

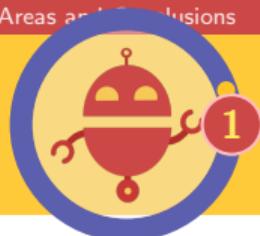
Neural Network Verification With Vehicle: Chapter 5 - Application Areas and Conclusions

VeTSS Summer School'23

Luca Arnaboldi¹ Ekaterina Komendantskaya² Matthew Daggitt (online) ³

¹University of Birmingham · ²University of Southampton · ³Heriot-Watt University

Some More Reasons to Verify AI... (Cars)



stopsign
87% accuracy

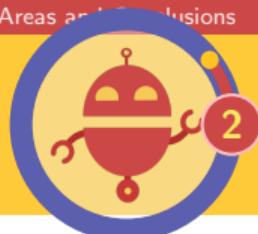
$$+ \quad \square =$$



post-it note
speedlimit 70 miles
99.7% accuracy

Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., ... & Song, D. (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).

EVEN MORE REASONS! (NLP)**



- ▶ Adversarial examples in NLP
 - ▶ Character perturbations
 - ▶ Word perturbations
 - ▶ Sentence perturbations

Are you a robot?

Casadio, M., Arnaboldi, L., Daggitt, M. L., Isac, O., Dinkar, T., Kienitz, D., ... & Komendantskaya, E. (2023). ANTONIO: Towards a Systematic Method of Generating NLP Benchmarks for Verification. arXiv preprint arXiv:2305.04003.

With slide contributions from M. Casadio (Thanks)

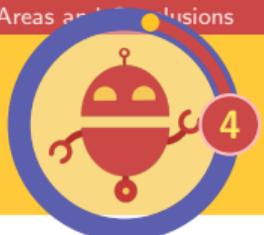


EVEN MORE REASONS! (NLP)

- ▶ Adversarial examples in NLP
 - ▶ Character perturbations
 - ▶ Word perturbations
 - ▶ Sentence perturbations

Are you a robot?
Are you a rpbott?
Are you an robot?

Casadio, M., Arnaboldi, L., Daggitt, M. L., Isac, O., Dinkar, T., Kienitz, D., ... & Komendantskaya, E. (2023). ANTONIO: Towards a Systematic Method of Generating NLP Benchmarks for Verification. arXiv preprint arXiv:2305.04003.



EVEN MORE REASONS! (NLP)

- ▶ Adversarial examples in NLP
 - ▶ Character perturbations
 - ▶ Word perturbations
 - ▶ Sentence perturbations

Are you a robot?

Are you **not** a robot?

Were you a robot?

Casadio, M., Arnaboldi, L., Daggitt, M. L., Isac, O., Dinkar, T., Kienitz, D., ... & Komendantskaya, E. (2023). ANTONIO: Towards a Systematic Method of Generating NLP Benchmarks for Verification. arXiv preprint arXiv:2305.04003.

EVEN MORE REASONS! (NLP)



- ▶ Adversarial examples in NLP
 - ▶ Character perturbations
 - ▶ Word perturbations
 - ▶ Sentence perturbations

Are you a robot?
Am I talking to a robot?
Can u tell me if you are a
chatbot?

Casadio, M., Arnaboldi, L., Daggitt, M. L., Isac, O., Dinkar, T., Kienitz, D., ... & Komendantskaya, E. (2023). ANTONIO: Towards a Systematic Method of Generating NLP Benchmarks for Verification. arXiv preprint arXiv:2305.04003.

