HERIOT-WATT UNIVERSITY

Masters Thesis

Formal Verification of Neural Networks in Go

Author: Supervisor:

Arran Dinsmore Ekaterina Komendanskaya

A thesis submitted in fulfilment of the requirements for the degree of MSc. Robotics

in the

School of Electrical, Electronic & Computer Engineering

School of Engineering & Physical Sciences

&

LAIV: Lab for AI Verification

March 2021



Declaration of Authorship

I, Arran DINSMORE, declare that this thesis titled, 'Formal Verification of Neural Networks in Go' and the work presented in it is my own. I confirm that this work submitted for assessment is my own and is expressed in my own words. Any uses made within it of the works of other authors in any form (e.g., ideas, equations, figures, text, tables, programs) are properly acknowledged at any point of their use. A list of the references employed is included.

Signed: Arran Dinsmore	
Date: March 2021	

"Program testing can be used to show the **presence** of bugs, but never to show their **absence**!"

Edsger W. Dijkstra

Abstract

As machine learning for safety critical applications such as autonomous vehicles are starting to be developed beyond proof of concepts, and enter into production within society, there is a need to ensure these systems do not fail.

Traditional rigorous testing methods are not a viable approach for such black box systems, and thus a need for formal verification methods that can prove the robustness of a system are required.

Additionally, the choice of programming language used for these tasks has grown with new machine learning extensions being developed on existing languages, and newly created languages dedicated for the use within machine learning.

This paper will investigate the feasibility of formal verification for a selection of the most commonly used languages within machine learning tasks, and aim to develop an understanding of generalised approaches and best practices for proving the integrity of systems that cannot afford to fail.

Contents

D	eclara	ation o	of Authorship	i
\mathbf{A}	bstra	\mathbf{ct}		iii
Co	onten	ıts		iv
Li	st of	Figur	es	vi
\mathbf{Li}	st of	Table	S	vii
Li	st of	Abbre	eviations	viii
1	Intr	oducti	ion	1
-	1.1	Conte		1
	1.2	Motiv	ation	3
	1.3		& Objectives	4
2	Bac	kgroui	nd & Literature Review	5
	2.1	Forma	al Verification of AI	5
		2.1.1	Background	5
		2.1.2	Z3	5
			2.1.2.1 Bindings for Z3	5
	0.0	2.1.3	Sapphire	6
	2.2		amming Language Rankings and Metrics	6 6
		2.2.1 $2.2.2$	RedMonk Programming Language Rankings	8
		2.2.2	IEEE Spectrum Ranking of Programming Languages	9
		2.2.4	TIOBE	10
	2.3		iew of Programming Languages	10
			Python	10
		2.3.2	CPP	11
		2.3.3	Matlab	11
		2.3.4	Julia	11
		2.3.5	Go	11
3	Met	hodol	\mathbf{ogy}	12

Contents		v
4	Implementation	13
5	Analysis	14
6	Conclusions	15
A	Appendix Title Here	16
Bi	ibliography	17

List of Figures

1.1	Google's Aversarial Patch	2
2.1	RedMonk Programming Language Rankings for first quarter of 2021	7
2.2	PYPL PopularitY of Programming Language Index	8
2.3	Top 10 Overall IEEE Spectrum Language Ranking	9
2.4	TIOBE Programming Language Index – Long Term History	10

List of Tables

List of Abbreviations

```
AI Artificial Intelligence. 3, 4
AIV Artificial Intelligence Verification. 3–5

MAS Multi-Agent System. 3
ML Machine Learning. 1–5, 10, 11

NN Neural Network. 1–5

PYPL PopularitY of Programming Language Index. 8

SMT Satisfiability Modulo Theories. 4
```

For/Dedicated to/To my...

Introduction

1.1 Context

Machine Learning (ML) algorithms are becoming increasingly present in systems that operate within shared environments with humans, or involve direct interaction with humans themselves [Pereira and Thomas, 2020]. These systems are often defined as safety-critical, such that their failures lead to unintended and potentially harmful behaviours [Amodei et al., 2016]. Examples of these systems include autonomous automotive systems, traffic control systems, medical devices, aviation software, industrial robotics, and many more cyber-physical systems that interact with our environment. Many of these systems have so far only existed as proof of concepts, but are steadily approaching commercial use within our society.

