

On Taming Differentiable Logics

Reynald Affeldt, Alessandro Bruni, Matthew Daggitt, Ekaterina Komendantskaya,
Natalia Slusarz, **Kathrin Stark**, Robert Stewart

This Work



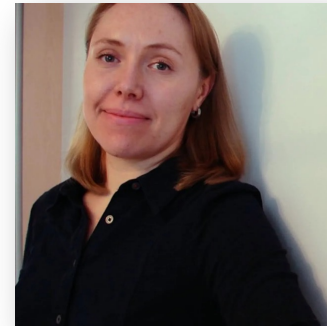
Reynald Affeldt



Alessandro Bruni



Matthew Daggitt



Katya
Komendantskaya



Natalia Ślusarz



Kathrin Stark



Rob Stewart

Towards Ensuring that Neural Networks Satisfy Verification Properties

We want a neural network N to satisfy certain properties P ; for example:

$\epsilon - \delta$ robustness

Given a vector v , N is said to be $\epsilon - \delta$ robust w.r.t v if:

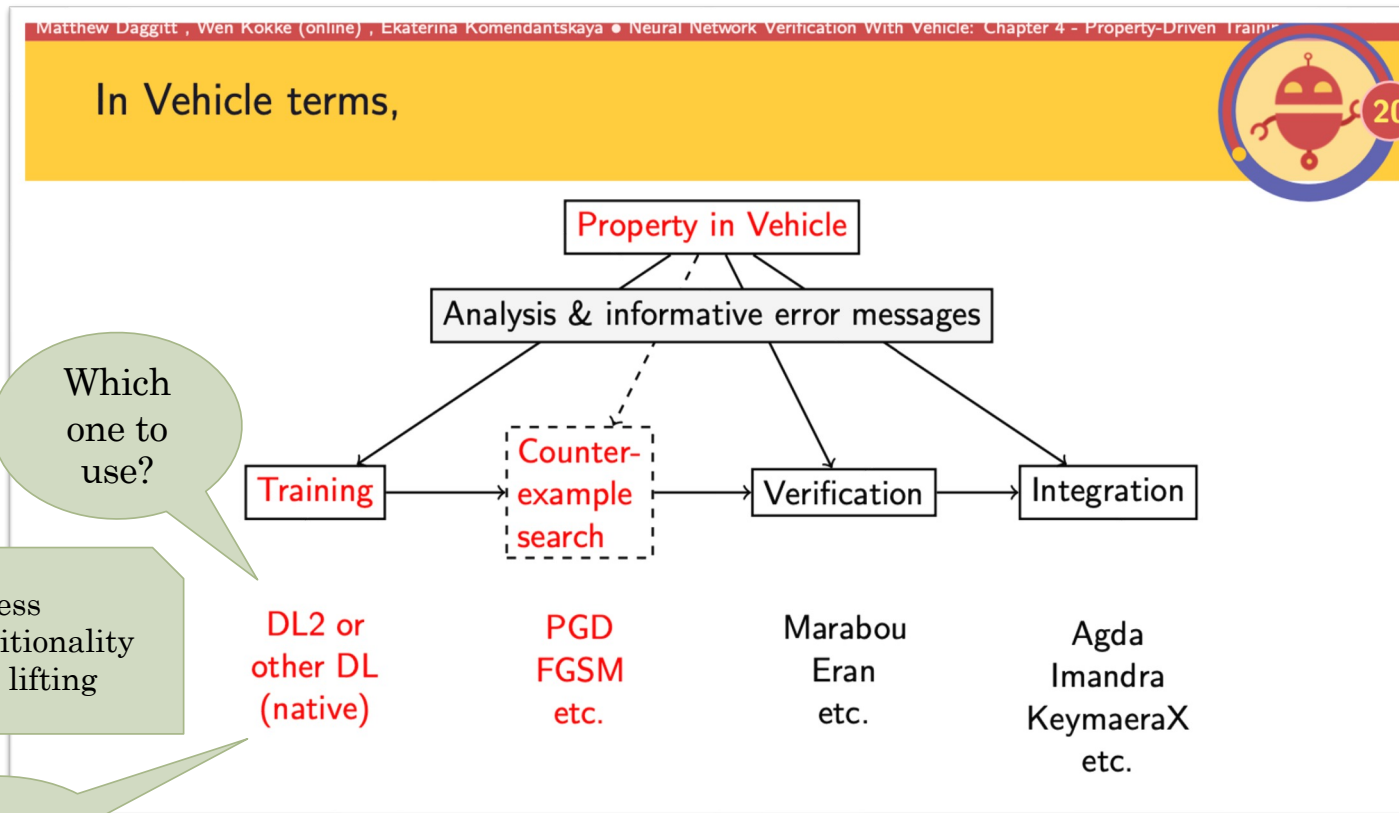
$$\forall x. |x - v|_{L_\infty} \leq \epsilon \Rightarrow |N(x) - N(v)|_{L_\infty} \leq \delta'$$

