

Finding Property Violations through Network Falsification: Challenges, Adaptations and Lessons Learned from OpenPilot

Meriel von Stein
meriel@virginia.edu
University of Virginia
Charlottesville, VA, United States

Sebastian G. Elbaum
selbaum@virginia.edu
University of Virginia
Charlottesville, VA, United States

ABSTRACT

OpenPilot is an open source system to assist drivers by providing features like automated lane centering and adaptive cruise control. Like most systems for autonomous vehicles, OpenPilot relies on a sophisticated deep neural network (DNN) to provide its functionality, one that is susceptible to safety property violations that can lead to crashes. To uncover such potential violations before deployment, we investigate the use of falsification, a form of directed testing that analyzes a DNN to generate an input that will cause a safety property violation. Specifically, we explore the application of a state-of-the-art falsifier to the DNN used in OpenPilot, which reflects recent trends in network design. Our investigation reveals the challenges in applying such falsifiers to real-world DNNs, conveys our engineering efforts to overcome such challenges, and showcases the potential of falsifiers to detect property violations and provide meaningful counterexamples. Finally, we summarize the lessons learned as well as the pending challenges for falsifiers to realize their potential on systems like OpenPilot.

KEYWORDS

testing, autonomous systems, adversarial examples, falsification, formal methods, neural nets

ACM Reference Format:

Meriel von Stein and Sebastian G. Elbaum. 2022. Finding Property Violations through Network Falsification: Challenges, Adaptations and Lessons Learned from OpenPilot. In *37th IEEE/ACM International Conference on Automated Software Engineering (ASE '22)*, October 10–14, 2022, Rochester, MI, USA. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/3551349.3559500>

1 INTRODUCTION

In 2016, comma.ai unveiled OpenPilot, one of the first open-source systems for self-driving vehicles that today forms part of a community [12, 14] competing with their commercial counterparts [7, 16]. OpenPilot enables vehicles with Automated Lane Centering, Adaptive Cruise Control, Forward Collision Warning, and Lane Departure Warning. These features assist drivers in limited contexts, prioritizing highway driving [11]. Since its inception, OpenPilot has grown to support over 150 commercial vehicles. Its userbase

has grown substantially, with over 2,750 new users a week by the end of 2020 [6] and the recent debut of the Comma 3 device.

OpenPilot autonomy relies heavily on the supercombo DNN, an end-to-end driving model for making predictions about various features of the driving environment, including: optimal paths over upcoming timesteps, position and probability of lead cars, locations of lane lines and road edges, and the probability of OpenPilot disengagement, braking, and acceleration. Model inputs are taken from camera sensors, human input to the system, system configuration, and recurrent system state. supercombo predictions are used by many different parts of the OpenPilot software, such as lateral (turning) and longitudinal (gas/brake) planning and identifying driving events such as laneChange or roadCameraError.

The functionality of OpenPilot depends on images that can be reliably interpreted by the supercombo network. To maintain safe operation, network output must conform to certain safety properties. One such property might specify that adding a low level of noise to a camera image input should not change lane line predictions by more than $\frac{1}{4}$ standard lane width. To test these properties, DNN falsification aims to uncover violations related to the robustness of DNN predictions. Falsification techniques extend adversarial example generation approaches to find violations to safety properties within some measure of local robustness [1–5, 8, 15, 17]. Falsifiers take in a network and a safety property and attempt to find a violation of that property, resulting in either a counterexample or a timeout. For the previous property associated with lane line prediction, a counterexample would consist of an image plus minimal noise that causes supercombo to predict lane lines that deviate by over $\frac{1}{4}$ lane width from predictions on the original image.

Recent studies involving falsifiers illustrate their potential to find counterexamples for valuable safety properties [13]. However, the networks in these studies do not capture the complexity of DNNs encountered in autonomous vehicles such as supercombo, as well as other real-world applications of DNNs at scale. Moreover, recent complex networks often involve multi-dimensional inputs and outputs with inherent dependencies that require more sophisticated properties than early falsifiers can represent. In this work we:

- Analyze the challenges in applying a state-of-the-art falsification framework, DNNF [13], to a state-of-the-art DNN for autonomous vehicles, supercombo. Our analysis reveals that recent DNNs introduce conceptual issues beyond scale and format that challenge today’s falsifiers, such as multiple inputs, multi-dimensional outputs, and complex operators and architecture.
- Extend DNNF to enable its application to a set of safety properties of supercombo and perform a study illustrating the



This work is licensed under a Creative Commons Attribution International 4.0 License.