Legal Requirement of NLP Verification



People have the right to know if and when they are interacting with a machine's algorithm instead of a human being, the AI Act introduces specific transparency obligations for both users and providers of AI system, such as bot disclosure. Limited Risk AI Systems such as chatbots necessitate specific transparency obligations as well [EU Legislation 2020]



..... Yet another one? (Malware Analysis)

BEFORE

```
1 import android.os.Bundle;
2 import android.view.View;
3 import android.widget.Button;
4 import android.widget.TextView;
5
6 public class MainActivity extends AppCompatActivity
7 {
8
9     private Button button;
10    private TextView .....
11
12    ...
13 }
```

AFTER

```
1 import android.os.Bundle;
2 import android.view.View;
3 import android.widget.Button;
4 import android.widget.TextView;
5 import androidx.appcompat.app.AppCompatActivity;
6 import com.example.randomlibrary1.RandomLibrary1;
7 import com.example.randomlibrary2.RandomLibrary2;
8
9 public class MainActivity extends AppCompatActivity
10 {
11
12     private Button button;
13     private TextView .....
14
15     ...
16 }
```

lines 5 to 7 (AFTER)....

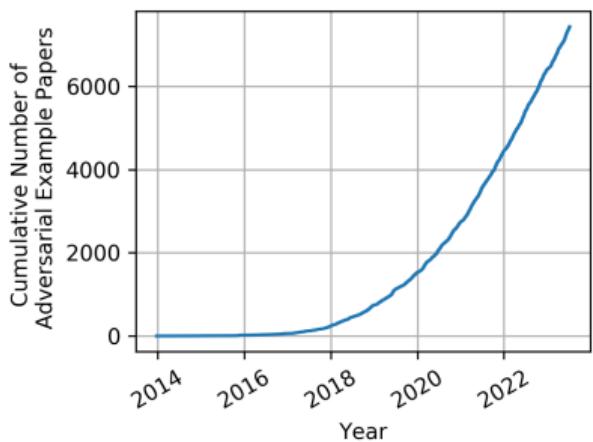
OK - I promise last one! (ML Network IDS)



- ▶ Identification fields: Src IP, Src Port, Dst IP, Dst Port, Protocol, Timestamp
- ▶ Features: Flow Duration, Fwd/Bwd Header Length, (Fwd/Bwd) Packet Length Min/Max/Mean/Std/Total, Total Fwd/Bwd Packets, (Fwd/Bwd) Inter-Arrival Time Min/Max/Mean/Std/Total, (Fwd/Bwd) SYN/FIN/ACK/RST/CWR/PSH/URG/ECE flags count, Packets/second, Bytes/second, Flow Active Duration Min/Max/Mean/Std, Subflow (Fwd/Bwd) Packets/Bytes, Up/Down Ratio
- ▶ Label: FlowType (should be mapped to 0 - BENIGN or 1 - MALICIOUS)
- ▶ **Attacker Objective:** Can packets be manipulated in such a way that the classification switches?

Apruzzese, G., Andreolini, M., Ferretti, L., Marchetti, M., & Colajanni, M. (2022). Modeling realistic adversarial attacks against network intrusion detection systems. *Digital Threats: Research and Practice (DTRAP)*, 3(3), 1-19.