But: Even the most accurate neural networks fail even the most natural verification properties, such as robustness [Fischer et al., '19].

A possible solution: Translate the desired property P into a differentiable **loss function** that punishes not satisfying P (**Differentiable Logic/DL**) in the tradition of property-based training for neural networks [Giunchiglia et al., IJCAI '22]

For example:
 $\llbracket e_1 \leq e_2 \rrbracket$
 $= -\max(\llbracket e_1 \rrbracket - \llbracket e_2 \rrbracket, 0)$

Towards PL Tools for Neural Networks



<https://vehicle-lang.github.io/tutorial/#introduction>

This Talk

A **common framework** for different differentiable logics that allows to give a uniform semantics:

LDL (Logic of Differentiable Logic)

[Ślusarz et al., LPAR '23]

Goal: A generic framework in which logical/geometric properties of different DLs can be formalized/proven.

A common framework to **compare properties** of different DLs

A **mechanization*** in Coq/MathComp, allowing to easily test out new extensions/different approaches
[Affeldt et al., ITP '24]

* So far: Without quantifiers

This Talk

A **common framework** for different differentiable logics that allows to give a uniform semantics:

LDL (Logic of Differentiable Logic)

[Ślusarz et al., LPAR '23]

Goal: A generic framework in which logical/geometric properties of different DLs can be formalized/proven.

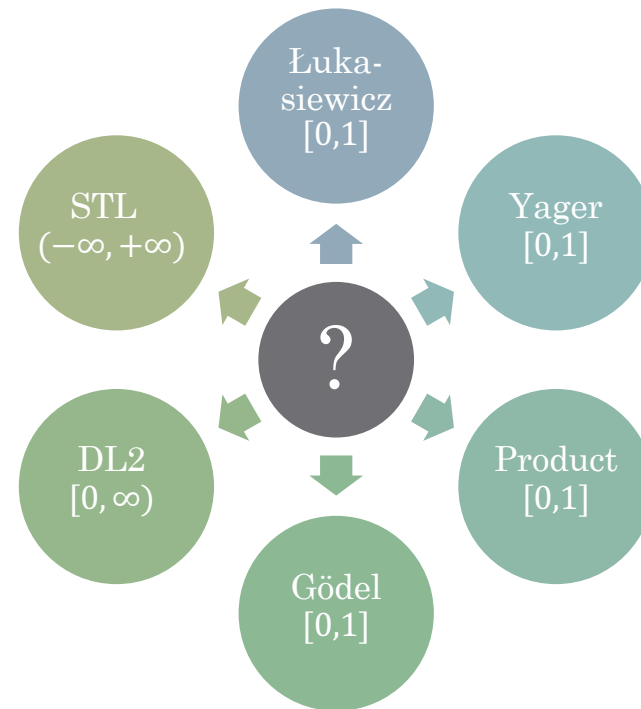
A common framework to **compare properties** of different DLs

A **mechanization*** in Coq/MathComp, allowing to easily test out new extensions/different approaches
[Affeldt et al., ITP '24]

* So far: Without quantifiers

Possible Differentiable Logics

- DL2 [Fischer et al., '19]
- STL [Varnai, Dimarogonas '20]
- Fuzzy logics [Krieken et al. '21]
 - Gödel
 - Łukasiewicz
 - Yager
 - Product



But:

1. Only propositional fragment
2. Syntax/semantics/pragmatics partially not well-separated
3. Lack of unified general syntax and semantics