ASE '22, October 10–14, 2022, Rochester, MI, USA
© 2022 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-9475-8/22/10.
<https://doi.org/10.1145/3551349.3559500>

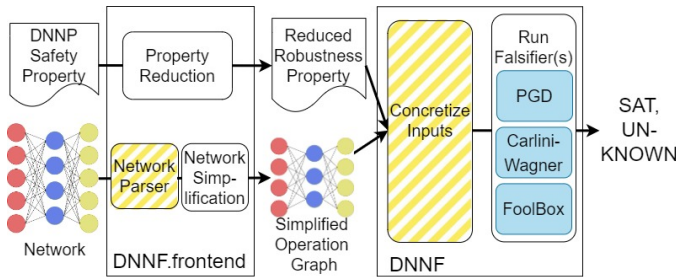


Figure 1: Approach Overview. Extended functionality is shown in yellow stripes.

potential for falsification to uncover meaningful counterexamples that will lead to violations.

- Collect lessons learned from our analysis, application, and tool extension, as well as a recognition of the challenges that remain for falsification to be more broadly deployed on DNNs used in autonomous vehicles.

2 OBSERVATIONS AND CHALLENGES IN FALSIFYING SUPERCOMBO

The super combo model is a large-scale, multi-input network with a relatively complex architecture. It consists of a series of 70 convolutional layers, followed by a block of 23 layers with multiple connections between each layer. `supercombo` takes in four input tensors: `input_imgs`, `desire`, `traffic_convention`, and `initial_state`, all of which enter at different points in the network. Only `input_imgs` is processed by the convolutions at the beginning of the network. The `input_imgs` tensor is a composite of two reordered image tensors collected consecutively by an onboard camera. The network outputs a single tensor with shape=(1, 6472). The output provides information on many aspects of the driving environment, including traffic car locations, path waypoints over upcoming timesteps, and probabilities associated with certain predictions.

Among state-of-the-art falsifiers, Shriver et al. have developed the framework DNNF [13] that enables the use of many existing adversarial example generators to falsify more sophisticated safety properties. DNNF architecture is shown in Figure 1. The DNNF frontend consumes a safety property written using property specification language DNNP and a network. It outputs a simplified network in an intermediate representation *operation graph* and a robustness property that is equivalent to the safety property. These outputs are then consumed by DNNF, which concretizes inputs that are not being manipulated, according to the property specification. Input that can be manipulated is referred to as *symbolic input*. DNNF then performs falsification, converting the operation graph according to the deep learning (DL) framework utilized by the falsifier, and manipulating the symbolic input to falsify the safety property. In this way, DNNF functions as a framework to allow existing falsification techniques to run on more sophisticated properties and a wider variety of networks regardless of DL framework.

Early falsifiers [8, 17] did not leverage the internal structure of the DNN, instead fuzzing it as a black box. Treating the DNN as a black box makes them applicable to any DNN independent of internal complexity, but they are significantly slower at finding deeper violations because they lack an understanding of the DNN structure.

	Benchmark Networks[13]	Supercombo Network
Total Operators	74	356
Unique Operators	11	15
Unique Activation Functions	2	4
Inputs	1	4
Input Dimensionality	(200,200,1)	(1,12,128,256) (1,8) (1,2) (1,512)
Outputs	1	1
Output Dimensionality	(1,10)	(1,6472)

Table 1: Complexity and scale of DNNF benchmark models versus the industrial supercombo network.

Moreover, though these early falsifiers could be extended to support networks with complex relationships between inputs and/or outputs, they currently do not, in part because of the conventions of falsification benchmark networks. Subsequent falsifiers [3, 5, 15] used adversarial attacks, enabling these techniques to work on any type of network. However, these falsifiers were limited to only local robustness properties, instead of falsifying specific properties written across inputs and outputs using first-order logic.

