

Explainable Deep Learning and its potential use in TSF

Tirtharaj Dash

University of Cambridge, UK

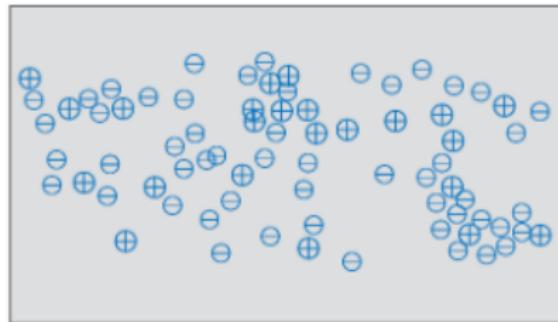
September 28, 2024



Today

- ▶ Explainability and Deep Learning
- ▶ “Logically” explainable Deep Neural Nets
- ▶ Application in TSF, esp. Healthcare

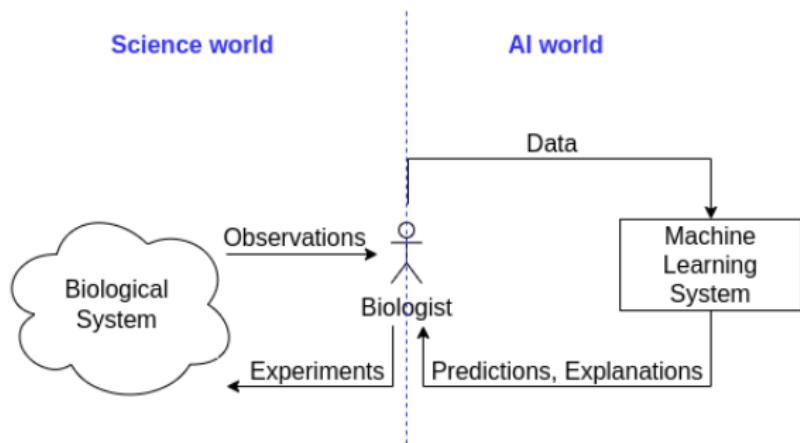
Explainability and Deep Neural Nets



A classical ML problem: classification.

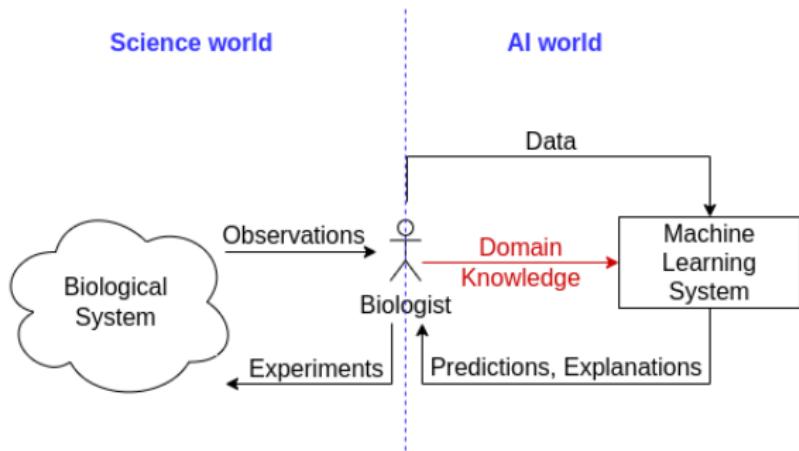
Explainability and Deep Neural Nets

Classical ML:



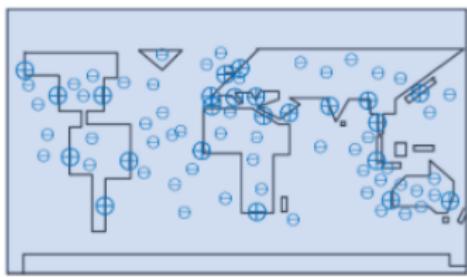
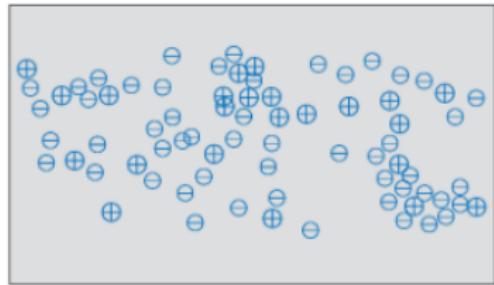
Explainability and Deep Neural Nets

Human-in-the-Loop ML:



Background knowledge comes from domain-expert(s).

