hypothesis, and "*modeb*" indicates that it should be in the body of the induced hypothesis. According to Table 1, six predicates can be used in the body of the induced hypothesis which. The meaning of each predicate is defined as follows:

- $traffic_sign(a, \#class)$, which holds when the sign "a" belongs to a specific category of traffic sign determined by #class (e.g. a stop sign).
- $colour(a, \#colour)$, which holds when a certain #colour(e.g. red) exists in the sign "a".
- $shape(a, \#shape)$, which holds when the shape of sign "a" is a specific shape determined by #shape(e.g. circle).
- $has_word(a, a_w1)$, which holds when the sign "a" has the word a_w1 on it.
- $closely_match(a_w1, w)$, which holds when the word " a_w1 " closely matches the word "w" (e.g. stop).
- $number(a, a_n1)$, which holds when the sign "a" has the number " a_n1 ".
- $digits(a_n1, n)$, which holds when the number " a_n1 " includes "n" (e.g. 60)

Table 2: Employed metarules in Metagol experiment.

Name	Metarule
Identify	$P(x, y) \leftarrow Q(x, y)$
Inverse	$P(x, y) \leftarrow Q(y, x)$
Precon	$P(x, y) \leftarrow Q(x), R(x, y)$
Postcon	$P(x, y) \leftarrow Q(x, y), R(y)$
Chain	$P(x, y) \leftarrow Q(x, z), R(z, y)$
Recursion	$P(x, y) \leftarrow Q(x, z), P(z, y)$

Table 3: Extracted features for a positive (p1) and negative (n1) examples.

Pos example(p1)	Neg example(n1)
<i>color(p1, red)</i> .	<i>color(n1, red)</i> .
<i>color(p1, white)</i> .	<i>color(n1, white)</i> .
<i>shape(p1, octagon)</i> .	<i>shape(n1, Circle)</i> .
<i>has_word(p1, p1_w1)</i> .	<i>number(n1, n1_d1)</i> .
<i>closely_match(p1_w1, stop)</i> .	<i>digits(n1_d1, 30)</i> .

To further explain the process of converting images into a set of logical facts, we take one positive example, a stop sign named p1, and one negative example, a speed limit sign named n1. In the pre-processing stage, the high-level features of these traffic signs were extracted to be included in the BK. The details of these features and their corresponding logical representation are presented in Table 3.

These logical facts, together with the names of the positive and negative examples (such as sign(p1) as a positive and sign(n1) as a negative example), will enable the ILP system to induce a hypothesis (logical rule). Finally, the ILP system recognises the new traffic signs using this induced rule.

4.1 Material and Method

Base data set. Our base data set includes traffic sign images without any adversarial perturbation. They have been downloaded from Wikimedia Commons as no-restriction images. Positive examples contain stop sign images, and negative examples include other traffic signs excluding stop sign images. Normal positive and negative examples are shown in Fig. 3.

Adversarial test data set. To evaluate the robustness of the ILP stop sign detector against adversarial attacks, we used the targetted traffic signs attacked by RP_2 and AdvCam.

RP_2 is a general attack algorithm for misleading standard-architecture road sign classifiers. It generates visual adversarial perturbations, such as black and white stickers attached to a traffic sign to mislead the classifier. The RP_2 data set contains three types of perturbation, stop signs perturbed by subtle, camouflage graffiti and camouflage art attacks viewed from different angles.

AdvCam is an approach for creating camouflage physical adversarial images to fool state-of-the-art DNN-based image classifiers. This approach can make the classifier identify a stop sign as a "barber shop" with high confidence. This paper will utilise targetted stop signs with Advcam with different stain styles to evaluate the proposed ILP sign classifier.

The feature recognition element should extract high-level features, including the traffic sign's border shape, colour and text; symbol extraction can be added in future; *OpenCV* is employed for this purpose. First, the image background is removed utilising *rembg* and then pre-processed by *cv2.bilateralFilter* for noise reduction.

Colour masks were then employed using *cv2.inRange* for colour detection, and small areas were ignored. After that, by applying morphological operations, colour masks were post-processed.

For text and digit detection, *EasyOCR* is utilised; if the detected item is a word, it will be investigated if the detected word closely matches some common words in traffic signs; for example, it should have at least three letters in common with the word STOP to recognise as a stop word.

For shape detection, *cv2.findContours* is applied on detected colour masks, and *cv2.approxPolyDP* is utilised for polygon detection.