Summary so far

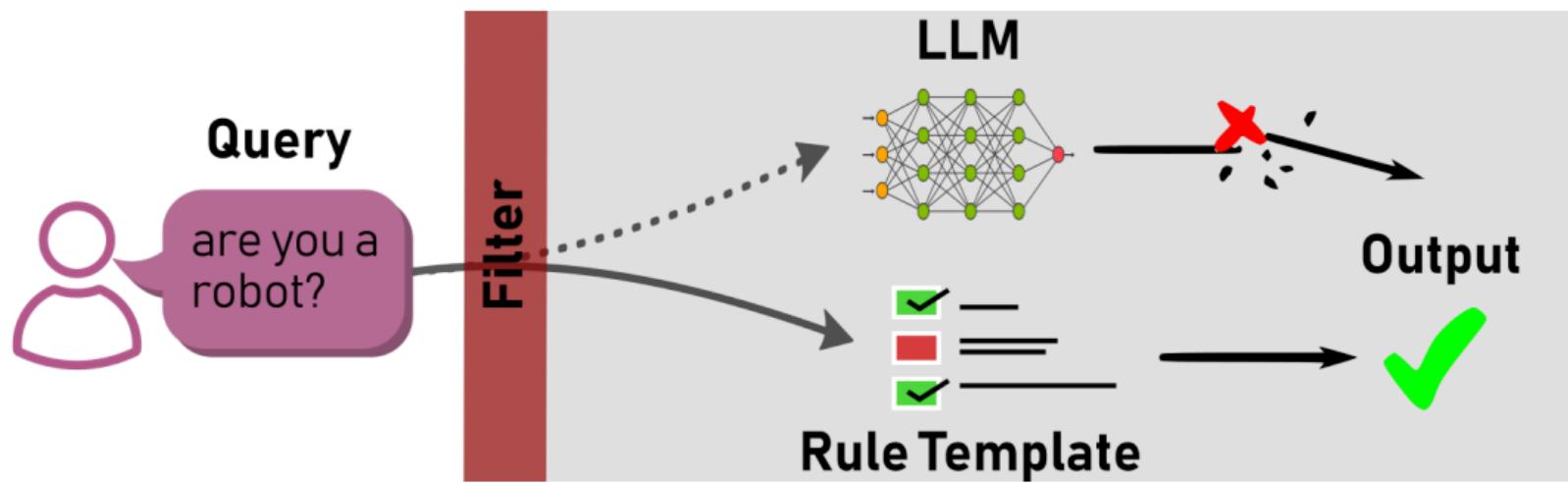


- ▶ Adversarial attacks are here to stay
- ▶ Verification is a promising way to protect against them
- ▶ We have a tool to specify properties and verify them
- ▶ So what are the open problems?
- ▶ **Remember NLP? (Malware, Text, Dialogue etc....)**

NLP Verification - Our approach



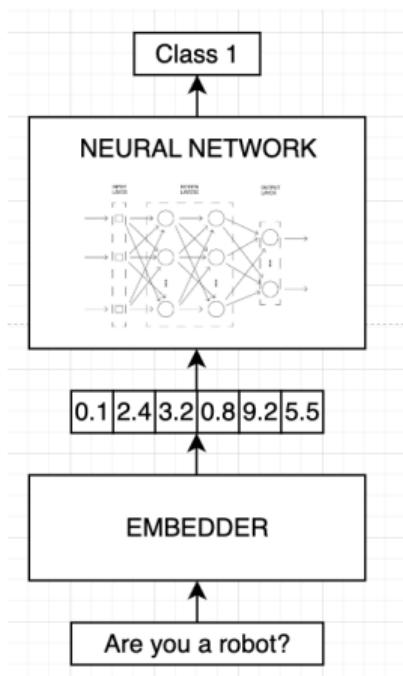
NOT LLMs. (Too big) - Setup a filter network instead



NLP Verification - Our approach



- ▶ Verify the NLP system
- ▶ ϵ -ball
- ▶ Naive approach (ϵ -ball verification)