Explainability and Deep Neural Nets



The \oplus points are the port cities, inferable from boundaries.

Explainability and Deep Neural Nets

Inclusion of domain-knowledge into deep neural networks significantly improves its predictive performance.

DNN Type	Comparative Performance (with domain-knowledge)		
	Better	Same	Worse
MLP	71	0	2
GNN	63	9	1

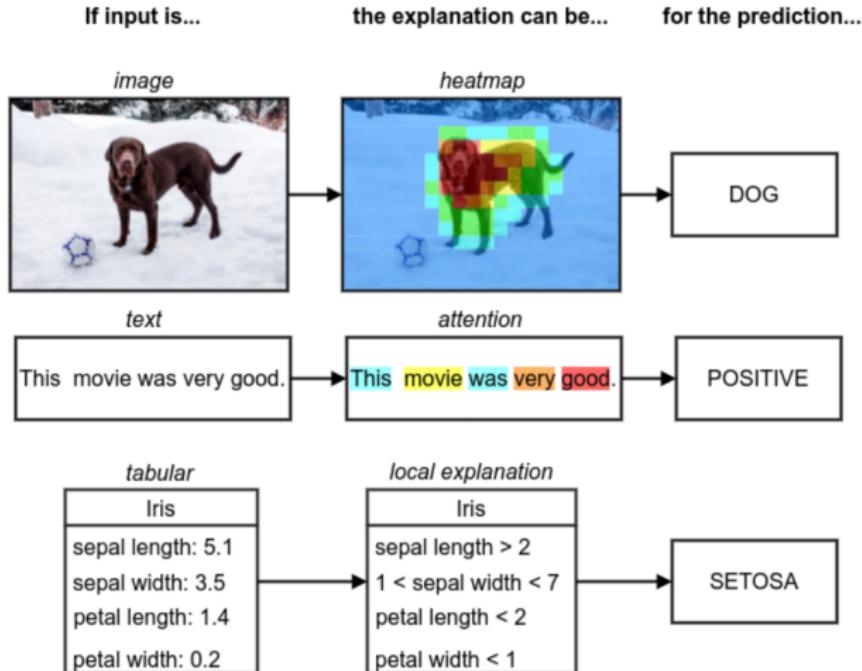
Dash et al.: *MLJ* (2021, 2022, 2023), *ILP* (2018, 2021), *Sci.Rep.* (2022).

Dash, *PhD Thesis*, BITS Pilani, 2022.