Applying DNNF to a network like `supercombo` is not possible out-of-the-box due to assumptions about the networks it consumes. Four benchmarks containing 52 networks in total were used to benchmark DNNF. Each benchmark contains trained models categorized by application: ACAS-Xu, CIFAR-EQ, GHPR, and Neurify-DAVE. Table 1 shows the differences between the largest networks used to benchmark DNNF and `supercombo` compared in the ONNX DL framework [9]. The first row compares the maximum number of operators¹ per benchmark network versus `supercombo`, with `supercombo` containing almost 5 times as many. `supercombo` contains 4 more unique operators and 2 more unique activation functions than the most varied networks in the DNNF benchmarks. Significantly, `supercombo` takes 4 inputs whereas all benchmark networks take 1, with the largest dimensions of those inputs being a 1-channel 200×200 image. Although all networks produce one output, the largest output dimensions in the benchmark are only (1,10) compared to `supercombo`'s (1, 6472).

The differences in Table 1 highlight the types of networks for which DNNF was designed. First, multi-input networks are not supported, despite many network architectures requiring multiple inputs. For example, recurrent neural networks are multi-input networks that take in previous output from the network to supplement new input. `supercombo` is itself a recurrent network. Similarly, long-short term memory (LSTM) networks are a specific kind of recurrent network for timeseries data and have appeared in real-world self-driving vehicle software [12]. The falsification step of DNNF presents further challenges for multi-input networks. Falsification relies on manipulating symbolic input in order to falsify properties. So, multi-input networks may need to restrict which inputs are used for falsification, while keeping others constant.

Second, DNNF needs to be able to convert all operators in a network into an intermediate representation. If operator types are not represented internally by DNNF, the network cannot be consumed

¹ONNX operators [9] are a weighted function or matrix operation such as Add, Gemm, or Conv2D, or an activation function. Unlike neurons, operators consider weighted functions or matrix operations and activation functions separately.

for falsification. More recent large-scale networks involve operators that are not supported by the reduction process. For example, `supercombo` involves the `Split` operator that is more common in large networks to direct the flow of information. In total, 20 unique operators are present across all networks in this set of benchmarks. However, the size of the individual networks means that some operators endemic to large-scale networks are not present.

Finally, the scale and complexity of benchmark network outputs stand in contrast to that of `supercombo`. The largest benchmark output dimensionality is (1,10), which comes from ResNet networks' classification score output. Although these scores are correlated, they are low-dimensional in that a limited number of correlations can be inferred between them. In contrast, `supercombo` output consists of far more features with high correlation, such as the predicted path and anticipated acceleration, or locations of lead vehicles and lane line predictions. Falsifying a property for one set of outputs may have unintended consequences on many other outputs. Moreover, many of these individual pieces of information have dimensionality in the x-y-z plane and over time, further complicating the relationships between predicted features. The interconnected and high-dimensional nature of `supercombo` output requires complex properties with low overhead to falsify.

3 ADAPTING DNNF FOR SUPERCOMBO

Given these challenges, we explore extending DNNF to make falsification of `supercombo` viable in order to produce valuable counterexamples of safety properties. We enacted the following workflow to extend DNNF and determine relevant properties for `supercombo`: (1) Extending the DNNF frontend to handle new operators in DL frameworks PyTorch and Tensorflow to support parsing of the input network; and (2) Extending DNNF to support multiple inputs and the specification of a subset of them as constants.

Figure 1 displays the extended portions of DNNF and DNNF frontend in diagonal yellow stripes. The DNNF frontend takes in a network and parses it recursively to create an intermediate representation. This requires all operators in the network to be analyzed, determining the shape of the input and output for each operator, and tracing them throughout the analysis. This is performed in the network parser by a visitor pattern visiting each operator and converting it into a DL framework-agnostic representation to form the operation graph as an intermediate representation of the network.