LDL: A Logic of Differentiable Logics

Dependently Typed Version from [Affeldt et al., ITP '24]

type $\ni t ::= \text{Bool} \mid \text{Index } n \text{ for } n \in \mathbb{N} \mid \text{Real} \mid \text{Vec } n \mid \text{Fun } n \ m \text{ for } n, m \in \mathbb{N}$

No implication

$\text{exprInd } \ni i ::= i \in \mathbb{N}$
 $\text{exprR } \ni r, r_1, r_2 ::= r \in \mathbb{R} \mid [v]_i$
 $\text{exprFun } \ni f ::= f \in \mathbb{R}^n \rightarrow \mathbb{R}^m$
 $\text{exprVec } \ni v ::= v \in \mathbb{R}^n \mid f \ v$

$\text{exprB } \ni p, p_0, \dots, p_n ::= \text{True} \mid \text{False}$

Partially not supported

Required by (non-assoc.) STL

$r_1 \leq r_2$
 $r_1 \neq r_2$
 $\bigwedge_M (p_0, \dots, p_M)$
 $\bigvee_M (p_0, \dots, p_M)$
 $\neg p$

$\text{expr } \ni e ::= x \mid \lambda (x : t). e$

These parts are not yet in the Coq formalisation

$\exists (x : t). p$
 $\forall (x : t). p$

Example: For $\epsilon - \delta$ robustness, given concrete values for $\epsilon/\delta/v/N$ both the L_∞ norm and the right hand of the implication can be expressed via LDL.

One Logic – Several Interpretations

	$\llbracket \bigwedge_M s \rrbracket$	$\llbracket \bigvee_M s \rrbracket$	$\llbracket \neg e \rrbracket$
Gödel	$\min \llbracket s \rrbracket_G$	$\max \llbracket s \rrbracket_G$	$1 - \llbracket e \rrbracket_G$
Łukasiewicz	$\max \left[\sum_{a \in \llbracket s \rrbracket_L} a - s + 1, 0 \right]$	$\min \left[\sum_{a \in \llbracket s \rrbracket_L} a, 1 \right]$	$1 - \llbracket e \rrbracket_L$
Yager	$\max \left[1 - \left(\sum_{a \in \llbracket s \rrbracket_Y} (1-a)^p \right)^{1/p}, 0 \right]$	$\min \left[\left(\sum_{a \in \llbracket s \rrbracket_Y} a^p \right)^{1/p}, 1 \right]$	$1 - \llbracket e \rrbracket_Y$
product	$\prod_{a \in \llbracket s \rrbracket_P} a$	$\text{fold } (\lambda x y . x + y - xy) \ 0 \ \llbracket s \rrbracket_P$	$1 - \llbracket e \rrbracket_P$
DL2	$\sum_{a \in \llbracket s \rrbracket_{DL2}} a$	$(-1)^{ s +1} \cdot \prod_{a \in \llbracket s \rrbracket_{DL2}} a$	undefined [†]
STL	$\text{and}_S \llbracket s \rrbracket_{STL}$	$\text{or}_S \llbracket s \rrbracket_{STL}$	$-\llbracket e \rrbracket_{STL}$
Bool	$\bigwedge_M \llbracket s \rrbracket_B$	$\bigvee_M \llbracket s \rrbracket_B$	$-\llbracket e \rrbracket_B$

