

# Vehicle: Interfacing Neural Network Verifiers with Interactive Theorem Provers

Matthew L. Daggitt<sup>1</sup> ✉ 🏠 📧

Department of Computer Science, Heriot-Watt University, Edinburgh, UK

Wen Kokke ✉ 📧

Mathematically Structured Programming Group, University of Strathclyde, Glasgow, UK

Robert Atkey ✉ 📧

Mathematically Structured Programming Group, University of Strathclyde, Glasgow, UK

Luca Arnaboldi ✉ 📧

Laboratory for Foundations of Computer Science, University of Edinburgh, Edinburgh, UK

Ekaterina Komendantskya ✉ 📧

Department of Computer Science, Heriot-Watt University, Edinburgh, UK

## Abstract

Verification of neural networks is currently a hot topic in automated theorem proving. Progress has been rapid and there are now a wide range of tools available that can verify properties of networks with hundreds of thousands of nodes. In theory this opens the door to the verification of larger control systems that make use of neural network components. However, although work has managed to incorporate the results of these verifiers to prove larger properties of individual systems, there is currently no general methodology for bridging the gap between verifiers and interactive theorem provers (ITPs).

In this paper we present Vehicle, our solution to this problem. Vehicle is equipped with an expressive domain specific language for stating neural network specifications which can be compiled to both verifiers and ITPs. It overcomes previous issues with maintainability and scalability in similar ITP formalisations by using a standard ONNX file as the single canonical representation of the network. We demonstrate its utility by using it to connect the neural network verifier Marabou to Agda and then formally verifying that a car steered by a neural network never leaves the road, even in the face of an unpredictable cross wind and imperfect sensors. The network has over 20,000 nodes, and therefore this proof represents an improvement of 3 orders of magnitude over prior proofs about neural network enhanced systems in ITPs.

**2012 ACM Subject Classification** Software and its engineering → Software verification; Computing methodologies → Machine learning

**Keywords and phrases** Neural networks, Verification, Interactive Theorem Provers, Agda, Marabou

**Digital Object Identifier** 10.4230/LIPIcs.ITP.2022.23

**Funding** This work was funded by the AISEC grant under EPSRC numbers EP/T026952/1, EP/T026960/1, and EP/T027037/1.

**Acknowledgements** We would like to thank the Marabou development team for their support and advice with integrating Vehicle with Marabou.

## 1 Introduction

In the last decade deep neural networks have made their way into systems used in everyday life and, as with any system that interacts directly with humans, it is highly desirable to

<sup>1</sup> Corresponding author

41 have formal guarantees about their behaviour. However, these systems present a challenge  
42 for the verification community as typically the neural networks are used in domains where a  
43 formal specification of the desired behaviour remains elusive. Furthermore, their size and  
44 inability to be decomposed into components that tackle identifiable sub-tasks mean that they  
45 are usually viewed as black-box components, which makes traditional verification difficult to  
46 apply.

47 Nonetheless, the discovery of adversarial attacks on neural networks in 2014 [27], has  
48 spurred the Automated Theorem Proving (ATP) community to develop neural network  
49 verifiers. These are based on SMT solvers [9, 19], abstract interpretation [26] or interval  
50 bounded arithmetic [29] and aim to prove linear (or occasionally semi-definite) relationships  
51 between the inputs and the outputs of the network. Progress has been rapid and they are  
52 now reaching the point where they are powerful enough to verify properties of networks with  
53 tens or even hundreds of thousands of nodes.

54 The majority of the ATP community's attention has been focused on using them to prove  
55 the robustness of the networks against adversarial attack, and there has been comparatively  
56 little work on using them to prove the functional correctness of larger systems that incorporate  
57 neural network components. Although there are several reasons for this, including the difficulty  
58 in coming up with a specification for the neural network in the first place, we believe that a  
59 key missing component on the technical side is that no one has yet integrated them with  
60 interactive theorem provers (ITPs).

61 As when incorporating other ATP tools such as SMT solvers into ITPs [4, 10], one must  
62 translate high-level statements in the ITP into low level queries for the verifier. However,  
63 there are several additional challenges unique to the integration of neural network verifiers:

- 64 **1. Modelling mismatch** - in an ITP a neural network is usually modelled as a function  
65 which can be composed with a larger system. However, neural network verifiers model a  
66 network as a relation between its inputs and outputs. As discussed in Sections 2.2 & 4.2.4,  
67 the translation from the former to the latter is not straightforward.
- 68 **2. Scalability** - modern networks can contain millions or even billions of nodes, whereas  
69 most ITPs will consume excessive amounts of memory when representing even very  
70 small networks. For example, recent formalisations in Coq [3, 8] have worked with  
71 networks of just 10 or 20 nodes.
- 72 **3. Maintainability** - most neural networks are not static artefacts, and regularly undergo  
73 further training as new data is collected. For obvious practical reasons, a formally verified  
74 representation of the network within an ITP is unlikely to be the canonical representation  
75 deployed in a production system. Instead, specialised file formats are used to distribute  
76 and deploy networks. This raises the problem of how to maintain the faithfulness of the  
77 model (and the proofs) within the ITP with the rapidly evolving implementation stored  
78 elsewhere.
- 79 **4. Performance** - even with domain specific verifiers, verification of large neural networks  
80 can be exceedingly expensive, often taking hours or even days to complete. This has the  
81 potential to be very disruptive in ITPs whose workflow encourages users to regularly  
82 recheck the validity of their proof during development.
- 83 **5. Integration into other parts of the neural network lifecycle** - verification is only a  
84 small part of constructing a neural network with formal guarantees. When trained using  
85 traditional methods, a network is highly unlikely to satisfy the required specifications.  
86 Recent work has shown how specifications may be integrated into the training of the  
87 network [11] and gradient-based counter-example search [21]. Therefore ideally the  
88 specification should be usable in these other tools as well. However, current ITPs are

89 a poor environment in which to perform these complex and computationally intensive  
90 operations.

91 We believe that in order for verification of neural network-enhanced systems to achieve  
92 wide-spread adoption, all of these issues must be addressed. Our vision is as follows.  
93 The specification for a neural network should be stated once, in a high-level, human-  
94 readable format. This specification should then be automatically translated to work with  
95 the appropriate tools in the different stages in the life-cycle of the network. Importantly,  
96 there should only ever be *one* representation of the network, stored in a format usable by  
97 the mainstream machine learning community. For performance and modularity reasons,  
98 we also argue that an appropriate level of abstraction should be maintained at each stage.  
99 Concretely, the training and verification tools should be able to inspect the internal neural  
100 network structure. However, when writing the specification and using it to prove properties  
101 of the larger system in an ITP, the view of the network as a black-box function should be  
102 maintained where possible. Finally to maintain interactivity in the ITP, the checking of the  
103 proof of system correctness should be decoupled from the checking of the proof of the neural  
104 network specification, so that former does not automatically trigger the latter.