For `supercombo`, only one operator needed to be added to the DNNF frontend, alongside extensions to existing operators. The `Split` operator splits the input tensor into a set of chunks of parameterized size along a parameterized dimension. The output of this operator is those "chunks" of the original tensor. It is a means to reroute the flow of information in a network. With the introduction of `Split` and the accompanying re-routing of output information from `Split` nodes, the `OutputSelect` operator had to be parameterized to the number of input and output channels and internal logic extended. The `OutputSelect` operator is endemic to the DNNF frontend to assist the operation graph visitor in directing output from a preceding operator with multiple outputs to subsequent operators. While the `OutputSelect` operator was already defined,

it did not have a visitor implemented for any DL framework converter. Besides the implementation, we also added tests to meet the continuous integration requirements of the DNNF framework.

Concretization of inputs allows for a property to be falsified through one symbolic input such as an image, which can be changed by the falsifier, while using other concrete inputs that do not change for the purpose of testing, such as recurrent state. This functionality must be added to DNNF for each falsifier. For multi-input networks, not all inputs may be relevant to falsification. Inputs that are not used for falsification must be concretized in order to be kept constant. Anytime a back-end falsifier analyzes the network, either to determine a gradient or to validate a counterexample, it must do so using a set of inputs. By default, all of those inputs were represented symbolically, because typical networks were only concerned with one input. Code was written for the PGD falsifier to separate symbolic and concrete inputs, determine concrete values, and order them such that they were fed to the model at the correct points. The values corresponding to these inputs are defined via the property file and passed to the model alongside the symbolic input within the appropriate falsifier, which uses them for gradient attachment to the symbolic input and generation of candidate counterexamples.

Understanding, developing, and testing these two extensions to DNNF took approximately 80 hours of engineering time, resulting in 8 commits and approximately 300 lines of code, most for testing the changes. Besides those larger changes, we carried out multiple smaller changes in DNNF. For example, `supercombo`'s size required more memory and a higher parsing recursion limit so the parser could fit the model into memory. Another instance was the need for more sophisticated converters between DL frameworks. For example, ONNX supports named inputs, but since most other frameworks rely on the ordering of inputs, we had to support re-ordering of inputs based on anticipated size when converting from an ONNX framework. Moreover, the default behavior of some existing operators had to be changed. For example, operator `Conv2D` had hardcoded values that assumed less complex operator behavior. Hardcoded parameter `groups=1` affected the behavior of the convolutions with respect to the number of input and output channels, as well as the output shape of this operator.

4 RUNNING DNNF ON SUPERCOMBO

We now explore the effectiveness of our extended version of DNNF v0.1.3 at finding counterexamples for three safety properties in the `supercombo` model of OpenPilot version with commit hash 54d6d9.

4.1 Inputs, Outputs, and Properties

We define safety properties in terms of the inputs and outputs of `supercombo`². The input of interest `input_imgs` is a tensor comprised of two consecutive 256×512 images, collected in 6-channel YUV420 format. `supercombo` output is a tensor of shape=(1, 6472) containing multiple estimates. Among those, we focus on the location of the lane and road edges, the confidence in those estimates, and the probabilities of lead vehicle positions. Figure 2a exemplifies a well-formed image from the OpenPilot dataset we describe in

²A full breakdown of OpenPilot models and input/output is available here: <https://github.com/commaai/openpilot/tree/90af436a121164a51da9fa48d093c29f738ad6a/selfdrive/modeld/models>.

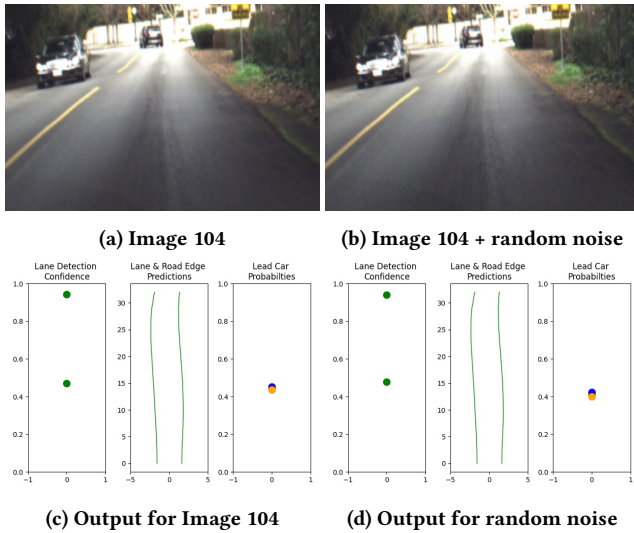


Figure 2: Original and baseline images from the comma.ai dataset and the corresponding output from the supercombo model.