	$\llbracket e_1 = e_2 \rrbracket$	$\llbracket e_1 \leq e_2 \rrbracket$	$\llbracket \text{True} \rrbracket$	$\llbracket \text{False} \rrbracket$
f	if $\llbracket e_1 \rrbracket = -\llbracket e_2 \rrbracket$ then $\llbracket e_1 \rrbracket = \llbracket e_2 \rrbracket$ else $\max \left[1 - \left \frac{\llbracket e_1 \rrbracket - \llbracket e_2 \rrbracket}{\llbracket e_1 \rrbracket + \llbracket e_2 \rrbracket} \right , 0 \right]$	if $\llbracket e_1 \rrbracket = -\llbracket e_2 \rrbracket$ then $\llbracket e_1 \rrbracket \leq \llbracket e_2 \rrbracket$ else $\max \left[1 - \max \left[\frac{\llbracket e_1 \rrbracket - \llbracket e_2 \rrbracket}{\llbracket e_1 \rrbracket + \llbracket e_2 \rrbracket}, 0 \right], 0 \right]$	1	0
DL2	$-\llbracket e_2 \rrbracket_{DL2} - \llbracket e_1 \rrbracket_{DL2}$	$-\max \left[\llbracket e_1 \rrbracket_{DL2} - \llbracket e_2 \rrbracket_{DL2}, 0 \right]$	0	$-\infty$
STL	$-\llbracket e_2 \rrbracket_{STL} - \llbracket e_1 \rrbracket_{STL}$	$\llbracket e_2 \rrbracket_{STL} - \llbracket e_1 \rrbracket_{STL}$	$+\infty$	$-\infty$
Bool	$\llbracket e_1 \rrbracket_B = \llbracket e_2 \rrbracket_B$	$\llbracket e_1 \rrbracket_B \leq \llbracket e_2 \rrbracket_B$	True	False

Generalization to n-ary

Strong negation

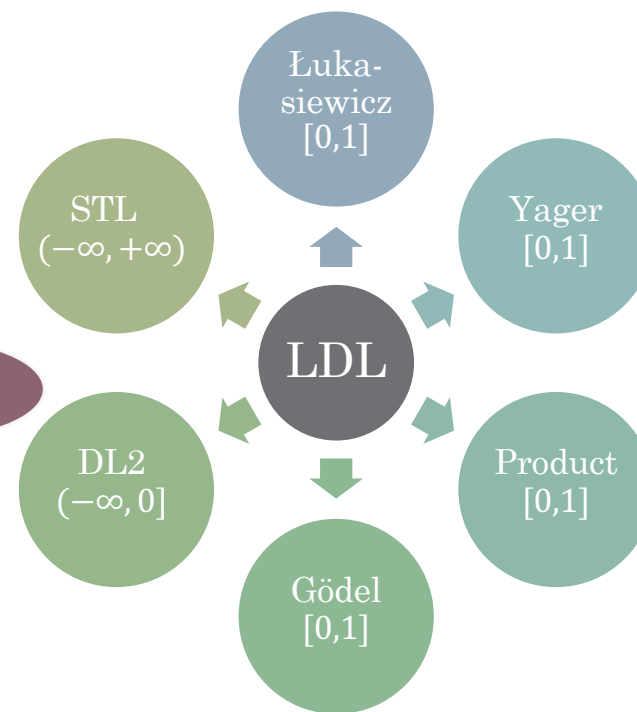
Not part of original definition

Not part of original definition

Originally: $[0, \infty)$

$$\text{and}_S [a_1, \dots, a_M] = \begin{cases} \frac{\sum_i a_{\min} e^{\tilde{a}_i} e^{\nu \tilde{a}_i}}{\sum_i e^{\nu \tilde{a}_i}} & \text{if } a_{\min} < 0 \\ \frac{\sum_i a_i e^{-\nu \tilde{a}_i}}{\sum_i e^{-\nu \tilde{a}_i}} & \text{if } a_{\min} > 0 \\ 0 & \text{if } a_{\min} = 0 \end{cases} \quad \text{where} \quad \begin{cases} \nu \in \mathbb{R}^+ \text{ (constant)} \\ a_{\min} = \min [a_1, \dots, a_M] \\ \tilde{a}_i = \frac{a_i - a_{\min}}{a_{\min}} \end{cases}$$

or_S is analogous to and_S



This Talk

A **common framework** for different differentiable logics that allows to give a uniform semantics:

LDL (Logic of Differentiable Logic)

[Ślusarz et al., LPAR '23]

Goal: A generic framework in which logical/geometric properties of different DLs can be formalized/proven.

A common framework to **compare properties** of different DLs

A **mechanization*** in Coq/MathComp, allowing to easily test out new extensions/different approaches
[Affeldt et al., ITP '24]

* So far: Without quantifiers

This Talk

A **common framework** for different differentiable logics that allows to give a uniform semantics:

LDL (Logic of Differentiable Logic)

[Ślusarz et al., LPAR '23]

Goal: A generic framework in which logical/geometric properties of different DLs can be formalized/proven.