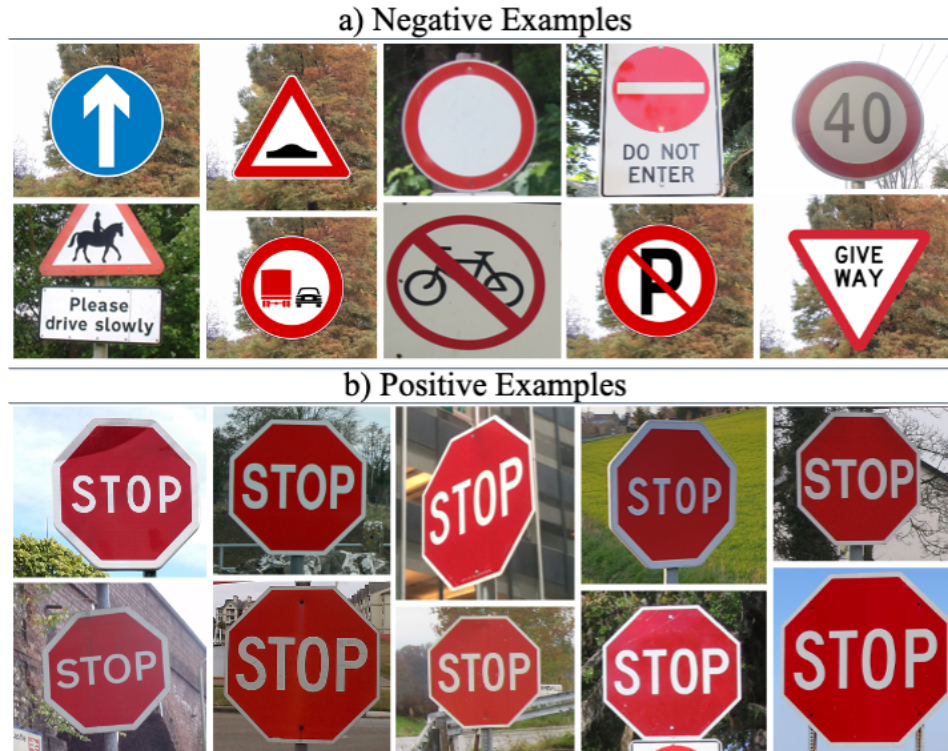


Figure 3: Base data set for training and testing

Convolutional Neural Network. To compare the results, a well-known CNN classifier [34] is utilised that is trained on the German Traffic Sign Recognition Benchmark (GTSRB) [30]. The evaluation of this architecture achieved 97.6% accuracy on the GTSRB test data set.

The base data set is utilised for training the ILP systems (Aleph and Metagol). First, we randomly select an equal number of positive and negative examples in each run, so the default accuracy is 50% for this training data set. Next, the ILP-based systems try to find a hypothesis that covers as many positive and as few negative examples as possible. Then the remaining examples in the data set are used as a test data set for evaluation to determine the accuracy. This process is repeated one hundred times, and average accuracy is calculated for each certain number of positive and negative examples in the training set. Therefore we have a fair comparison between Aleph and Metagol regarding the size of the required training data set.

The data and the code used in this experiment are available on GitHub [6].

4.2 Results and Discussion

Fig. 4 illustrates the average accuracy of Aleph and Metagol-based ILP systems with increasing training examples. According to this figure, Metagol can find a hypothesis with 100% accuracy on the test data set including only one positive and one negative example. In comparison, Aleph starts learning with at least two positive and two negative examples with around 65% accuracy. Aleph can reach the same level of accuracy as Metagol by learning from eight positive and negative examples. According to these results, Metagol is more data-efficient than Aleph. In this figure, the orange curve shows the default accuracy, which is equal to 50% because the number of negative and positive examples are equal in each

run.