Section 4.2. Three of the outputs of supercombo for this image are shown in Figure 2c: predicted nearest lane lines, the probability for the 2 nearest lane lines, and probabilities for lead car positions identified in the image. Figure 2b shows the same image from Figure 2a with Gaussian noise applied such that the distance from the original image is $\epsilon=10$, a value derived from example properties in the DNNF benchmark. supercombo output for this image (Figure 2d) shows minimal deviation from Figure 2c. A 10-image subset of the dataset with $\epsilon=10$ Gaussian noise applied showed similar results. This suggests random noise is not sufficient to foul supercombo, and that more sophisticated approaches like falsification are needed to cause significant deviations in output.

Given these features of the input and output of the network, we design three properties that can be defined informally as ³:

Safety Property 1: lane line confidence. – “Given an original image x and the lane confidence prediction of that image y , there is no image x' within $\epsilon=10$ from x that renders supercombo lane confidence prediction y' such that y' is more than a 10% margin from y .” A violation in this property could cause OpenPilot to confuse lane lines or disengage following a decrease in confidence.

Safety Property 2: location of nearest lane lines. – “Given an original image x and the lane prediction of that image y , there is no image x' within $\epsilon=10$ for which the last y -value of the near left and near right lane lines in y is farther than a quarter of a standard lane width from the output of x' .” This property determines the overall shape of the lane occupied by the vehicle in the upcoming timesteps. A violation could cause a lane departure.

Safety Property 3: confidence of locations of lead car(s). – “Given an original image x and the lead car confidence prediction of that image being y , there is no image x' within $\epsilon=10$ for which the lead car location confidence in y is farther than a 10% margin from

³The formal properties represented in DNNF, and the code to run them, can be found in the paper repository: <https://github.com/MissMeriel/openpilot-falsification>.

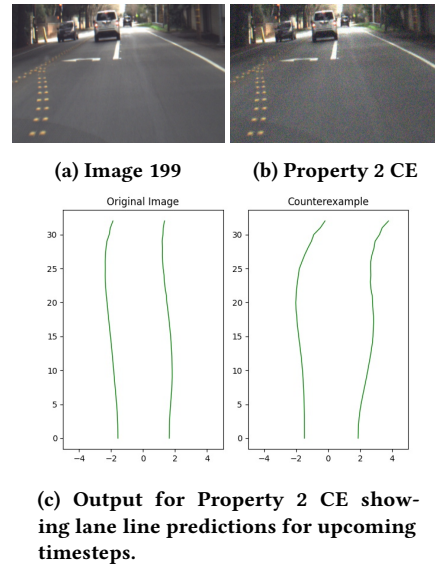


Figure 3: Original dataset image 199 and CE within $\epsilon=10$ and the corresponding output from the supercombo model.

the corresponding lead car location confidence in the output of x' . A violation of this property could cause a collision, a correlated shift in lead car locations, or a disengagement of OpenPilot.

4.2 Finding Counterexamples

Given the specified properties, we run DNNF v0.1.3 on supercombo network using input images from the OpenPilot public dataset [10]. The dataset consists of over 7 hours of driving data and 522,434 images, predominantly on highways.

Figure 3 shows the results of falsification for $\epsilon=10$ over Property 2. These images look relatively similar, in that all of the features of the original image are still discernible to the human eye⁴. Additionally, the counterexample (CE) produced for Property 2, shown in Figure 3b, seems “grainy” but is still human-interpretable.

Property 2 CE output in Figure 3c shows the rightward shift in the last index of the near right lane line pushes it 1.118 meters to the right. This large delta in CE output exceeds the quarter of a lane width specified in the property (0.925 meters) by 0.193 meters. The CE has also pushed the near left lane line slightly closer to the center by a much smaller amount (0.302 meters to the right). The shift in near left and right lanelines creates the appearance of a slight right curve in the road not visible in the original image.