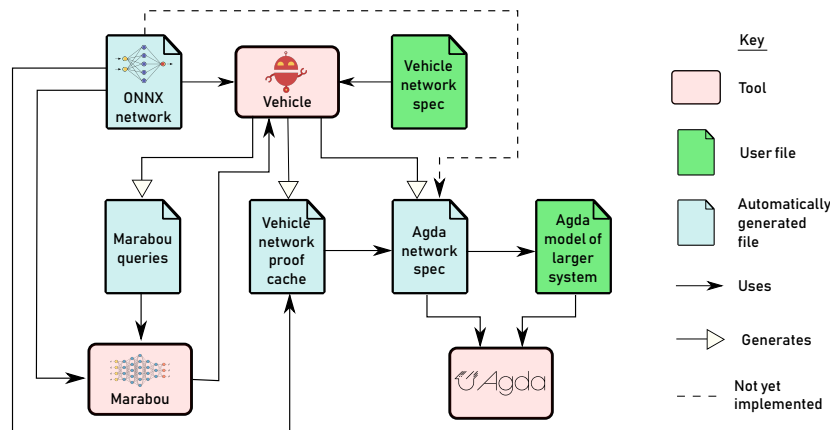
## 105 1.1 Contributions

106 In this paper we present *Vehicle*, a tool that implements the interaction between network  
107 specifications, verifiers and ITPs proposed in our vision above. In particular:

- 108 1. *Vehicle* is equipped with a high-level, domain specific language (DSL) for writing neural  
109 network specifications. The *Vehicle* compiler automatically translates these specifications  
110 down to low-level queries for neural network verifiers, and then subsequently to high-level  
111 ITP code. The latter can be used as an interface upon which proofs about the larger  
112 systems can be constructed.
- 113 2. Instead of modelling the network directly inside the ITP, *Vehicle* side-steps the maintain-  
114 ability issue outlined above by using an Open Neural Network Exchange (ONNX) [1] file  
115 as the single canonical representation of the network. The ONNX format is a widely-  
116 supported, cross-platform, training-framework independent, representation of neural  
117 networks [1]. The *Vehicle* compiler extracts the necessary internal implementation details  
118 from the ONNX file, and uses hashing to maintain the integrity of the specification in  
119 the ITP with respect to the underlying ONNX file.
- 120 3. As *Vehicle* stores the neural network externally and uses specialised neural network  
121 verifiers, its performance is dependent on that of the underlying verifier rather than the  
122 ITP. This means that *Vehicle* can potentially be used to verify systems that use networks  
123 with hundreds of thousands of nodes.
- 124 4. *Vehicle* maintains interactivity in the ITP when checking the proof, as the generated  
125 interface code calls back to *Vehicle* rather than directly calling the neural network verifier.  
126 *Vehicle* then checks its proof cache to ascertain the verification status of the specification,  
127 thereby preventing costly and unnecessary re-verification of the network.
- 128 5. *Vehicle* maintains the abstraction of the neural network as a black-box function. When  
129 writing the specification the user is only required to provide the type of the function  
130 implemented by the network. Similarly in the ITP interface, the neural network is  
131 presented as a function that can be used, but not decomposed.

132 In our current implementation, we target the SMT-based neural network verifier Marabou [19]  
133 and the interactive theorem prover Agda [23]. The latter was chosen due to the authors'

## 23:4 Vehicle: NNs to ITPs



■ **Figure 1** The architecture of a Vehicle proof about a neural network enhanced system. The ability to evaluate the network directly in Agda has not yet been implemented. However, the integrity of the proof with respect to the network is maintained via the proof cache file, which also maintains the interactivity of the Agda file.

134 expertise in it, and we acknowledge it is better suited to modelling systems rather than ex-  
 135 tracting formally verified executable code. However, the integration with Agda is deliberately  
 136 very lightweight, with a single file in the compiler defining the translation, and a single Agda  
 137 file defining macros for calling back to Vehicle. It would therefore be relatively simple to  
 138 extend support to other ITPs such as Coq.

139 As an example of the utility of Vehicle, we use it to formally verify that a car controlled  
 140 by a neural network will never leave the road, even in the presence of noisy sensor data  
 141 and an unpredictable cross-wind. The neural network has over 20,000 nodes and therefore,  
 142 as far as we aware, this proof represents an improvement of over 3 orders of magnitude  
 143 when compared to previous proofs about neural network enhanced systems in ITPs. The  
 144 architecture of the proof is shown in Figure 1. All accompanying code, including Vehicle  
 145 itself, is available online [7].

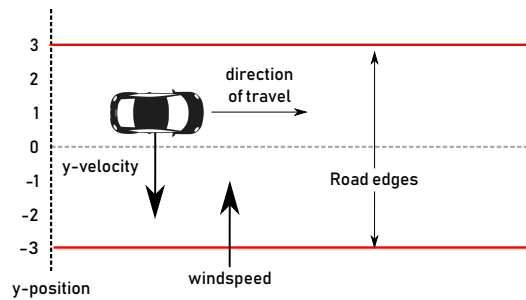
146 Vehicle will also be of use to people who are not interested in integrating with ITPs. The  
 147 existing interfaces to verifiers are very low-level, usually involving the user specifying a set  
 148 of equalities or inequalities over the individual inputs and outputs of the neural network.  
 149 As a large neural network typically can have thousands of inputs and outputs, creating  
 150 such inequalities is both time consuming and error prone. Furthermore, the result is almost  
 151 completely unintelligible to a non-technical domain expert. In contrast Vehicle can provide a  
 152 much higher-level, human readable version of the specification.

153 The paper is laid out as follows. In Section 2 we describe our running example problem  
 154 and related work on verifiers, neural network verification efforts and other similar specification  
 155 languages. In Section 3 we propose an extension to the query language for Marabou to allow  
 156 it to support specifications involving multiple networks. In Section 4, we describe the Vehicle  
 157 DSL and the novel passes in the Vehicle compiler used to generate verifier queries and the  
 158 Agda code. Finally, in Section 5 we discuss future work.