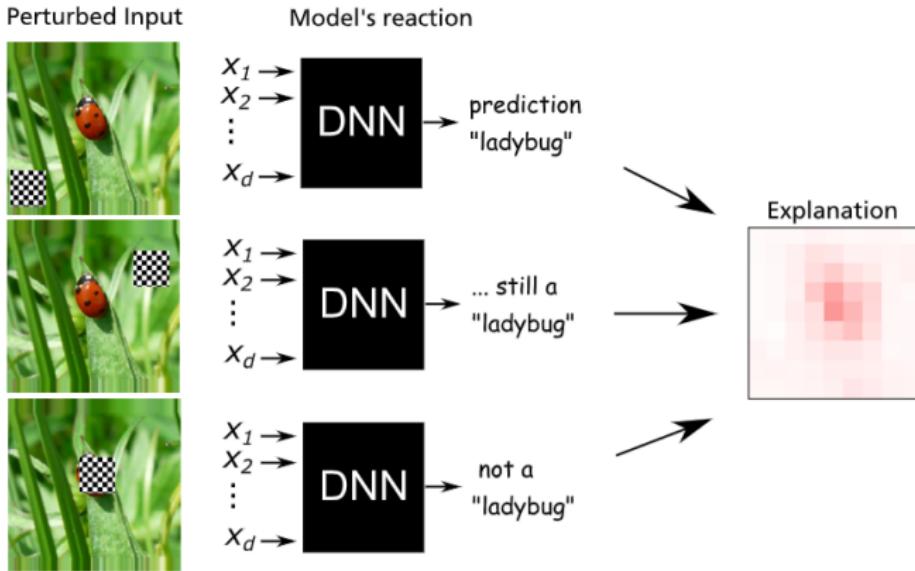
Explainability and Deep Neural Nets

Explainability or Interpretability is the concept that a machine learning model and its output can be explained in a way that “makes sense” to a human being at an acceptable level.

Explainability and Deep Neural Nets

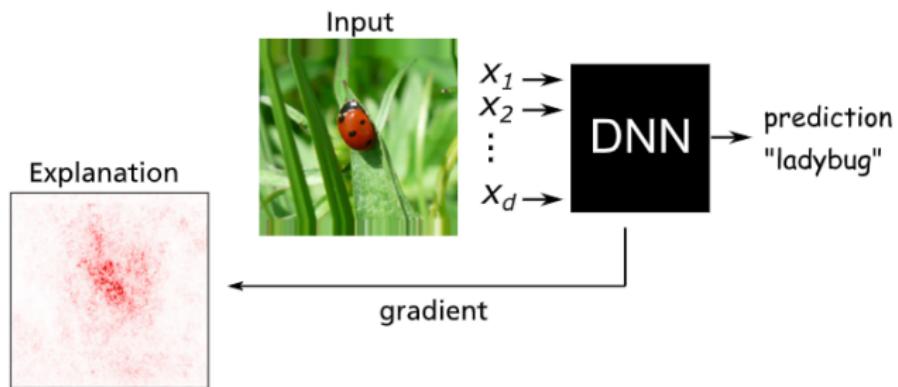


Explainability and Deep Neural Nets



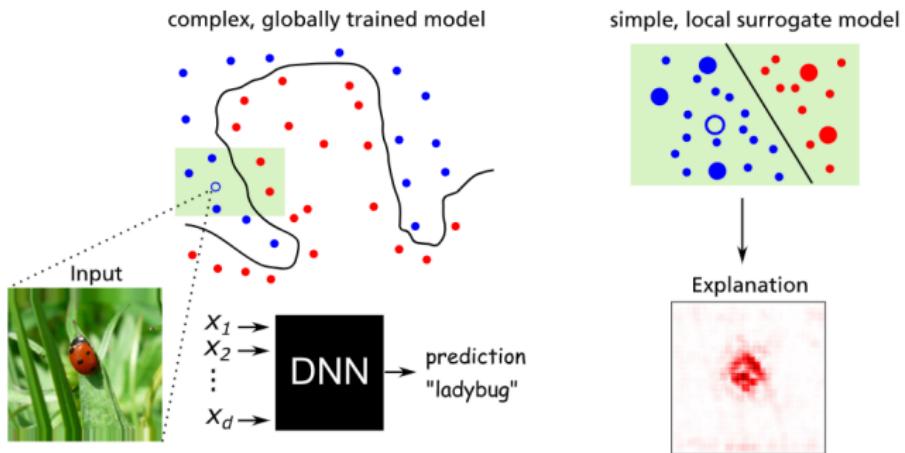
Perturbation-based explanation generation.