Table 2 shows the performance of falsification for the three properties on a subset of 10 images with well-formed output, averaged over 10 runs for each image. Each run was given 100 random restarts of 100 steps each to find a CE and quits once a CE is found. Column 1 shows nonzero falsification rates for all 3 properties, indicating how many of the 100 runs per property resulted in CEs. The low falsification rate for Property 1, coupled with frequent restarts, suggests that Property 1 is the most difficult to falsify. In column 2, average

⁴Note that, because consecutive images are near-identical, we only show the first image from the two-image composite tensor input to the network. Due to space constraints, we only demonstrate property falsification on image 199 in this paper. Further results can be found in the paper repository.

	Falsification Rate	Avg. Total Time (s)	Avg. Restarts	Images with CEs
Property 1	17%	222.6	88.2	4
Property 2	47%	123.6	60.3	6
Property 3	58%	197.4	46.5	6

Table 2: Average falsification performance for ten dataset images within $\epsilon=10$

total time is the time to reduce the network, simplify the property, and find a CE using falsifier PGD. Average falsification time for each property, respectively 142.8s, 44.6s, and 46.5s, compared to total time shows that network reduction and property simplification are by far the most expensive aspects of falsification. The high frontend time cost suggests avenues for future work in optimization, especially for larger networks and more complex properties. In column 3, average restarts reflect the amount of searching of the input space within $\epsilon=10$ was needed to find a CE, suggesting the properties are ordered in descending difficulty. The confidence of nearest lane line locations appears the most difficult to manipulate, and confidence of lead car positions the easiest. This is supported by falsification rate. Lastly, column 4 shows how many of the 10 images for which a CE was found. Only Image 104 (Figure 2a) produced CEs for all three properties, while 5 images produced CEs for two of the three properties. Nine of the 10 images showed a CE for at least one property. The lower time to run falsification on Property 2 than Property 3, despite the higher number of restarts, can be attributed to the higher number of bounds on output values for Property 3 as Property 3 bounds three values instead of two. Similarly, Property 1 bounds only two output values, and has a comparable average total time to run despite almost double as many restarts as Property 3.

5 LESSONS LEARNED

This study provides several insights into running falsifiers on “real-world” networks. With an engineering cost of tens of hours, we were able to perform falsification on a complex, real-world network and find interesting counterexamples for 3 properties that leverage aspects of the system enclosing the network. Falsification results show counterexamples exist within the estimated input space.

As suggested by Table 1, complex networks have unique architectural considerations, and falsifiers need additional functionality to support them. Adding new operators is low-cost, though operators that complicate the internal data flow of the network may require careful refactoring and testing. For example, the `Split` operator is present in the OpenPilot `dmonitoring` model and is also part of the BigGAN architecture. This suggests that there may be classes of operators that are only used on large-scale networks. Other operators, such as `Slice` and `LSTM`, appear in large and/or complex networks such as `AdmiralNet` and need support or extended functionality in the DNNF framework. `Slice` supports weights passed as functions, requiring special parsing on the part of the DNNF frontend.

This project has inspired several recommendations for future work. Further investigation of large- versus small-scale network considerations might facilitate the development and deployment of falsification tools. Managing network inputs during falsification is a larger effort, as it affects property construction and the functionality of each falsifier differently. Additionally, more complex

networks could benefit from the extension of DNNF to more falsifiers. For supporting most network architectures, we expect that the dominant issue is engineering cost rather than technical barriers.

There are pending challenges regarding input space complexity. Our current solution for input concretization supports properties with one symbolic and multiple concrete inputs, but not multiple symbolic inputs, and may not generalize to multiple concrete inputs of identical size. Inputs with dependencies, like an image and a recurrent state, would need careful property construction and likely require both inputs to be symbolic, which DNNF does not currently support. Furthermore, concrete inputs have the potential to better simplify networks and reduce properties, and may ultimately reduce overhead. Constant propagation and partial input execution for multi-input networks also promise improvement of DNNF.