**2 Background**

**2.1 An example: staying on the road**

We will use a modified version of the verification problem presented by Boyer, Green and Moore [5] as a running example throughout this paper. In the scenario an autonomous vehicle is travelling along a straight road of width 6 parallel to the x-axis, with a varying cross-wind that blows perpendicular to the x-axis. The vehicle has an imperfect sensor that it can use to get a (possibly noisy) reading on its position on the y-axis, and can change its velocity on the y-axis in response. The car’s controller takes in both the current sensor reading and the previous sensor reading and its goal is to keep the car on the road. The setup is illustrated in Figure 2.



**Figure 2** A simple model of an autonomous vehicle compensating against a cross-wind.

For simplicity, we assume that both the wind-speed and the car’s velocity in the y-direction can grow arbitrarily large. As in [5] we discretise the model, and then formalise it in Agda as follows. The state of the system consists of the current wind speed, the position and velocity of the car and the most recent sensor reading. An oracle provides updates in the form of observations consisting of the shift in wind speed and the error on the sensor reading.

```

174 record State : Set where
175   constructor state
176   field
177     windSpeed : ℚ
178     position  : ℚ
179     velocity  : ℚ
180     sensor    : ℚ
181
182 record Observation : Set where
183   constructor observe
184   field
185     windShift : ℚ
186     sensorError : ℚ
187

```

For the moment, we assume that we have some controller that takes in the current and the previous sensor reading and produces a resulting change in velocity:

```

190 controller : ℚ → ℚ → ℚ

```

Given this, we can define the evolution of the system as follows:

```

192 initialState : State
193 initialState = state 0ℚ 0ℚ 0ℚ 0ℚ

```

## 23:6 Vehicle: NNs to ITPs

```

195   nextState : Observation → State → State
196   nextState o s = state newWindSpeed newPosition newVelocity newSensor
197   where
198     newWindSpeed = windSpeed s + windShift o
199     newPosition = position s + velocity s + newWindSpeed
200     newSensor   = newPosition + sensorError o
201     newVelocity = velocity s + controller newSensor (sensor s)
202
203   finalState : List Observation → State
204   finalState xs = foldr nextState initialState xs

```

205 Given this setup we would like to prove the following property of the system:

206 ► **Theorem 1.** *Assuming that the wind-speed can shift by no more than 1 per unit time and*  
207 *that the sensor is never off by more than 0.25 then the car will never leave the road.*

208 We define the pre-conditions of the theorem in Agda as follows:

```

209   ValidObservation : Observation → Set
210   ValidObservation o = | sensorError o | ≤ 0.25ℚ × | windShift o | ≤ 1ℚ

```

211 and the post-condition as:

```

212   OnRoad : State → Set
213   OnRoad s = | position s | ≤ 3ℚ

```

214 which allows us to formalise the theorem as:

```

215   finalState-onRoad : ∀ xs → All ValidObservation xs → OnRoad (finalState xs)

```

216 As is standard when proving properties of large systems in a top-down manner, one eventually  
217 deduces the properties that must hold of the sub-components. In this case Theorem 1 can  
218 be proved by formulating a suitable inductive hypothesis about the state of the system at  
219 each time-step. The full inductive proof can be found online [], but the crucial part is that  
220 the inductive step requires the `controller` function to satisfy the following specification:

```

221   controller-lemma : ∀ x y → | x | ≤ 3.25ℚ → | y | ≤ 3.25ℚ → | controller x y + 2ℚ * x - y | < 1.25ℚ

```

222 We implement the controller with a 3-layer densely connected neural network with over  
223 20,000 nodes. The network is constructed in Tensorflow and all weights in the network are  
224 unique and non-zero. The challenge is then to implement the Agda function `controller` with  
225 this neural network, and prove `controller-lemma`, without having to directly represent either  
226 the network or the proof directly in Agda.

227 In the interests of full disclosure, the problem above as stated is hardly very challenging  
228 to solve with a neural network and the size of the network used is complete overkill. The  
229 scenario could be made more realistic by adding further, possibly conflicting, objectives  
230 that the controller should optimise for, thereby rendering the problem intractable to solve  
231 analytically. For example, adding waypoints on the road that the car had to pass through,  
232 or regions it had to avoid. In such scenarios the statement of Theorem 1 would remain the  
233 same, but the number of inputs to the controller would increase. We do not explore these  
234 more complicated scenarios in this paper due to space limitations. Nonetheless, the simple  
235 scenario presented above is sufficient to illustrate the methodology and the utility of Vehicle.

## 2.2 Neural network verifiers

Neural network verifiers can be roughly split into two families: those that are both sound and complete, like Marabou [19], and those that are only sound, based on techniques such as abstract interpretation [26] or bounded interval arithmetic [29]. While the latter are more performant, their incompleteness means that when the verifier fails to show that the property holds, it is unknown whether there truly exists a counter-example.

There is no consensus on a single input format for verifications queries. Each solver defines its own, which makes interfacing with multiple solvers difficult. However, one commonality is that they implicitly model the neural network as a *relation* between its inputs and outputs, where each of the network's inputs and outputs. For example when writing a Marabou query for a network with  $m$  inputs and  $n$  outputs, the input are labelled  $x_0, x_1, \dots, x_{[m-1]}$  and the outputs are labelled  $y_0, y_1, \dots, y_{[n-1]}$ . Queries are then written as a series of inequalities involving these variables. For example, the Marabou queries required to verify [controller-lemma](#) are shown in Figure 3.

```
x0 >= -3.25
x0 <= 3.25
x1 >= -3.25
x1 <= 3.25
-y0 -2.0x0 +x1 >= 1.25
```

(a) Query 1

```
x0 >= -3.25
x0 <= 3.25
x1 >= -3.25
x1 <= 3.25
y0 2.0x0 -x1 >= 1.25
```

(b) Query 2

**Figure 3** Marabou queries for [controller-lemma](#). The lemma is true iff Marabou cannot find an assignment of variables that satisfy the either set of inequalities.

In contrast to the solvers based on SMT technology, solvers based on abstract interpretation and interval arithmetic can only handle properties that involve reasoning about how regions in the input space get mapped to regions in the output space. Concretely this means that they are unable to solve inequalities that involve both input and output variables, such as the last line of the queries in Figure 3. In order to maximise the number of interesting properties solvable, we therefore choose Marabou as the first solver to integrate into Vehicle.

## 2.3 Verified neural network properties: the state of the art

As discussed in the introduction, most of the work using neural network verifiers has focused on the verification of the robustness of the network. Informally, a network is robust if when you move no more than a small distance in the input space, then the result of the neural network should only move a small distance in the output space. Several different types of robustness have been studied including, classification robustness, Lipschitz robustness, standard robustness and approximate classification robustness [6].