A common framework to **compare properties** of different DLs

A **mechanization*** in Coq/MathComp, allowing to easily test out new extensions/different approaches
[Affeldt et al., ITP '24]

* So far: Without quantifiers

Properties of Interest

- **Soundness:** If a property interprets as true/false in the DL; it is true in Boolean logic
- **Compositionality:** Composition of negation with conjunction/disjunction; idempotence, commutativity, and associativity of conjunction/disjunction
- **Shadow-lifting:** Gradual improvement in training

The bad news: None of the existing DLs satisfies all of these requirements

[Varnai, Dimarogonas '20]

Conclusion: We might need to provide support for incorporating a range of DLs for different scenarios

Properties of Interest

Compositionality

► **Definition 5** (Commutativity, idempotence, and associativity of \wedge_M). *Given a DL, the interpretation function of conjunction is commutative if for any permutation π of the integers $i \in \{1, \dots, M\}$ we have*

$$\llbracket \wedge_M(p_0, \dots, p_M) \rrbracket_{DL} = \llbracket \wedge_M(p_{\pi(0)}, \dots, p_{\pi(M)}) \rrbracket_{DL}.$$

It is idempotent and associative if we have

$$\begin{aligned} \llbracket \wedge_M(p, \dots, p) \rrbracket_{DL} &= \llbracket p \rrbracket_{DL}, \\ \llbracket \wedge_M(\wedge_M(p_0, p_1), p_2) \rrbracket_{DL} &= \llbracket \wedge_M(p_0, \wedge_M(p_1, p_2)) \rrbracket_{DL}. \end{aligned}$$

Properties of Interest

Soundness

Soundness. Given a DL, an expression e , and a Boolean value b , the DL is sound if:

$$\llbracket e \rrbracket_{DL} = \llbracket b \rrbracket_{DL} \implies \llbracket e \rrbracket_B = b$$

	Soundness
Gödel	Yes
Łukasiewicz	No
Yager	No
Product	Yes
DL2	Yes*
STL	?



In [Ślusarz et al., LPAR '23]

*negation-free fragment

On Soundness

- Clear if it's on closed intervals – see [Ślusarz et al., LPAR '23]
- How to even state soundness for open intervals – i.e., for DL2/STL?

Attempt 1: Add $-\infty$ and $+\infty$ as constants to the domain; keep the previous soundness statement.

- Vacuous proof – no formula evaluates to $-\infty$ or ∞

Not a problem in fuzzy logics – e.g.,
 $\llbracket 3 = 3 \rrbracket_P = 1$

Attempt 2: Keep the open interval intact; re-define soundness in terms of intervals: If the interpretation of the formula e is greater or equal to 0, then $\llbracket e \rrbracket_B = \text{True}$, else $\llbracket e \rrbracket_B = \text{False}$.

- But: Negation is no longer sound. If $\llbracket 3 = 3 \rrbracket_{STL} = 0$, then also $\llbracket \neg (3 = 3) \rrbracket_P = 0$

Exclude 0? =>
Complicates interpretations/ not differentiable

Attempt 3: Use intervals to define truth/ remove negation.

Properties of Interest

Shadow Lifting [Varnai, Dimarogonas '20]

► **Definition 6** (Shadow-lifting property [27]). *The DL satisfies the shadow-lifting property if, for any $\llbracket p \rrbracket_{DL} \neq 0$:*

$$\left. \frac{\partial \left[\left[\Lambda_M(p_0, \dots, p_i, \dots, p_M) \right]_{DL} \right]}{\partial \llbracket p_i \rrbracket_{DL}} \right|_{p_j = p \text{ where } i \neq j} > 0$$

holds for all $0 \leq i \leq M$, where ∂ denotes partial differentiation.