Additionally, recent research has exposed broad vulnerabilities to adversarial attacks within data driven ML algorithms, including Neural Networks (NNs); where applying small but intentional perturbations to an input which are not noticible to humans, can lead to a model outputting an incorrect classification with high confidence. [Goodfellow et al., 2014]. An example of such an attack can be seen in Fig. 1.1. Consequently, the testing and verification of ML for the use of controlling safety-critical systems has become a focused area of research in recent years.

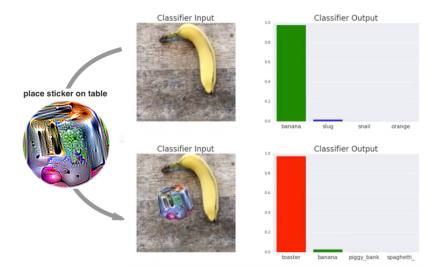


Figure 1.1: Google's Adversarial Patch – An example of a method to create targeted adversarial attacks on NNs by adding carefully designed noise via a physical patch [Brown et al., 2018].

This thesis will use the following definitions for software testing and verification. Software testing, or validation, is defined as the evaluation of a system under various conditions and observing its behaviour while looking for defects [Pereira and Thomas, 2020]. In the context of ML development, testing is used to ensure that a trained model generalises accurately to some previously unseen test data.

Verification is defined as the process of determining whether the products of a phase of the software development process fulfill the requirements established during the previous phase [Ammann and Offutt, 2008]. Formal verification in other words, formulates logical arguments that a system will not act abnormally under a wide range of circumstances, and can be used to determine not only generality, but also the robustness and correctness of a system.

The challenges regarding verification of ML models stem from the typically less deterministic and more statistically-oriented nature of their algorithms, which lead to a lower degree of understanding than software that is explicitly programmed to perform a specific task [Bishop, 2006]. These types of systems are commonly referred to as *black box* systems, where the internal mechanisms are not revealed; in other words, it is impossible to understand a model just by looking at its parameters [Molnar, 2019].

1.2 Motivation

Public calls for *sensible* or *verifiable Artificial Intelligence (AI)* have been raised in recent years due to ever increasing development of complex and pervasive systems that are entering into our everyday lives [Russell et al., 2016].

Formal verification of deterministic software systems has seen significant progress since the early verification systems developed in languages such as Lisp, Ada, and Pascal [Boyer and Moore, 1990, Guaspari et al., 1993, Polak, 1979]. However, verification of non-deterministic systems has seen relatively little progress, with the exception of Multi-Agent Systems (MASs) [Kouvaros and Lomuscio, 2016, Lomuscio et al., 2017].

Indeed, due to the neoteric nature of Artificial Intelligence Verification (AIV) research, there are limited resources with regard to the programming tools available for researchers in this area. This is especially true for work within ML, as the programming languages and tools commonly used for *traditional* verification are often disparate from those widely adopted by the ML communities.

Popular programming languages used for ML such as Python or Matlab currently have comparitively less formal verification tools available than those concerned with system infrastructure or embedded applications. Additionally, AIV toolkits for ML tasks in these languages are still in early stages of development, and mainly focused on the verification of NNs [Kokke, 2020].

Furthermore, the landscape of ML programming itself is forever shifting, and while there is yet a programming language dedicated for ML tasks, huge efforts from programming language designers have been made in developing ML libraries for existing languages. This is necessary in order to handle the extremely high computational demands, and to simplify model languages to make them easier to add domain-specific optimisations and features [Innes et al., 2017].