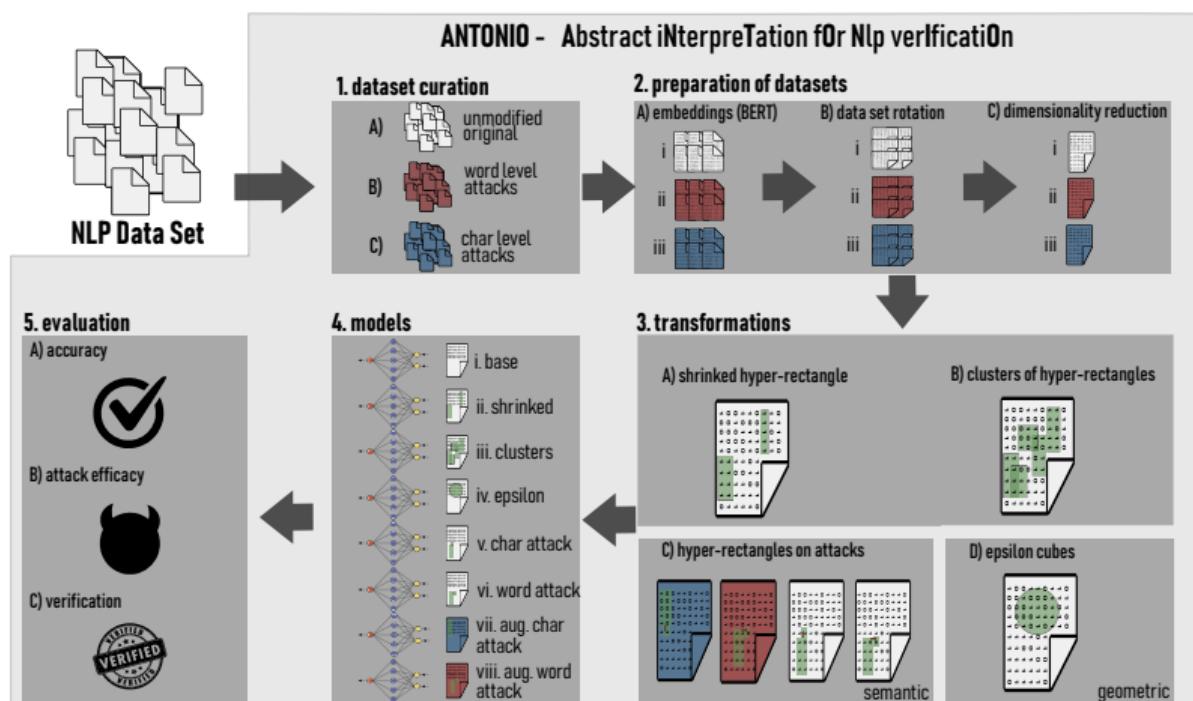


NLP Verification - Obstacles to Verification



1. **Continuous vs discrete space** changes to sentences do not correspond linearly to changes in embedding
2. **Perceptibility by humans.** semantic preservation in NLP vs Perceptibility of Image Space
3. **Difference of the data support.** from range of pixels (RGB 255) to variety of embedders

NLP Verification - ANTONIO



NLP Verification - Results



Model	Test Accuracy	Attack Accuracy	Verification		
			$\mathbb{H}_{\epsilon=0.005}$	$\mathbb{H}_{\epsilon=0.05}$	\mathbb{H}_{pert}
N_{base}	93.87	89.68	88.67	1.79	11.69
N_{adv}	93.38	90.27	98.22	12.17	45.12

Table: Accuracy on test set and attacks and verificaton results using Marabou.

Hyper-rectangles	Avg. Volume	Contained U.S. (%)	Contained U.S. (#)	Total U.S.
$\mathbb{H}_{\epsilon=0.005}$	1.00e-60	1.95	2821	144500
$\mathbb{H}_{\epsilon=0.05}$	1.00e-30	38.47	55592	144500
\mathbb{H}_{pert}	1.28e-30	47.67	68882	144500

Table: Number of unseen sentences inside each collection of hyper-rectangles.