One of the first neural network verifiers, Reluplex [18], was also one of the first to verify non-trivial domain-specific properties, proving several results about ACAS Xu, a collision avoidance system for unmanned aircraft. The ACAS Xu neural networks map inputs such as the aircraft's own speed and heading, and the angle and distance to the intruder aircraft to 5 different output actions, ranking each with a confidence score. The action with the highest confidence score being the one that the system will perform. In their verification the authors checked that the ACAS Xu worked as expected and avoided collisions, properties including, checking that a distant aircraft's path will mean that it remains clear of collision

Language	Paradigm	Host language	Verifiers supported	Typed	ITP integration	Multiple networks	Probabilistic properties
Vehicle	Functional	None	1	Yes	Yes	Yes	No
DNNV [25]	Imperative	Python	13	No	No	No	No
Socrates [24]	Declarative	None	2	No	No	No	Yes

■ **Table 1** Comparison of existing property languages for neural networks

271 given various conditions such as angle and speed, checking that despite previous actions it  
 272 will still perform the safe response given the new presence of a possible collision object etc.

273 More recent work in verification focused on the verification of a reinforcement learning  
 274 based neural network controller [16]. The authors show that they can verify that the function  
 275 will always terminate with certain guarantees. In their case study, the authors use an example  
 276 of a car climbing up a hill and verify that the car will always reach the top of the hill with at  
 277 least 90 reward. This is the first paper of its kind formally verifying reachability properties  
 278 for neural network components. Building on this same techniques the authors expand their  
 279 earlier work by formally verifying a reinforcement learning based controller for an autonomous  
 280 vehicle [15]. The authors first train some controllers for the vehicle and then verify its safety.  
 281 Specifically they are able to verify that the vehicle will never be less than 30cm away from a  
 282 wall and consequently will not crash. Unlike previous work by Katz [18] which deals with  
 283 ReLu activation functions, their tool Verisig supports neural networks with smooth activation  
 284 functions (e.g. sigmoid), however it only scales to small networks of about 100 neurons.

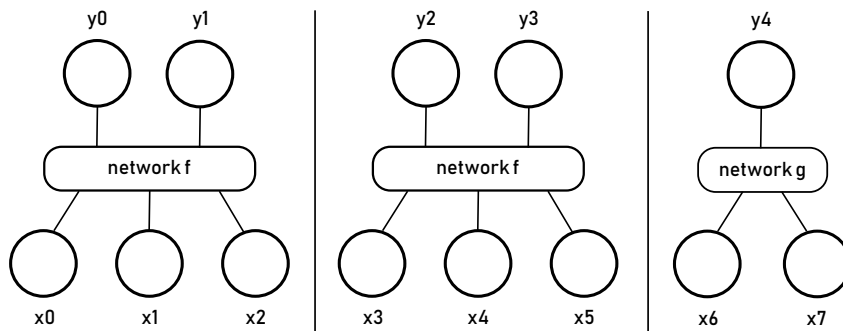
## 285 2.4 Other neural network specification languages

286 Given the low-level input formats supported by the verifiers described in Section 2.2, it  
 287 is unsurprising there have been other attempts at coming up with a high-level property  
 288 language. A comparison is shown in Table 1.

289 The first is the Deep Neural Network Verifier (DNNV) toolbox [25], which has an internal  
 290 Python DSL called DNNP. However its aims are somewhat orthogonal to that of Vehicle, as  
 291 its primary focus is on providing a unified interface for many different verifiers. In particular,  
 292 it has the ability to refactor the structure of the neural network to eliminate unsupported  
 293 operators. However, DNNP is untyped and relies on Python semantics and therefore would  
 294 be challenging to integrate into ITPs. Its dependence on Python also makes it difficult to  
 295 use in other languages commonly used with neural networks such as C++.

296 The second, yet unpublished attempt, is Socrates [24]. Again it positions itself as a  
 297 platform for neural network analysis, which aims to unify different tools. The DSL is  
 298 comparatively limited, primarily supporting different forms of robustness properties using  
 299 a structured JSON file. It also has the disadvantage that one must redefine the internal  
 300 structure of the network within the specification. However, one notable feature is its ability  
 301 to specify probabilistic queries, for example no more than 10% of inputs violate the property.





■ **Figure 4** Proposed extension to the Marabou query language to support properties involving multiple networks and multiple applications of the same network. Input and output variables are labelled sequentially in the order that the networks are passed to Marabou. The diagram shows the proposed labelling of input and output variables for a property that applies network  $f$  to two different inputs and network  $g$  to one input.

### 3 Multi-network specifications and the Marabou query language

As described in Section 2.2, Marabou queries use the variables  $x_1, \dots, x_m$  to represent the inputs to the network and  $y_1, \dots, y_n$  to represent the outputs of the network, and one consequence of this is that it is unable to represent queries that involve multiple networks or applying the same network to more than one input. This situation is suboptimal as there are several such queries that one might be interested in verifying. For example when using teacher-student training [14], one might want to prove that the output of the *student* network is approximately equal to that of the *teacher* network.

$$\forall x : | \text{student}(x) - \text{teacher}(x) | \leq \epsilon$$

Alternatively one might want to prove that a neural network  $f$  is a monotonic function with respect to one or more of its inputs [28], which requires reasoning about the output of the network when applied to two distinct inputs.

$$\forall x_1, x_2 : x_1 \leq x_2 \Rightarrow f(x_1) \leq f(x_2)$$

Luckily, the inability to solve such queries is a limitation of the query language rather than a fundamental shortcoming of the verification engine. This is because Marabou uses an SMT-based approach, and internally has no concept of a network. Instead it represents the nodes and edges simply as a set of variables and constraints between them. One of our purposes in designing Vehicle is to explore what a high-level interface for neural network verification should look like, and consequently we think there is significant value to be gained in targeting maximally expressive verifiers, even if they do not yet exist. We therefore now propose a conservative, backwards-compatible extension to the Marabou query language which will be targeted by the Vehicle compiler.

Conceptually the idea is very simple: if a property involves multiple neural network applications, then the set of applications is assigned an order, and additional input and output variables are then assigned sequentially using this order. An example illustration can be seen in Figure 4. In theory this can be seen as composing the networks in parallel, although in practice the different networks will still be stored in individual ONNX files.

We should stress that this extension is not yet implemented by Marabou, although we hope to do so in the near future. Therefore only Vehicle properties that involve a single application of a single network can currently be verified by Marabou.

```

type InputVector = Tensor Rat [2]

network controller : InputVector -> Rat

currentPosition : InputVector -> Rat
currentPosition x = x ! 0

previousPosition : InputVector -> Rat
previousPosition x = x ! 1

safeInput : InputVector -> Bool
safeInput x = -3.25 <= currentPosition x <= 3.25 and
              -3.25 <= previousPosition x <= 3.25

safeOutput : InputVector -> Bool
safeOutput x = let y = controller x in
               -1.25 < y + 2 * currentPosition x - previousPosition x < 1.25

safe : Prop
safe = forall x . safeInput x => safeOutput x

```

■ **Figure 5** The specification of the safety property expressed in Vehicle for car’s neural network controller.