ACKNOWLEDGMENTS

We are thankful to David Shriver, author of the DNNF toolchain. This work was supported in part by NSF Award #1924777 and AFOSR Award #FA9550-21-1-0164.

REFERENCES

- [1] Arvind Adimoolam, Thao Dang, Alexandre Donzé, James Kapinski, and Xiaoping Jin. 2017. Classification and coverage-based falsification for embedded control systems. In *International Conference on Computer Aided Verification*. Springer, 483–503.
- [2] Cas JF Cremers. 2008. The Scyther Tool: Verification, falsification, and analysis of security protocols. In *International conference on computer aided verification*. Springer, 414–418.
- [3] Tommaso Dreossi, Alexandre Donzé, and Sanjit A Seshia. 2019. Compositional falsification of cyber-physical systems with machine learning components. *Journal of Automated Reasoning* 63, 4 (2019), 1031–1053.
- [4] Daniel J Fremont, Johnathan Chiu, Dragos D Margineantu, Denis Osipchev, and Sanjit A Seshia. 2020. Formal analysis and redesign of a neural network-based aircraft taxiing system with VerifAI. In *International Conference on Computer Aided Verification*. Springer, 122–134.
- [5] Xingwu Guo, Wenjie Wan, Zhaodi Zhang, Min Zhang, Fu Song, and Xuejun Wen. 2021. Eager Falsification for Accelerating Robustness Verification of Deep Neural Networks. In *2021 IEEE 32nd International Symposium on Software Reliability Engineering (ISSRE)*. IEEE, 345–356.
- [6] Greg Hogan. [2021]. Scaling for 10X User Growth. *comma.ai blog* (2021). <https://blog.comma.ai/scaling-for-10x-user-growth/>
- [7] Chris Isidore and Peter Valdes-Dapena. 2021. Tesla is under investigation because its cars keep hitting emergency vehicles. *CNN* (2021). <https://www.cnn.com/2021/08/16/business/tesla-autopilot-federal-safety-probe/index.html>.
- [8] Augustus Odena, Catherine Olsson, David Andersen, and Ian Goodfellow. 2019. Tensorfuzz: Debugging neural networks with coverage-guided fuzzing. In *International Conference on Machine Learning*. PMLR, 4901–4911.
- [9] ONNX. 2019. ONNX: Open Neural Network Exchange. <https://onnx.ai/>.
- [10] OpenPilot. 2022. comma.ai driving dataset. <http://github.com/commaai/research>.
- [11] OpenPilot. 2022. OpenPilot open source driver assistance system, SHA256 tag 54d6d9. <https://github.com/commaai/openpilot>. *OpenPilot* (2022).
- [12] Zi Peng, Jinqiu Yang, Tse-Hsun (Peter) Chen, and Lei Ma. 2020. *A First Look at the Integration of Machine Learning Models in Complex Autonomous Driving Systems: A Case Study on Apollo*. Association for Computing Machinery, New York, NY, USA, 1240–1250. <https://doi.org/10.1145/3368089.3417063>
- [13] David Shriver, Sebastian Elbaum, and Matthew B. Dwyer. 2021. Reducing DNN Properties to Enable Falsification with Adversarial Attacks. In *2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE)*. 275–287. <https://doi.org/10.1109/ICSE43902.2021.00036>
- [14] The Autoware Foundation. 2021. Autoware. <https://www.autoware.org/>.
- [15] Xiao Wang, Saasha Nair, and Matthias Althoff. 2020. Falsification-based robust adversarial reinforcement learning. In *2020 19th IEEE International Conference on Machine Learning and Applications (ICMLA)*. IEEE, 205–212.
- [16] Waymo. 2021. Waymo Driver. <https://waymo.com/waymo-driver/>.
- [17] Xiaofei Xie, Lei Ma, Felix Juefei-Xu, Minhui Xue, Hongxu Chen, Yang Liu, Jianjun Zhao, Bo Li, Jianxiang Yin, and Simon See. 2019. *DeepHunter: A Coverage-Guided Fuzz Testing Framework for Deep Neural Networks*. Association for Computing Machinery, New York, NY, USA, 146–157. <https://doi.org/10.1145/3293882.3330579>