An Overview

	DL2	Gödel	Łuka- siewicz	Yager	Product	STL
Weak Smoothness	Yes*	No	No	No	Yes*	Yes
Shadow-Lifting	Yes	No	No	No	Yes	Yes [Varnai et al. '20]
Scale Invariance	Yes	Yes	No	No	No	Yes [Varnai et al. '20]
Negation	No	Yes	Yes	Yes	Yes	Yes
Idempotence	No	Yes [Krieken et al. '21]	No [Cintula et al. '11]	No [Klement et al. '04]	No [Cintula et al. '11]	Yes [Varnai et al. '20]
Commutativity	Yes	Yes [Krieken et al. '21]	Yes [Krieken et al. '21]	Yes [Krieken et al. '21]	Yes [Krieken et al. '21]	Yes [Varnai et al. '20]
Associativity	Yes	Yes [Krieken et al. '21]	Yes [Krieken et al. '21]	Yes [Krieken et al. '21]	Yes [Krieken et al. '21]	No [Varnai et al. '20]
Soundness	Yes [†]	Yes	No	No	Yes	Yes [†]

* Smoothness in propositional case [†] Negation-free fragment

This Talk

A **common framework** for different differentiable logics that allows to give a uniform semantics:

LDL (Logic of Differentiable Logic)

[Ślusarz et al., LPAR '23]

Goal: A generic framework in which logical/geometric properties of different DLs can be formalized/proven.

A common framework to **compare properties** of different DLs

A **mechanization*** in Coq/MathComp, allowing to easily test out new extensions/different approaches
[Affeldt et al., ITP '24]

* So far: Without quantifiers

This Talk

A **common framework** for different differentiable logics that allows to give a uniform semantics:

LDL (Logic of Differentiable Logic)

[Ślusarz et al., LPAR '23]

Goal: A generic framework in which logical/geometric properties of different DLs can be formalized/proven.

A common framework to **compare properties** of different DLs

A **mechanization*** in Coq/MathComp, allowing to easily test out new extensions/different approaches
[Affeldt et al., ITP '24]

* So far: Without quantifiers

Syntax

```
Inductive flag := def | undef.
```

```
Inductive ldl_type :=  
  Bool_T of flag | Index_T of nat | Real_T | Vector_T of nat | Fun_T of nat & nat.
```

```
1 Definition Bool_T_undef := Bool_T undef.  
2 Definition Bool_T_def := Bool_T def.  
3 Inductive comparison := cmp_le | cmp_eq.  
4  
5 Inductive expr : ldl_type -> Type :=  
6   | ldl_real   : R -> expr Real_T  
7   | ldl_bool   : forall p, bool -> expr (Bool_T p)  
8   | ldl_idx    : forall n, 'I_n -> expr (Index_T n)  
9   | ldl_vec    : forall n, n.-tuple R -> expr (Vector_T n)  
10  | ldl_and    : forall x, seq (expr (Bool_T x)) -> expr (Bool_T x)  
11  | ldl_or     : forall x, seq (expr (Bool_T x)) -> expr (Bool_T x)  
12  | ldl_not    : expr Bool_T def -> expr Bool_T def  
13  | ldl_cmp    : forall x, comparison -> expr Real_T -> expr Real_T -> expr (Bool_T x)  
14  | ldl_fun    : forall n m, (n.-tuple R -> m.-tuple R) -> expr (Fun_T n m)  
15  | ldl_app    : forall n m, expr (Fun_T n m) -> expr (Vector_T n) -> expr (Vector_T m)  
16  | ldl_lookup : forall n, expr (Vector_T n) -> expr (Index_T n) -> expr Real_T.
```

Translation

```
1 Definition type_translation (t : ldl_type) : Type :=
2 match t with
3 | Bool_T x => R
4 | Real_T => R
5 | Vector_T n => n.-tuple R
6 | Index_T n => 'I_n
7 | Fun_T n m => n.-tuple R -> m.-tuple R
8 end.
```

```
Fixpoint stl_translation {t} (e : expr t) : type_translation t :=
  match e in expr t return type_translation t with
  | ldl_and _ (e0 :: s) => let A      := map stl_translation s in
                           let a0    := stl_translation e0 in
                           let a_min := \big[minr/a0]_(i <- A) i in
                           if a_min < 0 then stl_and_lt0 (a0 :: A) else
                           if a_min > 0 then stl_and_gt0 (a0 :: A) else
                           0
  | `~ E1                => - {[ E1 ]}
  | E1 `<= E2            => {[ E2 ]} - {[ E1 ]}
  ... (* see [29] for omitted connectives *)
  end where "{[ e ]}" := (stl_translation e).
```