A prime example of such development can be seen in the Go programming language, or GoLang. A relatively new language, originally developed by Google in 2009 with the intention of creating a modern general-purpose language similar to C. GoLang has seen a surge in popularity within the ML community since the release of its first extensive ML package, Gorgonia, in 2016, which heavily relies on the use of expression graphs [Chew, 2016]. This package allows GoLang developers to take advantage of automatic and symbolic differentiation, gradient descent

optimisations, numerical stabilisation, added support for CUDA/GPGPU computation, and comparitively quick speeds than its Python counterparts (Theano and TensorFlow) [GoLang, 2020].

Consequently, as programming languages continue to develop ML capabilities, there is a need for exploring new and scalable approaches for developing AIV tools in these languages. This is especially important for programming languages which are being adopted by industry to implement ML models for the use within safety-critical or pervasive systems.

1.3 Aims & Objectives

The aim of this project is to investigate the current programming paradigms within ML development, and to explore the suitability of current formal verification toolkits available to them. Subsequently, this thesis will aim to design and implement a GoLang formal methods framework for Gorgonia NNs, providing GoLang ML developers with a set of tools which will allow them to produce safe and fair AI applications.

This framework will extend upon the work made by [Kokke, 2020], and the Sapphire library implemented in Python which successfully translates TensorFlow feed-forward NN models to the Z3 Satisfiability Modulo Theories (SMT) solver created by Microsoft Research [De Moura and Bjørner, 2008].

To achieve this project's aims, the following objectives should be met:

- Objective 1 Implement bindings that map the parameters of a Gorgonia NN model to Z3 variables.
- Objective 2 Select data in order to train and verify NN models using this project's formal methods framework.
- Objective 3 Implement a series of NN models in Gorgonia using the data sets mentioned in Objective 2.
- Objective 4 Verify the correctness of Gorgonia NNs using the bindings mentioned in Objective 1.

Background & Literature Review

This chapter will provide a background understanding to the important concepts that are required by this thesis, and explore the current trends within AIV research. This includes an introduction to formal verification tools and their use for verifying NNs, along with a discussion about the importance of developing AIV tools across a wide range of programming languages.

A high-level overview of programming languages that are used for ML tasks, as well as a discussion of their popularity and usage within industry and academia will conclude this chapter.

2.1 Formal Verification of AI

- 2.1.1 Background
- 2.1.2 Z3
- 2.1.2.1 Bindings for **Z**3

Go-Z3 etc.

2.1.3 Sapphire

2.2 Programming Language Rankings and Metrics

The following section outlines the various programming language ranking systems and metrics used to determine the popularity of, and overall usage of programming languages. A combination of the following rankings will be used to assess the current trends for each programming language discussed in this paper.

2.2.1 RedMonk Programming Language Rankings

RedMonk is a developer-focused industry analyst firm that curates a quarterly ranking of programming languages. The rankings are created by looking at a programming language's presence on GitHub and Stack Overflow, attempting to reflect both the usage of the language from GitHub, and the amount of discussion regarding a language from Stack Overflow [O'Grady, 2021].

This ranking uses a simple metric that can be used for gaining an overview of a programming language's popularity within industry. However, due to quantifying discussions from Stack Overflow, older languages with a small set of built-in library functions such as C can be at a disadvantage compared to newer languages and those with a large set of built-in library functions [Benson, 2017].

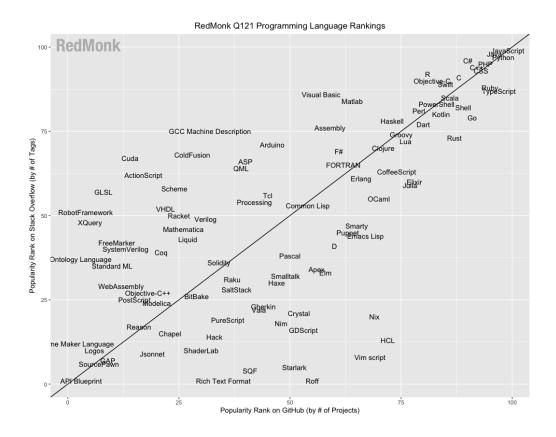


FIGURE 2.1: **RedMonk Rankings Q121** – Comparing the presence of programming languages from GitHub and Stack Overflow [O'Grady, 2021].