## 332 4 Vehicle

### 333 4.1 Specification language

334 The Vehicle specification language is a functional language with Haskell-like syntax. At its  
 335 centre is a small dependently-typed core, upon which various built-in operators and types  
 336 are then added. Figure 5 shows one possible formulation of the specification for the running  
 337 example, and will be used to explain the key features of the language. The full BNF grammar  
 338 of the language can be found online.

339 In order to better abstract away the representation of the inputs of our network, the  
 340 first line of the specification declares `InputVector` to be a synonym for the type of 1-  
 341 dimensional rational tensors of length 2. Vehicle has a set of builtin types that includes  
 342 `Bool`, `Int`, `Rat`, `Real` and `Tensor`. An observant reader may note that neural networks use  
 343 floating point arithmetic whereas we are using rationals in our specification. We acknowledge  
 344 this compromises soundness, and aim for Vehicle to support floating point types in the  
 345 future. Some neural network verifiers have recently been found to have similar unsoundness  
 346 problems [17].

347 Next, the car’s controller is bound to the name `controller` using the `network` keyword.  
 348 As previously discussed, we are only required to provide a name and a type for the network  
 349 in the specification, and the implementation of the network is provided later to the Vehicle  
 350 compiler in the form of an ONNX file. Consequently the network remains a black-box  
 351 function from the perspective of the specification, while still allowing us to write expressive  
 352 properties which can be statically type-checked.

353 The local functions `currentPosition` and `previousPosition` assign meaningful names  
 354 to the first and second components of the input vector. The `!` operator looks up the value  
 355 of the tensor at the provided index. The `safeInput` and `safeOutput` declarations then use

356 these to provide human-readable statements of the pre-condition and post-conditions of  
 357 `controller-lemma`.

358 Finally, the `safe` declaration assembles these pieces together to complete the specification.  
 359 It uses the universal quantifier `forall` to bind a new variable `x` representing an arbitrary  
 360 input to the network and then states that whenever `safeInput x` is true then `safeOutput`  
 361 `x` is true as well. Note that the type of `safe` is `Prop` rather than `Bool`. The `Prop` type  
 362 represents the type of boolean expressions whose value cannot be decided within `Vehicle`  
 363 itself. Most of the built-ins that use booleans are polymorphic with respect to either `Bool`  
 364 and `Prop`. The exceptions are the quantifiers `forall` and `exists` which always return `Prop`,  
 365 and `if then else` which always requires that the condition must be of type `Bool`.

## 366 4.2 Compilation to Marabou

367 The `Vehicle` compiler is implemented in Haskell and uses a lexer and parser generated by  
 368 BNFC [12]. The type-checking algorithm is based on the one presented in [20], with the  
 369 addition of type classes and unification-based term inference.

370 We will now describe the compiler passes that translate a type-checked `Vehicle` program  
 371 into verification queries suitable for Marabou. Note that only the very last pass in the  
 372 pipeline does anything that is specific to Marabou, and therefore it should be relatively easy  
 373 to target further verifiers.

### 374 4.2.1 Network type analysis

375 The first step is to check the networks declared in the specification against their implementa-  
 376 tions in the ONNX files provided by the user. Using a custom-written Haskell bindings for  
 377 the C implementation of ONNX, `Vehicle` reads the type information from the ONNX file and  
 378 checks that it matches that declared by the user.

379 Next it checks that the network type is supported by `Vehicle`. Although the ONNX  
 380 format is significantly more expressive, supporting multiple tensor inputs with different sizes  
 381 and element types, at the moment `Vehicle` supports only networks of following type, where `A`,  
 382 `B` are one of `Nat/Int/Rat/Real` and `m` and `n` are literals:

383 `Tensor A [m] -> Tensor B [n]`

384 However it also allows syntactic sugar for this pattern in the form of:

385 `A -> ... -> A -> B`

386 If the type is in this form, then during this pass, it is normalised to `Tensor A [n] -> Tensor`  
 387 `B [1]`. This necessitates traversing the program to find any applications of the network,  
 388 `f x1 ... xn` and replacing them with `f [x1, ... , xn] ! 0`, where the `[ ]` syntax  
 389 constructs a tensor from the comma-separated list of elements contained within the brackets.

390 Finally, the network declarations are removed from the program, and the names and  
 391 types of the networks are stored in the *network context* which is passed along to subsequent  
 392 stages in the compiler.

### 393 4.2.2 Normalisation

394 The next step is normalisation. As well as performing the standard operations such as  
 395 beta-reduction, normalisation of builtin operations applied to constants and substituting  
 396 through any references to top-level functions, some domain-specific operations are required.

## 23:12 Vehicle: NNs to ITPs

397 Firstly, the verifier input format described in Section 2.2, assigns variables to each  
398 individual input and therefore quantifiers over tensor variables must be converted to multiple  
399 quantifiers over its elements, e.g. `forall (x : Tensor A [2, 2])` is normalised to `forall`  
400 `(x11 x12 x21 x22 : A)`. Secondly, after normalisation, only top-level declarations with  
401 type `Prop` are of interest, as the rest should have been substituted through. Declarations  
402 which do not have type `Prop` are therefore removed from the program. After normalisation,  
403 we are therefore left with the following Vehicle program:

```
safe : Prop  
safe : forall (p0 p1 : Rat) .  
404   (3.25 <= p0 <= -3.25) and (3.25 <= p1 <= -3.25) =>  
   -1.25 < controller [p0, p1] + 2 * p0 - p1 < 1.25
```

### 405 4.2.3 Subdivision of queries

406 Neural network verifiers only support solving existential queries involving conjunctions of  
407 numeric equalities and inequalities. We now describe how a Vehicle property is reduced to a  
408 set of such queries.

409 Initially, Vehicle traverses the property making note of the set of quantifiers used. If only  
410 existential quantifiers are used then the property passes through this stage untouched. If  
411 only universal quantifiers are used then the property is negated. If both types of quantifier  
412 are used then the compiler will emit an error.