Vehicle Sensor Verification - Reminder

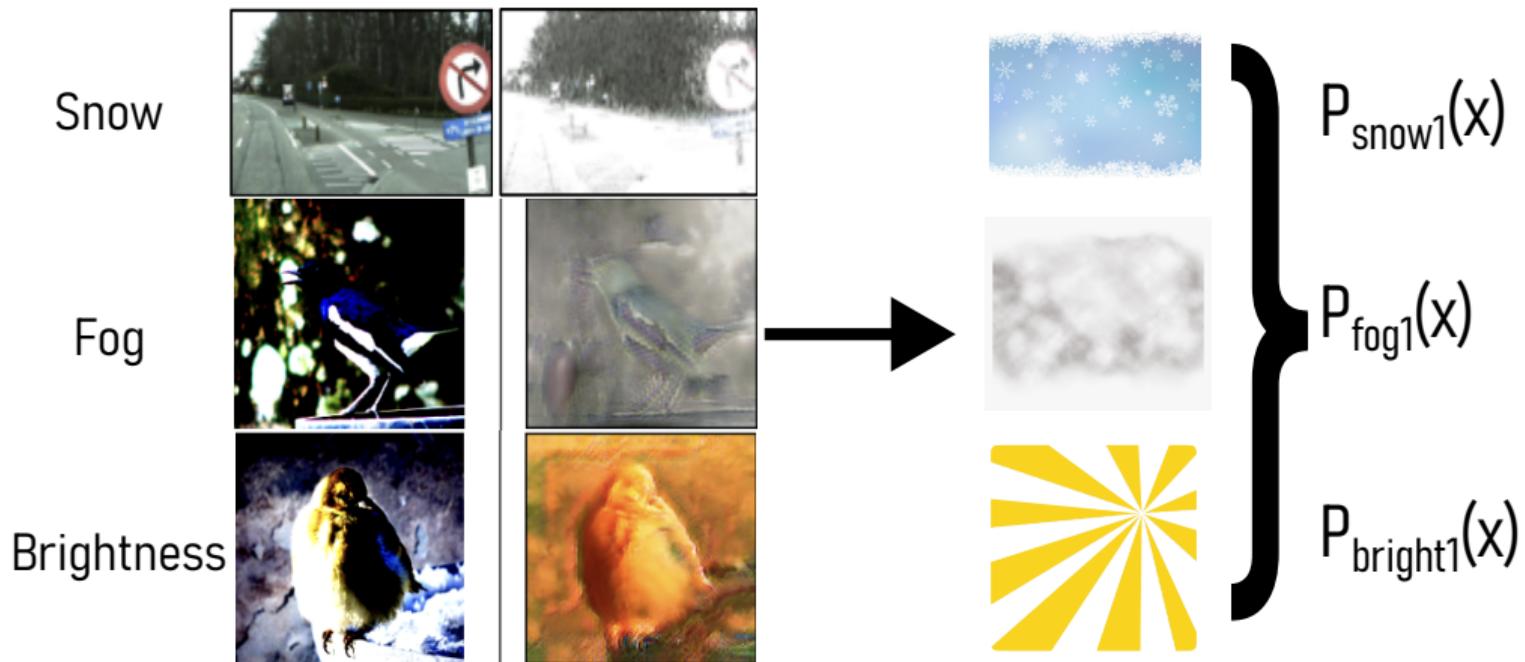


Definition of Verification for a Black Box Model

For a neural network $N : \bar{x} \rightarrow \bar{y}$, the input property $P(\bar{x})$ and the output property $Q(\bar{y})$, does there exist an input \bar{x}_0 which satisfies $P(\bar{x}_0)$ such that its corresponding output \bar{y}_0 satisfies $Q(\bar{y}_0)$?

- ▶ $P(\bar{x})$ characterises inputs checked
- ▶ $Q(\bar{y})$ characterises the behaviour we **DO NOT** wish for
- ▶ if satisfied, counterexample is returned, else property holds
- ▶ the P for traditional adversarial robustness is $\|\bar{x} - \bar{x}_0\|_{L_\infty} \leq \delta$
- ▶ the Q is, $\bigvee_i (\bar{y}[i_0] \leq \bar{y}[i])$, where $\bar{y}[i_0]$ is the desired label