Explainability and Deep Neural Nets



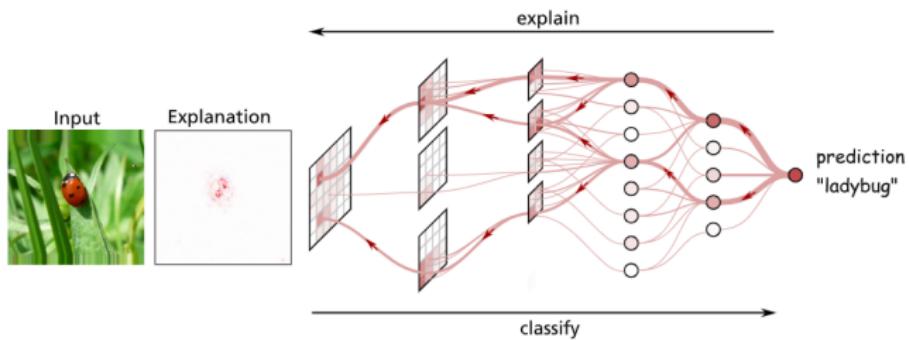
Gradient-based explanation generation.

Explainability and Deep Neural Nets



Surrogate-based explanation generation.

Explainability and Deep Neural Nets



Relevance propagation based explanation generation.

Explainability and Deep Neural Nets

So, what did we see?

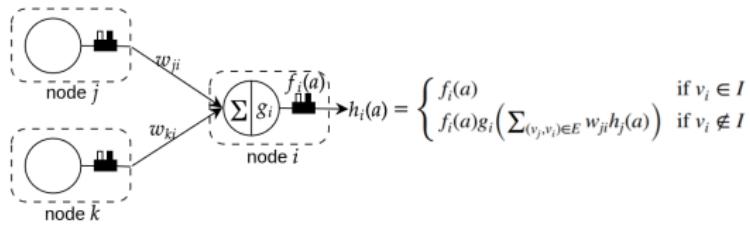
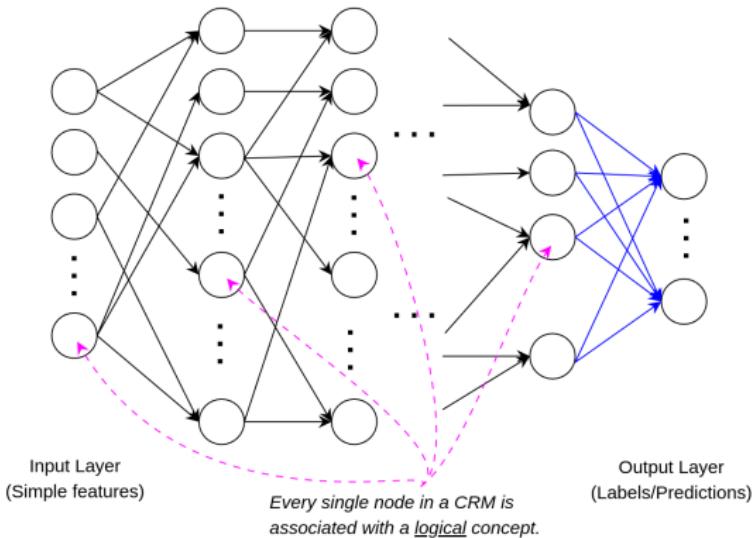
Explainability and Deep Neural Nets

Methods to explain “black-box” deep networks *post hoc*.

Explainability and Deep Neural Nets

- Q1. But, is a deep neural network “truly” explainable?
- Q2. Can we build deep networks that are “explainable by design”?

Logically explainable DNNs



CRM nodes as Gated nodes

Logically explainable DNNs

Relational Features:

A relational feature takes a clausal form:

$$C : \forall X \ (p(X) \leftarrow \exists \mathbf{Y} \ Body(X, \mathbf{Y}))$$

or,

$$C : (p(X) \leftarrow Body(X, \mathbf{Y}))$$

Here,

$p(X)$: Head literal

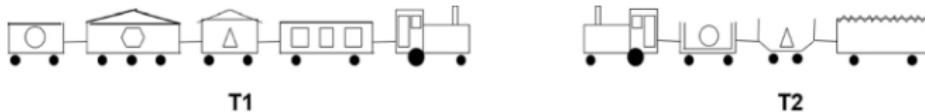
$Body(X, \mathbf{Y})$: Conjunction of body literals

Assumption: C is not self-recursive. We call C a “feature-clause”.