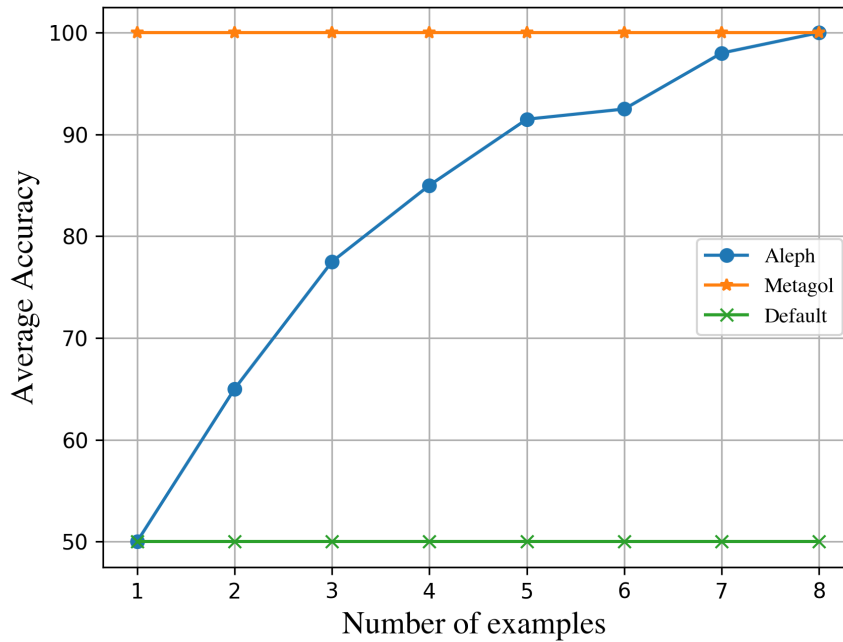


Figure 4: Average accuracy of Aleph vs Metagol with an increasing number of training examples from the base data set (equal positive and negative sets)

The hypothesis (a logic program) induced by Metagol with only one set of positive and negative examples is the same as the learned rule by Aleph with eight positive and negative examples. It is entirely accurate on the base test data set and is shown below:

```
traffic_sign(A, stop_sign):-
    has_word(A, A_w1),
    closely_match(A_w1, stop).
```

This learned rule is completely explainable and matches human interpretation. The rule says the traffic sign "A" is a stop sign when the two literals *has_word(A, A_w1)* and *closely_match(A_w1, stop)* hold, i.e. if the sign contains a word and that word closely matches stop, that sign would be predicted a stop sign.

The performance of this hypothesis is evaluated on the base data set and attack data sets, including RP_2 (subtle, camouflage graffiti and camouflage art attacks) and AdvCam with different stains. The accuracy of this rule on all test data sets is 100%. While the ILP-based sign detector shows a perfect performance, the DNN-based classifier shows abysmal performance facing manipulated images.

Table 4 compares the results of the DNN-based classifier and the ILP-based classifier on different data sets. The DNN-based classifier is trained on the GTSRB data set, which contains more than 50,000 images. The Aleph-based classifier is trained on eight positive and negative examples, while the Metagol approach is trained on only one negative and one positive example. It shows that while ILP-based systems can learn from small amounts of data, they are more resilient to noise and adversarial attacks.

data set	DNN-based	ILP-based
Base	100%	100%
subtle	0	100%
RP_2 camouflage graffiti	0%	100%
camouflage art attacks	6.6%	100%
AdvCam	66.6%	100%

Table 4: Comparing the results of the hypothesis induced by the proposed ILP-based approach (Aleph and Metagol) on different test data sets with a DNN classifier.

5 Conclusions

DNN-based traffic sign classifiers need a large amount of data for training, and it has been shown that they are vulnerable to adversarial attacks or natural noise. They are also not explainable; consequently, there is no way to debug them. While DNN-based classifiers suffer from these problems, we propose an ILP-based approach for traffic sign detection in autonomous vehicles to address these issues.

Our proposed technique mimics humans in traffic sign detection and uses high-level features of a sign, such as colour and shape, for detection. Therefore this method is data efficient, explainable and able to withstand adversarial attacks that cannot easily deceive humans.

The results indicate that our approach with only a handful of training data can induce logical rules easily understandable by humans. Furthermore, it significantly outperforms the deep learning approach regarding adversarial attacks. It shows a 100% accuracy on the data set targetted with RP_2 and AdvCam attacking approaches, while a DNN-based classifier performs poorly on these data sets.

For future works, we suggest employing DNN for high-level feature extraction (shapes or symbols) in traffic signs. Integrate Machine Learning and Logic programming in AV applications will use both the strengths of machine learning and symbolic AI (knowledge and reasoning) to address the AI obstacles.

Acknowledgments

The first author would like to acknowledge her PhD grant funding from the Breaking Barriers Studentship Award at the University of Surrey. Also, we would like to acknowledge the support of Dany Varghese with Aleph PyILP [33].