413 Next, any `if` statements contained within the property are eliminated, using the trans-  
414 formation:

```
415   if a then b else c   ⇒   a => b and not a => c
```

416 Note that this transformation is only valid if the arguments of `if` statement have type `Prop`  
417 or `Bool`. However this may not be the case, e.g. `exists x . (if a then x else x + 2)`  
418 `>= 8`. Nonetheless, as the overall type of the Vehicle property is guaranteed to be `Prop`, it is  
419 always possible to lift the `if` expression recursively until the arguments have type `Prop`, e.g.  
420 `exists x . if a then x >= 8 else x + 2 >= 8`, and then perform the elimination.

421 Next, the expression is converted to disjunctive normal form, with implications being  
422 converted to `ors` and `or` expressions being lifted to the top-level. At this point the following  
423 invariants should hold of the property: only existential quantifiers, no negations, no `if`  
424 statements, and the only tensor literals present should be the input to the networks. Our  
425 running example is therefore:

```
exists (p0 p1 : Rat) .  
  (3.25 <= p0 <= -3.25) and (3.25 <= p1 <= -3.25) and  
  controller [p0, p1] + 2 * p0 - p1 < -1.25  
426 or  
exists (p0 p1 : Rat) .  
  (3.25 <= p0 <= -3.25) and (3.25 <= p1 <= -3.25) and  
  controller [p0, p1] + 2 * p0 - p1 > 1.25
```

427 Each disjunction is now split up into its own query. From this point onwards, we will  
428 only follow the compilation of the first query in the running example.

### 429 4.2.4 Moving from a functional to relational model of networks

430 The next stage is to move from our model of the network as a function to the relational  
431 model used by the verifiers. As discussed in Section 3, we aim to support multiple neural

432 networks applications, and therefore this is a little more involved than one might at first  
433 suspect.

434 The first step is to construct a list of the network applications in the query. Despite  
435 supporting multiple applications of the same network, we must be careful not to duplicate  
436 applications of the same network to the same input, as doing so would result in an exponential  
437 decrease in performance during verification. Therefore in order to avoid this, we first perform  
438 a common-sub-expression elimination pass, binding all network applications to fresh variables  
439 using let-expressions. We use an efficient hash-based approach procedure using co-de-Bruijn  
440 indices [22]. Once this pass is complete, we can generate the required list of network  
441 applications simply by calculating the set of free variables that occur in the query.

442 In the next step, we recurse downwards into the expression, keeping track of our position  
443 in the list of network applications and replacing every let-bound neural network application  
444 with an expression that a) equates the inputs to the application with the input variables,  
445 and b) substitutes a list of the output variables for the bound variable in the body of the  
446 let-binding. For example, if we are looking at an application of network  $f$  with  $n$  inputs and  
447  $m$  outputs, and the list of network applications that we have traversed so far has  $k$  inputs  
448 and  $l$  outputs in total then the following transformation would be performed:

```

449   let y = f [e1, ..., en] in e
                                     ⇒
   [e1, ..., en] == [X<k+1> ... X<k+n>] and e{y/[Y<l+1>, ..., Y<l+m>]}

```

450 As with if-elimination in Section 4.2.3, this is only a valid transformation if the type of  $e$   
451 is Prop. However, this is guaranteed as the common sub-expression elimination inserts the  
452 lets at the top-most position, i.e. just before the quantifiers of the variables bound in the  
453 expression.

454 However, there is one more niggle, as it is necessary to eliminate the variables quantified  
455 over by the user. Therefore, when performing the above substitution the compiler tracks  
456 which user variables are equated to which introduced input and output variables. Upon  
457 recursing back up the tree, when a quantifier is reached, it checks for an equated input or  
458 output variable. If there is, we remove the quantifier and substitute the latter through in its  
459 place. If there is not, then at the moment, the compiler errors. Note that in future it should  
460 be possible to refine this analysis, to reduce the number of such errors. For example the  
461 property `exists v . f (v + 2) <= 0` would generate the constraint  $x_0 = v + 2$  which  
462 currently errors, but could in fact be rearranged to the form  $v = x_0 - 2$  and then substituted  
463 through.

464 Finally, the compiler then reruns the normalisation pass in order to simplify the introduced  
465 tensor expressions. This leaves the first query from the running example as:

```

466   (3.25 <= x0 <= -3.25) and
   (3.25 <= x1 <= -3.25) and
   y0 + 2 * x0 - x1 < -1.25

```

## 467 4.2.5 Conversion to Marabou syntax

468 Finally we now specifically target Marabou's query language, recursing through the query  
469 and splitting conjunctions into separate assertions. Each individual assertion is first checked  
470 for linearity and then rearranged to put the constant on the right hand side of the relation.  
471 The left-hand side is then transformed into the required syntax, resulting in the query in

## 23:14 Vehicle: NNs to ITPs

472 Figure 3a. This query is then written to a user-defined location ready to be fed into Marabou.  
473 Marabou proves both queries in approximately 20 seconds on a mid-range laptop.

### 474 4.3 Compilation to Agda

475 After compiling the Vehicle specification to Marabou and verifying the resulting queries, we  
476 can now use it to complete the overall proof that the car never leaves the road by compiling  
477 it to Agda. As before, the Vehicle program is type-checked and the network types compared  
478 against those in the ONNX file. The subsequent transformation to Agda code is relatively  
479 simple. The only non-trivial insight needed is that any Vehicle expression of type `Prop` must  
480 be lifted to the `Set` type in Agda. This leads to having two different methods of compiling  
481 each built-in, depending on whether it is being instantiated with the `Bool` or `Prop` type. The  
482 output of compiling the example specification is:

```
483 module ControllerSpec where
484
485 InputVector : Set
486 InputVector = Tensor ℚ (2 :: [])
487
488 postulate controller : InputVector → ℚ
489
490 currentPositon : InputVector → ℚ
491 currentPositon x = x (# 0)
492
493 previousPositon : InputVector → ℚ
494 previousPositon x = x (# 1)
495
496 SafeInput : InputVector → Set
497 SafeInput x = (ℚ - (ℤ + 13 ℚ / 4) ℚ ≤ currentPositon x × currentPositon x ℚ ≤ ℤ + 13 ℚ / 4)
498             × (ℚ - (ℤ + 13 ℚ / 4) ℚ ≤ previousPositon x × previousPositon x ℚ ≤ ℤ + 13 ℚ / 4)
499
500 SafeOutput : InputVector → Set
501 SafeOutput x =
502   ℚ - (ℤ + 5 ℚ / 4) ℚ < (controller x ℚ + (ℤ + 2 ℚ / 1) ℚ * currentPositon x) ℚ - previousPositon x
503   × (controller x ℚ + (ℤ + 2 ℚ / 1) ℚ * currentPositon x) ℚ - previousPositon x ℚ < ℤ + 5 ℚ / 4
504
505 abstract
506 safe : ∀ (x : Tensor ℚ (2 :: [])) → SafeInput x → SafeOutput x
507 safe = checkVehicleProperty record
508   { propertyFile = "path/to/property/file.vclp"
509     ; propertyName = "safe"
510   }
511
```

512 There are a few things to note. Firstly, the network is declared as a `postulate` and therefore  
513 cannot be evaluated within Agda. This is not a fundamental limitation, and with a bit of  
514 effort the Haskell bindings created for ONNX could be lifted to Agda. Secondly, the desired  
515 proof `safe` is within an `abstract` block which prevents code that uses the generated module  
516 from depending on the implementation of the proof.