Fig. 2.1 shows an example of the current visualisation tools available from this index.

2.2.2 PYPL PopularitY of Programming Language Index

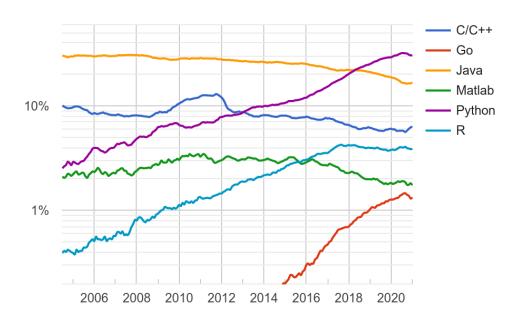


FIGURE 2.2: **PYPL Rankings** – An overview of the share of language tutorial Google searches over the last 15 years [Carbonelle, 2021]. *This chart uses a logarithmic scale.*

According to the PopularitY of Programming Language Index (PYPL), which uses the amount in which a language tutorial is searched on Google as a metric for popularity, Python has been the most popular language for the last few years (see Fig. 2.2), currently sitting on 30.44% of tutorial searches [Carbonelle, 2021].

Several reasons could be contributing to this popularity, including its versatility across a wide range of tasks, being platform agnostic, and easy to learn for beginners due to a large community and *readable* syntax. However, PYPL does not look at either the performance of languages or their popularity within specific tasks such as machine learning.

2.2.3 IEEE Spectrum Ranking of Programming Languages

Rank	Language	Туре			Score
1	Python▼	#	Ģ	0	100.0
2	Java▼	#	Ç		95.3
3	C▼		Ç	0	94.6
4	C++▼		Ç	@	87.0
5	JavaScript▼	#			79.5
6	R▼		Ç		78.6
7	Arduino▼			0	73.2
8	Go▼	#	Ç		73.1
9	Swift▼		Ç		70.5
10	Matlab▼		Ç		68.4

FIGURE 2.3: **Top 10 Overall IEEE Spectrum Language Ranking** – Default metrics ranking top 10 programming languages [Diakopoulous et al., 2020].

2.2.4 TIOBE

Programming Language	2021	2016	2011	2006	2001	1996	1991	1986
С	1	2	2	2	1	1	1	1
Java	2	1	1	1	3	28	-	-
Python	3	5	6	7	23	16	-	-
C++	4	3	3	3	2	2	2	8
C#	5	4	5	6	9	-	-	-
JavaScript	6	7	9	9	6	30		-
PHP	7	6	4	4	20	-	-	-
R	8	14	35	-	-	-	-	-
SQL	9		-	-				-
Go	10	56	15	-	-	-	-	-
Perl	14	8	7	5	4	3		-
Lisp	32	23	12	13	16	7	3	2
Ada	34	22	20	15	15	5	9	3

Figure 2.4: **TIOBE Very Long Term History Programming Language**Index – Long term history of programming languages, average positions for a period of 12 months [Jansen, 2021].

2.3 Overview of Programming Languages

A plethora of programming languages have been developed over time, some dedicated and others extended for the purpose of ML. This section will look at some of the most commonly used languages within the field of ML, and provide an overview of their strengths and weaknesses.

2.3.1 Python

Python is an interpreted, object-oriented, high-level, dynamically typed programming language with dynamic semantics. It was initially designed in 1991 by Guido Van Rossum and subsequently developed by Python Software Foundation. Python's simple syntax emphasising readability was the main purpose of its creation, making it very attractive for Rapid Application Development and reducing the cost of program maintenance [Python, 2021].

use in machine learning

Many versions of Python have since been released, and over time has seen the development of an extensive collection of community libraries used for a wide range of tasks. This has made Python one of the most versatile languages available today. Amongst these tasks, Python has become one of the most popular languages for ML and data science.

strengths

weaknesses

- 2.3.2 CPP
- **2.3.3** Matlab
- 2.3.4 Julia
- 2.3.5 Go

According to PYPL it has had the greatest increase in popularity since its release in 2015.