References

- [1] Greg Anderson, Abhinav Verma, Isil Dillig & Swarat Chaudhuri (2020): *Neurosymbolic reinforcement learning with formally verified exploration*. *Advances in neural information processing systems* 33, pp. 6172–6183.
- [2] Anurag Arnab, Ondrej Miksik & Philip HS Torr (2018): *On the robustness of semantic segmentation models to adversarial attacks*. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 888–897.
- [3] Ashwin Srinivasan (2001): *The aleph manual*. <https://www.cs.ox.ac.uk/activities/programinduction/Aleph/aleph.html>.

- [4] Yasin Bayzidi, Alen Smajic, Fabian Hüger, Ruby Moritz, Serin Varghese, Peter Schlicht & Alois Knoll (2022): *Traffic sign classifiers under physical world realistic sticker occlusions: A cross analysis study*. In: *2022 IEEE Intelligent Vehicles Symposium (IV)*, IEEE, pp. 644–650, doi:10.1109/CVPR.2017.634.
- [5] Adith Bolor, Xin He, Christopher Gill, Yevgeniy Vorobeychik & Xuan Zhang (2019): *Simple physical adversarial examples against end-to-end autonomous driving models*. In: *2019 IEEE International Conference on Embedded Software and Systems (ICCESS)*, IEEE, pp. 1–7, doi:10.1145/1081870.1081950.
- [6] Zahra Chaghazardi (2022): *Traffic Sign Detection using ILP*. <https://github.com/Chaghazardi/Traffic-Sign-Detection-using-ILP>. Available at <https://github.com/Chaghazardi/Traffic-Sign-Detection-using-ILP>.
- [7] Zahra Chaghazardi, Saber Fallah & Alireza Tamaddoni-Nezhad (2023): *A Logic-based Compositional Generalisation Approach for Robust Traffic Sign Detection*. In: *International Joint Conference on Artificial Intelligence 2023 Workshop on Knowledge-Based Compositional Generalization*.
- [8] Hongge Chen, Huan Zhang, Pin-Yu Chen, Jinfeng Yi & Cho-Jui Hsieh (2017): *Attacking visual language grounding with adversarial examples: A case study on neural image captioning*. arXiv preprint arXiv:1712.02051.
- [9] Andrew Cropper, Sebastijan Dumancic & Stephen H Muggleton (2020): *Turning 30: New Ideas in Inductive Logic Programming*. In: *IJCAI*, doi:10.24963/ijcai.2020/673.
- [10] Andrew Cropper & Stephen H. Muggleton (2016): *Metagol System*. <https://github.com/metagol/metagol>. Available at <https://github.com/metagol/metagol>.
- [11] Wang-Zhou Dai, Stephen Muggleton, Jing Wen, Alireza Tamaddoni-Nezhad & Zhi-Hua Zhou (2018): *Logical vision: One-shot meta-interpretive learning from real images*. In: *Inductive Logic Programming: 27th International Conference, ILP 2017, Orléans, France, September 4-6, 2017, Revised Selected Papers 27*, Springer, pp. 46–62, doi:10.1016/j.cviu.2007.08.003.
- [12] Wang-Zhou Dai, Stephen H Muggleton & Zhi-Hua Zhou (2015): *Logical Vision: Meta-Interpretive Learning for Simple Geometrical Concepts*. In: *ILP (Late Breaking Papers)*, pp. 1–16.
- [13] Yao Deng, Xi Zheng, Tianyi Zhang, Chen Chen, Guannan Lou & Miryung Kim (2020): *An analysis of adversarial attacks and defenses on autonomous driving models*. In: *2020 IEEE international conference on pervasive computing and communications (PerCom)*, IEEE, pp. 1–10, doi:10.1109/PerCom45495.2020.9127389.
- [14] Ranjie Duan, Xingjun Ma, Yisen Wang, James Bailey, A Kai Qin & Yun Yang (2020): *Adversarial camouflage: Hiding physical-world attacks with natural styles*. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1000–1008.
- [15] Kevin Eykholt, Ivan Evtimov, Earleence Fernandes, Bo Li, Amir Rahmati, Florian Tramer, Atul Prakash, Tadayoshi Kohno & Dawn Song (2018): *Physical adversarial examples for object detectors*. arXiv preprint arXiv:1807.07769 1(3), p. 4.
- [16] Kevin Eykholt, Ivan Evtimov, Earleence Fernandes, Bo Li, Amir Rahmati, Chaowei Xiao, Atul Prakash, Tadayoshi Kohno & Dawn Song (2018): *Robust physical-world attacks on deep learning visual classification*. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1625–1634.
- [17] Reza Farid & Claude Sammut (2014): *Plane-based object categorisation using relational learning*. *Machine Learning* 94, pp. 3–23, doi:10.1007/s10994-013-5352-9.
- [18] Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt & Siddharth Garg (2019): *Badnets: Evaluating backdoor attacks on deep neural networks*. *IEEE Access* 7, pp. 47230–47244, doi:10.1109/TKDE.2009.191.
- [19] Jie Gui, Zhenan Sun, Yonggang Wen, Dacheng Tao & Jieping Ye (2021): *A review on generative adversarial networks: Algorithms, theory, and applications*. *IEEE transactions on knowledge and data engineering*.
- [20] Jan Hendrik Metzen, Mummadi Chaithanya Kumar, Thomas Brox & Volker Fischer (2017): *Universal adversarial perturbations against semantic image segmentation*. In: *Proceedings of the IEEE international conference on computer vision*, pp. 2755–2764.