Formal Verification of ML/Sensors - For Resilient Autonomy



Formally Verified IDS Systems



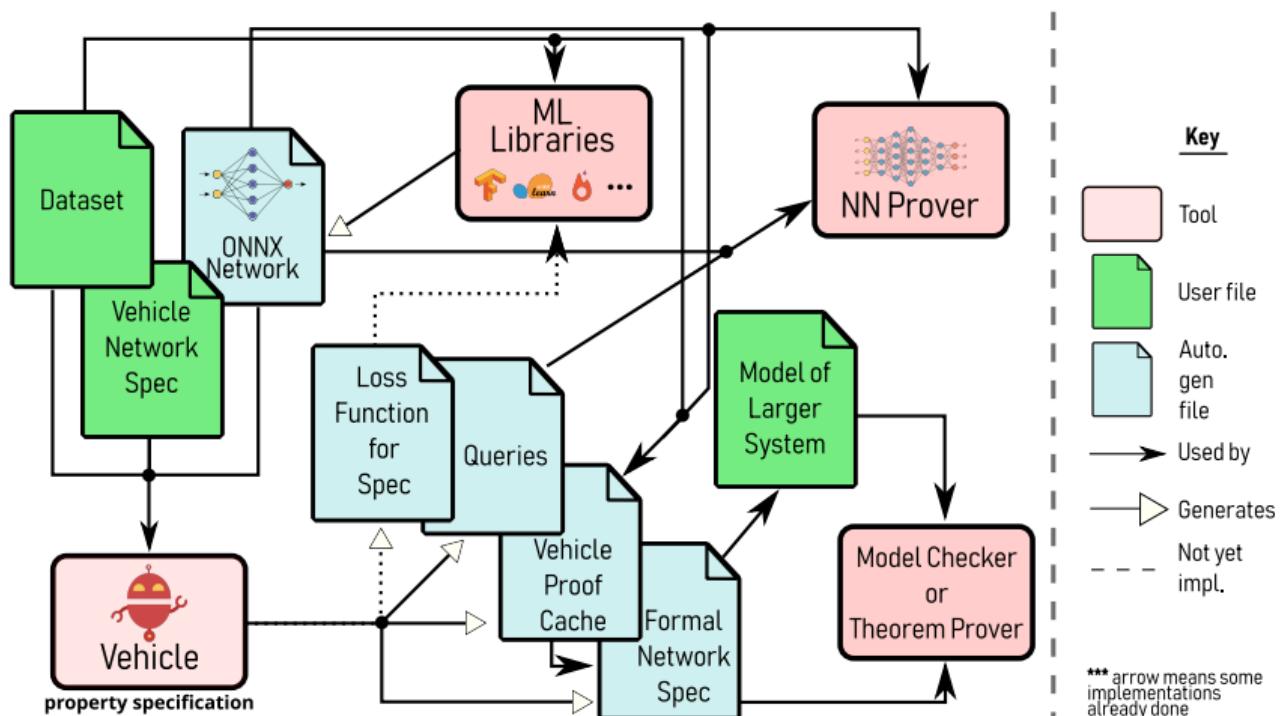
- ▶ Identification fields: Src IP, Src Port, Dst IP, Dst Port, Protocol, Timestamp
- ▶ Features: Flow Duration, Fwd/Bwd Header Length, (Fwd/Bwd) Packet Length Min/Max/Mean/Std/Total, Total Fwd/Bwd Packets, (Fwd/Bwd) Inter-Arrival Time Min/Max/Mean/Std/Total, (Fwd/Bwd) SYN/FIN/ACK/RST/CWR/PSH/URG/ECE flags count, Packets/second, Bytes/second, Flow Active Duration Min/Max/Mean/Std, Subflow (Fwd/Bwd) Packets/Bytes, Up/Down Ratio
- ▶ Label: FlowType (should be mapped to 0 - BENIGN or 1 - MALICIOUS)
- ▶ **Objective:** Given an attacker can perturb these, can we still correctly classify benign and malign traffic?

Vehicle-Tool

One specification, multiple verifications, and more!



18





Conclusions

- ▶ Verification of AI has tons of security case studies to investigate
- ▶ Some upcoming research work from the AISEC team:
 1. Create a detailed mathematical representation of different weather events
 2. Formally Verified ML based Network Intrusion Detection
 3. Continue Down NLP path to include Dialogues
(e.g. Q. Are you a robot? A. No Q2. Are you sure?)
 4. Formal verification of Soundwaves (e.g. Dolphin Attacks)

Thats all folks!