517 Finally, the definition of the proof is implemented via a macro `checkVehicleProperty` which  
518 calls out to the Vehicle compiler. A naive implementation would have Agda use a reference  
519 to the query files to call Marabou directly. However, while Marabou is running the user  
520 would be unable to interact with the file. As Marabou takes over 20 seconds to verify these  
521 queries, and potentially much longer for more complex queries, this would unacceptably  
522 degrade the user experience.

523 Instead we require users to explicitly ask Vehicle to use Marabou to verify the specification.  
524 If successful then Vehicle will write the result to a Vehicle proof cache file which contains  
525 locations and hashes of the ONNX files and the verification status of the specification. It  
526 is a reference to this file rather than the original source code that the Agda macro passes  
527 back to Vehicle. Vehicle reads this file, uses the hashes to check that the network has not  
528 changed on disk, and then returns the verification status of the property to Agda. There is  
529 however one missing piece of the puzzle, namely the integrity of the generated Agda code.  
530 In theory a malicious or careless user could change the type of the proof in the Agda code  
531 which would currently not be detectable by Vehicle. While Vehicle does store the hash of  
532 the generated Agda code in the proof file, Agda macros cannot query the location of the  
533 source file from which they are called and so Vehicle cannot locate the current Agda file to  
534 re-hash. Hopefully this short-coming in Agda can be fixed in the future.

535 Leaving aside this issue, in summary the generated Agda module is now linked to both  
536 the underlying ONNX file and the result of the Marabou verifier. The module can now be  
537 imported into the main development outlined in Section 2.1 and the `safe` proof can now be  
538 used in the proof of `controller-lemma`. This therefore completes the formal proof of Theorem 1  
539 that, subject to the assumptions that the sensor error is never more than 0.25 and the wind  
540 never shifts by more than 1 per unit time, the neural network controlled vehicle will never  
541 leave the road.

## 542 **5 Conclusions and future work**

543 This paper has described how the Vehicle tool can be used to express high-level specifications  
544 for neural networks. These can then be compiled to low-level queries for neural network  
545 verifiers, and to high-level postulates within an interactive theorem prover, linked by hashes, to  
546 provided end-to-end verification of neural-network powered systems. We have demonstrated  
547 this process by proving the safety of a system involving a 20,000 node neural network  
548 controller. All accompanying code, including Vehicle itself, is available online [7].

549 We expect that as the field of neural network verification matures, a system like Vehicle  
550 that integrates automated and interactive theorem proving will be required to ensure that  
551 verified neural network specifications support the specifications of the larger systems that  
552 contain them.

553 Before discussing possible future work, we should emphasise that building Vehicle has  
554 been a significant undertaking, as many of the necessary components are still relatively  
555 immature. Firstly, in order to allow the Vehicle compiler to read the ONNX files we have  
556 had to create Haskell bindings for the C version of the ONNX runtime library. Secondly,  
557 the current version of Marabou which is implemented in C++, only supports networks in  
558 the ONNX format via its Python bindings. We have therefore had to patch support into  
559 Marabou for reading ONNX files natively in C++. This involved using undocumented C++  
560 bindings to the ONNX format to manually traverse and parse the networks into Marabou's  
561 internal representation. Therefore, while we believe that neural network verification has  
562 a bright future, there is significant work to be done before mature, robust, flexible and

563 interoperable tools can be developed.

## 564 5.1 Improved integration

565 One possible future direction is to increase the number of ITPs that Vehicle connects to, and  
 566 in particular those with better support extracting verified code. While there are no reasons  
 567 why popular theorem provers such as Coq, Isabelle and Idris cannot be targeted, there are  
 568 also some more interesting potential targets. Among these is a relatively new ITP called  
 569 KeYmaera X [13] which is based on differential dynamic logic and is specifically designed to  
 570 verify the safety of real-world controllers using continuous dynamics. Alternatively, adding  
 571 the ability to compile to a suitable training framework such as Tensorflow would allow the  
 572 same specification to be used in training.

## 573 5.2 Language features

574 There are many additional features that could be added to the Vehicle language. This  
 575 includes adding top-level `parameter` declarations, which would allow users to pass arbitrary  
 576 values in at compile time rather than hard-coding them in the property. For example, the  
 577 epsilon value in specifications of robustness properties [6], or the training dataset with which  
 578 the network is robust to. Such a feature would facilitate the feedback loop between the  
 579 verifier and reinforcement learning algorithms implemented in [15].

580 Another improvement would be for the compiler to read the names of the input and  
 581 output variables from the ONNX file, and use them to automatically bind functions mapping  
 582 input and output vectors to their elements. In the running example this would remove the  
 583 need to explicitly define the functions `currentPosition` and `previousPosition`.

## 584 5.3 A call to arms

585 We would like to end with a call to arms to the neural network verifier community. Unlike  
 586 SMT solvers which have converged on SMTLib as a unified input format, every neural network  
 587 verifier currently uses their own incompatible input format. While it is impressive that  
 588 DNNV [25] unifies 13 different verifiers, it is high-level language and therefore is unsuitable  
 589 to build other tools such as Vehicle on top of. We argue that in order to have a thriving  
 590 ecosystem of both backends and frontends, a common *low-level* interface is needed for solvers.  
 591 There is currently a fledgling proposal out there in the form of VNNLib [2], but it needs  
 592 significantly more work to ensure it is suitably expressive to capture all possible properties  
 593 of interest.

## 594 ——— References ———

- 595 1 Open Neural Network Exchange format. Accessed on 30.01.2022. URL: <https://onnx.ai/>.
- 596 2 VNNLib format. Accessed on 01.12.2021. URL: <https://vnnlib.org/>.
- 597 3 Alexander Bagnall and Gordon Stewart. Certifying the true error: Machine learning in Coq  
 598 with verified generalization guarantees. In *Proceedings of the AAAI Conference on Artificial  
 599 Intelligence*, volume 33, pages 2662–2669, 2019.
- 600 4 Jasmin Christian Blanchette, Sascha Böhme, and Lawrence C Paulson. Extending Sledgeham-  
 601 mer with SMT solvers. *Journal of automated reasoning*, 51(1):109–128, 2013.
- 602 5 Robert S Boyer, Milton W Green, and J Strother Moore. The use of a formal simulator to  
 603 verify a simple real time control program. In *Beauty Is Our Business*, pages 54–66. Springer,  
 604 1990.