A relatively new language compared to the others used in machine learning.

Still in development, especially wrt machine learning tasks.

fast, concurrent, used a lot in system architecture and web applications.

Gorgonia

Gorgo

Go-ML

Methodology

Implementation

Analysis

Conclusions

Appendix A

Appendix Title Here

Write your Appendix content here.

Bibliography

- Ammann, P. and Offutt, A. J. (2008). Introduction to software testing.
- Amodei, D., Olah, C., Steinhardt, J., Christiano, P. F., Schulman, J., and Mané, D. (2016). Concrete problems in AI safety. CoRR, abs/1606.06565.
- Benson, D. (2017). 10 most popular programming languages open source for you. https://www.opensourceforu.com/2017/03/most-popular-programming-languages/. (Accessed on 03/09/2021).
- Bishop, C. (2006). Pattern Recognition and Machine Learning, volume 16, pages 140–155.
- Boyer, R. S. and Moore, J. S. (1990). A theorem prover for a computational logic. In Stickel, M. E., editor, 10th International Conference on Automated Deduction, pages 1–15, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Brown, T. B., Mané, D., Roy, A., Abadi, M., and Gilmer, J. (2018). Adversarial patch.
- Carbonelle, P. (2021). Pypl popularity of programming language index.
- Chew, X. (2016). Gorgonia.
- De Moura, L. and Bjørner, N. (2008). Z3: An efficient smt solver. In Proceedings of the Theory and Practice of Software, 14th International Conference on Tools and Algorithms for the Construction and Analysis of Systems, TACAS'08/ETAPS'08, page 337–340, Berlin, Heidelberg. Springer-Verlag.
- Diakopoulous, N., Bagavandas, M., Singh, G., and Kulka-(2020).Interactive: top programming languages rni, ieee spectrum. https://spectrum.ieee.org/static/ interactive-the-top-programming-languages-2020. (Accessed on 03/09/2021).
- GoLang (2020). Golang machine learning libraries.

Bibliography 18

Goodfellow, I., Shlens, J., and Szegedy, C. (2014). Explaining and harnessing adversarial examples. arXiv 1412.6572.

- Guaspari, D., Marceau, C., and Polak, W. (1993). Formal verification of ada programs. In Martin, U. and Wing, J. M., editors, *First International Workshop on Larch*, pages 104–141, London. Springer London.
- Innes, M., Barber, D., Besard, T., Bradbury, J., Churavy, V., Danisch, S., Edelman, A., Karpinski, S., Malmaud, J., Revels, J., and et al. (2017). On machine learning and programming languages.
- Jansen, P. (2021). index tiobe the software quality company. https://www.tiobe.com/tiobe-index/?20210304%20%20%20%20%20%20. (Accessed on 03/09/2021).
- Kokke, W. (2020). A library for translating tensorflow models to z3.
- Kouvaros, P. and Lomuscio, A. (2016). Parameterised verification for multi-agent systems. *Artif. Intell.*, 234(C):152–189.
- Lomuscio, A., Qu, H., and Raimondi, F. (2017). Mcmas: an open-source model checker for the verification of multi-agent systems.
- Molnar, C. (2019). *Interpretable Machine Learning*. https://christophm.github.io/interpretable-ml-book/.
- O'Grady, S. (2021). The redmonk programming language rankings: January 2021.
- Pereira, A. and Thomas, C. (2020). Challenges of machine learning applied to safety-critical cyber-physical systems. *Machine Learning and Knowledge Extraction*, 2(4):579–602.
- Polak, W. (1979). An exercise in automatic program verification. *IEEE Transactions on Software Engineering*, (5):453–458.
- Python (2021). What is python? executive summary python.org. https://www.python.org/doc/essays/blurb/. (Accessed on 03/10/2021).
- Russell, S., Dewey, D., and Tegmark, M. (2016). Research priorities for robust and beneficial artificial intelligence.