- [21] Mauritz Kop (2020): *Machine learning & EU data sharing practices*. In: *TTLF Newsletter on Transatlantic Antitrust and IPR Developments*, Stanford-Vienna Transatlantic Technology Law Forum, Transatlantic Antitrust
- [22] Alexey Kurakin, Ian J Goodfellow & Samy Bengio (2018): *Adversarial examples in the physical world*. In: *Artificial intelligence safety and security*, Chapman and Hall/CRC, pp. 99–112, doi:10.1201/9781351251389-8.
- [23] T. B. Lee (2018): *Software bug led to death in ubers self-driving crash*. Available at <https://arstechnica.com/tech-policy/2018/05/report-software-bug-led-to-death-in-ubers-self-driving-crash/>.
- [24] Yujie Li, Xing Xu, Jinhui Xiao, Siyuan Li & Heng Tao Shen (2020): *Adaptive square attack: Fooling autonomous cars with adversarial traffic signs*. *IEEE Internet of Things Journal* 8(8), pp. 6337–6347, doi:10.1109/CVPR.2018.00957.
- [25] Nir Morgulis, Alexander Kreines, Shachar Mendelowitz & Yuval Weisglass (2019): *Fooling a real car with adversarial traffic signs*. Available at <https://arxiv.org/abs/1907.00374>.
- [26] Stephen Muggleton (1991): *Inductive logic programming*. *New generation computing* 8(4), pp. 295–318, doi:10.1007/BF03037089.
- [27] Stephen H Muggleton, Dianhuan Lin & Alireza Tamaddoni-Nezhad (2015): *Meta-interpretive learning of higher-order dyadic datalog: Predicate invention revisited*. *Machine Learning* 100(1), pp. 49–73, doi:10.1007/s10994-014-5471-y.
- [28] Bukola Salami, Keijo Haataja & Pekka Toivanen (2021): *State-of-the-Art Techniques in Artificial Intelligence for Continual Learning: A Review*. *FedCSIS (Position Papers)*, pp. 23–32.
- [29] Murray Shanahan (2005): *Perception as abduction: Turning sensor data into meaningful representation*. *Cognitive science* 29(1), pp. 103–134, doi:10.1207/s15516709cog2901_5.
- [30] Johannes Stalkamp, Marc Schlipf, Jan Salmen & Christian Igel (2012): *Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition*. *Neural networks* 32, pp. 323–332, doi:10.1016/j.neunet.2012.02.016.
- [31] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow & Rob Fergus (2013): *Intriguing properties of neural networks*. *arXiv preprint arXiv:1312.6199*.
- [32] Simen Thys, Wiebe Van Ranst & Toon Goedemé (2019): *Fooling automated surveillance cameras: adversarial patches to attack person detection*. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pp. 0–0.
- [33] Dany Varghese (2022): *PyILP*. <https://github.com/danyvarghese/PyILP>. Available at <https://github.com/danyvarghese/PyILP>.
- [34] Vivek Yadav (2016): *German sign classification using deep learning neural networks*. <https://github.com/vxy10/p2-TrafficSigns>.
- [35] Chen Yan, Zhijian Xu, Zhanyuan Yin, Xiaoyu Ji & Wenyan Xu (2022): *Rolling Colors: Adversarial Laser Exploits against Traffic Light Recognition*. *arXiv preprint arXiv:2204.02675*.