- 605 6 Marco Casadio, Matthew L. Daggitt, Ekaterina Komendantskaya, Wen Kokke, Daniel Kienitz,  
606 and Rob Stewart. Property-driven training: All you (n) ever wanted to know about. *arXiv*  
607 *preprint arXiv:2104.01396*, 2021.
- 608 7 Matthew L. Daggitt and Wen Kokke. Vehicle. Accessed on 09.02.2022. URL: <https://github.com/vehicle-lang/vehicle>.  
609
- 610 8 Elisabetta De Maria, Abdorrahim Bahrami, Thibaud l'Yvonnet, Amy Felty, Daniel Gaffé,  
611 Annie Ressouche, and Franck Grammont. On the use of formal methods to model and verify  
612 neuronal archetypes. *Frontiers of Computer Science*, 16(3):1–22, 2022.
- 613 9 Ruediger Ehlers. Formal verification of piece-wise linear feed-forward neural networks. In  
614 *International Symposium on Automated Technology for Verification and Analysis*, pages 269–  
615 286. Springer, 2017.
- 616 10 Burak Ekici, Alain Mebsout, Cesare Tinelli, Chantal Keller, Guy Katz, Andrew Reynolds, and  
617 Clark Barrett. SMTCoq: A plug-in for integrating SMT solvers into Coq. In *International*  
618 *Conference on Computer Aided Verification*, pages 126–133. Springer, 2017.
- 619 11 Marc Fischer, Mislav Balunovic, Dana Drachler-Cohen, Timon Gehr, Ce Zhang, and Martin  
620 Vechev. DL2: training and querying neural networks with logic. In *International Conference*  
621 *on Machine Learning*, pages 1931–1941. PMLR, 2019.
- 622 12 Markus Forsberg and Aarne Ranta. BNF converter. In *Proceedings of the 2004 ACM SIGPLAN*  
623 *workshop on Haskell*, pages 94–95, 2004.
- 624 13 Nathan Fulton, Stefan Mitsch, Jan-David Quesel, Marcus Völp, and André Platzer. KeY-  
625 maeraX: An axiomatic tactical theorem prover for hybrid systems. In Amy P. Felty and  
626 Aart Middeldorp, editors, *CADE*, volume 9195 of *LNCS*, pages 527–538. Springer, 2015.  
627 doi:10.1007/978-3-319-21401-6\_36.
- 628 14 Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. Knowledge distillation:  
629 a survey. *International Journal of Computer Vision*, 129(6):1789–1819, 2021.
- 630 15 Radoslav Ivanov, Taylor J Carpenter, James Weimer, Rajeev Alur, George J Pappas, and  
631 Insup Lee. Case study: verifying the safety of an autonomous racing car with a neural  
632 network controller. In *Proceedings of the 23rd International Conference on Hybrid Systems:*  
633 *Computation and Control*, pages 1–7, 2020.
- 634 16 Radoslav Ivanov, James Weimer, Rajeev Alur, George J Pappas, and Insup Lee. Verisig:  
635 verifying safety properties of hybrid systems with neural network controllers. In *Proceedings*  
636 *of the 22nd ACM International Conference on Hybrid Systems: Computation and Control*,  
637 pages 169–178, 2019.
- 638 17 Kai Jia and Martin Rinard. Exploiting verified neural networks via floating point numerical  
639 error. In *International Static Analysis Symposium*, pages 191–205. Springer, 2021.
- 640 18 Guy Katz, Clark Barrett, David L Dill, Kyle Julian, and Mykel J Kochenderfer. Reluplex:  
641 An efficient smt solver for verifying deep neural networks. In *International Conference on*  
642 *Computer Aided Verification*, pages 97–117. Springer, 2017.
- 643 19 Guy Katz, Derek A Huang, Duligur Ibeling, Kyle Julian, Christopher Lazarus, Rachel Lim,  
644 Parth Shah, Shantanu Thakoor, Haoze Wu, Aleksandar Zeljić, et al. The Marabou framework  
645 for verification and analysis of deep neural networks. In *International Conference on Computer*  
646 *Aided Verification*, pages 443–452. Springer, 2019.
- 647 20 Andres Löb, Conor McBride, and Wouter Swierstra. A tutorial implementation of a dependently  
648 typed lambda calculus. *Fundamenta informaticae*, 102(2):177–207, 2010.
- 649 21 Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
650 Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*,  
651 2017.
- 652 22 Conor McBride. Everybody's got to be somewhere. *Electronic Proceedings in Theoretical*  
653 *Computer Science*, 275:53–69, Jul 2018. URL: <http://dx.doi.org/10.4204/EPTCS.275.6>,  
654 doi:10.4204/eptcs.275.6.
- 655 23 Ulf Norell. Dependently typed programming in Agda. In *International school on advanced*  
656 *functional programming*, pages 230–266. Springer, 2008.

- 657 24 Long H Pham, Jiaying Li, and Jun Sun. Socrates: Towards a unified platform for neural  
658 network analysis. *arXiv preprint arXiv:2007.11206*, 2020.
- 659 25 David Shriver, Sebastian Elbaum, and Matthew B. Dwyer. DNNV: A framework for deep  
660 neural network verification. In Alexandra Silva and K. Rustan M. Leino, editors, *Computer  
661 Aided Verification*, pages 137–150, Cham, 2021. Springer International Publishing.
- 662 26 Gagandeep Singh, Timon Gehr, Markus Püschel, and Martin Vechev. An abstract domain for  
663 certifying neural networks. *Proceedings of the ACM on Programming Languages*, 3(POPL):1–30,  
664 2019.
- 665 27 Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian  
666 Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International  
667 Conference on Learning Representations*, 2013. URL: <http://arxiv.org/abs/1312.6199>.
- 668 28 Antoine Wehenkel and Gilles Louppe. Unconstrained monotonic neural networks. *Advances in  
669 neural information processing systems*, 32, 2019.
- 670 29 Weiming Xiang, Hoang-Dung Tran, and Taylor T Johnson. Output reachable set estimation  
671 and verification for multilayer neural networks. *IEEE transactions on neural networks and  
672 learning systems*, 29(11):5777–5783, 2018.