On the Coq Formalization Overview

File	Contents	L.o.c.
<i>Additions to MATHCOMP libraries</i>		
mathcomp_extra.v	Lemmas iterated min/max, etc.	501
analysis_extra.v	L'Hôpital's rule, Cauchy's MVT (§ 5.3), etc.	820
<i>Generic logic and generic definitions of properties</i>		
ld1.v	LDL syntax and semantics (§ 2), shadow-lifting (Sect. 5.1)	417
<i>Soundness, logical and geometric properties of concrete logics</i>		
d12.v	DL2: logical (§ 4), geometric (§ 5.2)	250
fuzzy.v	Gödel, Łukasiewicz, Yager, product: logical (§ 4), geometric (§ 5.2)	731
st1.v	STL; logical (§ 4), geometric (§ 5.4)	977
<i>Alternative formalisations of logical properties/ soundness using extended reals</i>		
d12_ereal.v	DL2: logical (§ 4.2)	211
st1_ereal.v	STL: logical (§ 4.2)	362
Total		4281

Easily fixable gaps

Translate to corresponding properties of underlying definitions

Helped in closing previous gaps

Helped in finding the right definitions

https://github.com/ndslusarz/formal_LDL

This Talk

Future Work

A **common framework** for different differentiable logics that allows to give a uniform semantics:

LDL (Logic of Differentiable Logic)

[Ślusarz et al., LPAR '23]

Connection of logics with the logics of Lawvere quantale [Bacci et al., '23]

Goal: A generic framework in which logical/geometric properties of different DLs can be formalized/proven.

A common framework to **compare properties** of different DLs

Revised negation

A **mechanization*** in Coq/ MathComp, allowing to easily test out new extensions/ different approaches [Affeldt et al., ITP '24]

* So far: Without quantifiers

In progress

The Question of Quantifiers

- **Motivation:** More expressive/required to internally interpret different forms of robustness
- On **finite domains** such as Bool/ Index n (e.g., [Krieken et al. '21]):
 - Use finitely composed conjunction and disjunction:
Given $\llbracket \tau \rrbracket = \{d_1, \dots, d_n\}$, have $\llbracket \forall x : \tau. e \rrbracket = \llbracket e[x/d_1] \dots e[x/d_n] \rrbracket$
 - Analogously for $\exists x : \tau. e$ and disjunction.
- **Infinite domains**
 - [Fischer et al. '19]: Interpret universally quantified formulae via expectation maximization methods
 - In [Ślusarz et al., LPAR '23]: Suggestion of a language-independent quantifier
$$\llbracket \forall x : \tau. e \rrbracket_L^{N,Q,\Gamma} = \mathbb{E}_{\min} [(\lambda y. \llbracket e \rrbracket^{N,Q,\Gamma[x \rightarrow y]}) (Q[x])]$$
where $\mathbb{E}_{\min} [g(\mathbf{X})] = \lim_{\gamma \rightarrow 0} \int_{x \in \mathbb{B}_{x_{\min}}^\gamma} p_X(x) g(x) dx$

Practical algorithm for computing it?

- ... to be continued

Any Questions?

Published Papers

[Logic of differentiable logics: Towards a uniform semantics of DL](https://arxiv.org/pdf/2303.10650)

<https://arxiv.org/pdf/2303.10650>

Natalia Ślusarz, Ekaterina

Komendantskaya, Matthew L Daggitt,
Robert Stewart, KS (LPAR '23)

[Taming Differentiable Logics with Coq Formalisation](#)

Reynald Affeldt,

Alessandro Bruni, Ekaterina

Komendantskaya, Natalia Ślusarz, KS

(ITP '24)

Coq Development

https://github.com/ndslusarz/formal_L_DL

Parts of the project are looking for a PhD student:

Vacancy: PhD in Computer Science

Title: Formal Verification of AI Interfaces

Advisors: Ekaterina Komendantskaya (Southampton University, UK), Alessandro Bruni (IT University of Copenhagen, Denmark), Reynald Affeldt (AIST, Japan)

Start Date: As soon as the right candidate is found

Location: Southampton University, UK; with collaborative visits involving researchers at AIST, Japan and IT University of Copenhagen, Denmark.

September 9,
15:00, ITP