Logically explainable DNNs

Feature clauses:

Let's look at the classic trains problem:



Some feature-clauses for trains:

$C_1 : p(X) \leftarrow (\text{has_car}(X, Y), \text{short}(Y))$

$C_2 : p(X) \leftarrow (\text{has_car}(X, Y), \text{short}(Y), \text{closed}(Y))$

$C_3 : p(X) \leftarrow (\text{has_car}(X, Y), \text{has_car}(X, Z), \text{short}(Y), \text{closed}(Z))$

The predicates *has_car*/2, *short*/1, *closed*/1, etc. are defined as part of the background knowledge (B) about trains.

Logically explainable DNNs

Feature functions:

A feature function is defined, for $X = a$ as:

$$f_{C,B}(a) = \begin{cases} 1 & \text{if } B \cup (C\{X/a\}) \models p(a) \\ 0 & \text{otherwise} \end{cases}$$

Simply, for a feature-clause C_i , we refer to the corresponding feature-function as $f_i(X)$.

Example:

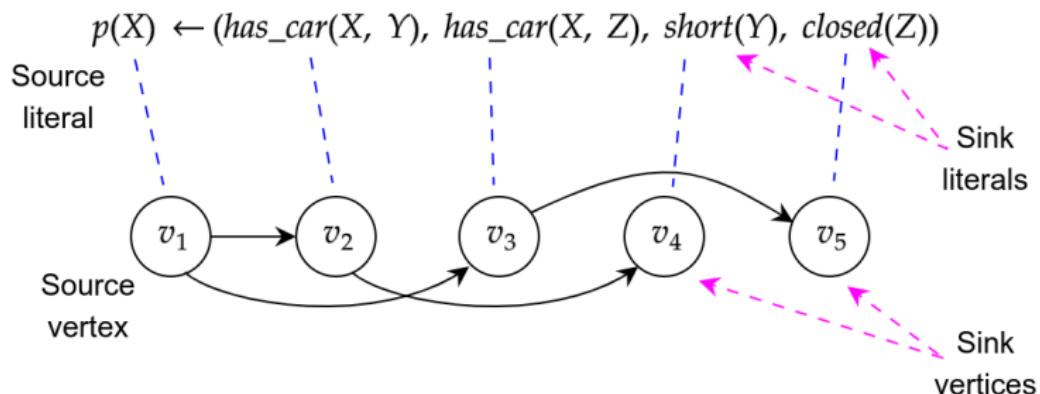


Some feature functions are: $f_1(t_1) = 1$, $f_2(t_1) = 1$, $f_2(t_2) = 0$.

Logically explainable DNNs

Ordered Clause:

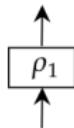
We impose an ordering of the literals in a clause. If C is a clause of the form $\lambda_1 \leftarrow \lambda_2, \dots, \lambda_k$, then the ordered clause is:
 $\langle C \rangle = \langle \lambda_1, \lambda_2, \dots, \lambda_k \rangle$.



Logically explainable DNNs

ρ -derivation of feature-clauses (Composition):

Example 1:

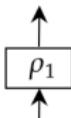
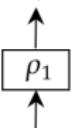
$$p(X) \leftarrow \text{has_car}(X, Y), \text{short}(Y), \\ \text{has_car}(X, Z), \text{closed}(Z), \\ Y = Z$$

$$p(X) \leftarrow \text{has_car}(X, Y), \text{short}(Y), \\ \text{has_car}(X, Z), \text{closed}(Z)$$


C1: $p(X) \leftarrow \text{has_car}(X, Y), \text{short}(Y)$

C2: $p(X) \leftarrow \text{has_car}(X, Z), \text{closed}(Z)$

Logically explainable DNNs

Example 2:

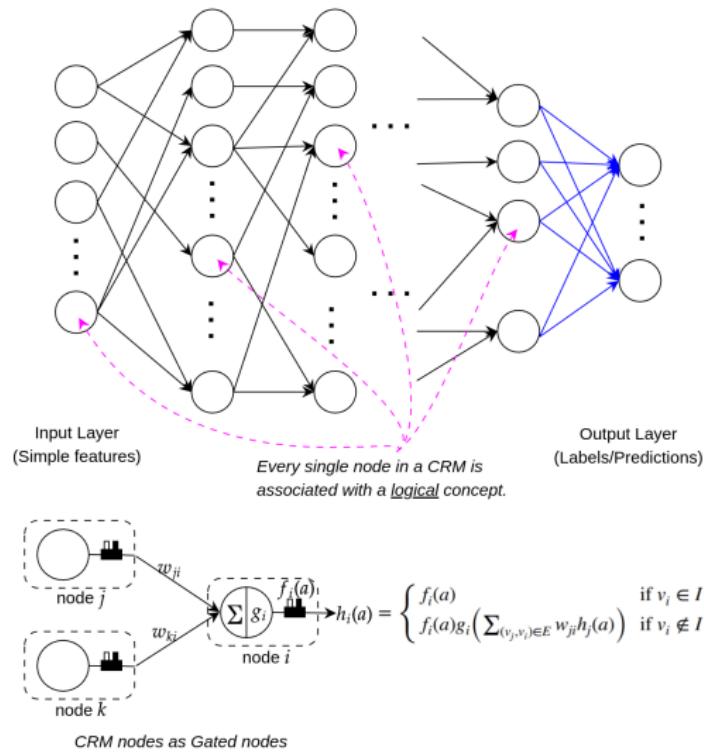
$$p(X) \leftarrow \text{has_car}(X, U), \text{has_car}(X, V), \text{smaller}(U, V), \\ \text{has_car}(X, Y), \text{short}(Y), U = V, U = Y$$

$$p(X) \leftarrow \text{has_car}(X, U), \text{has_car}(X, V), \text{smaller}(U, V), \\ \text{has_car}(X, Y), \text{short}(Y), U = V$$

$$p(X) \leftarrow \text{has_car}(X, U), \text{has_car}(X, V), \text{smaller}(U, V), \\ \text{has_car}(X, Y), \text{short}(Y)$$


C1: $p(X) \leftarrow \text{has_car}(X, Y), \text{short}(Y)$

C2: $p(X) \leftarrow \text{has_car}(X, U), \text{has_car}(X, V), \text{smaller}(U, V)$

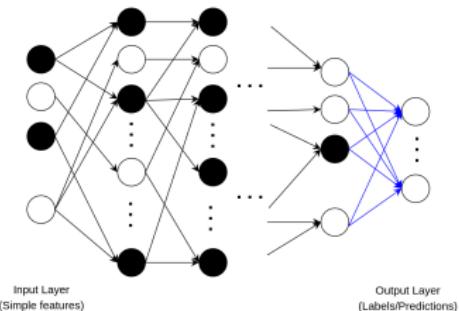
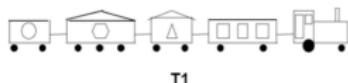
Logically explainable DNNs

Compositional Relational Machines (CRMs):

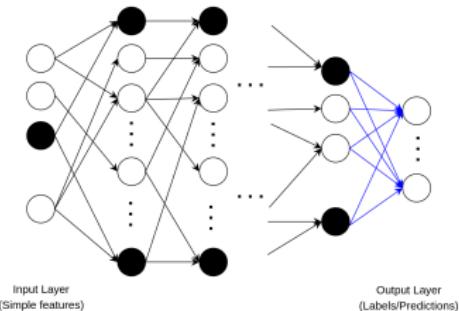
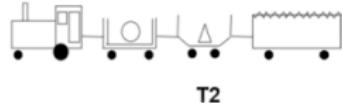


Logically explainable DNNs

Relational instance 1:

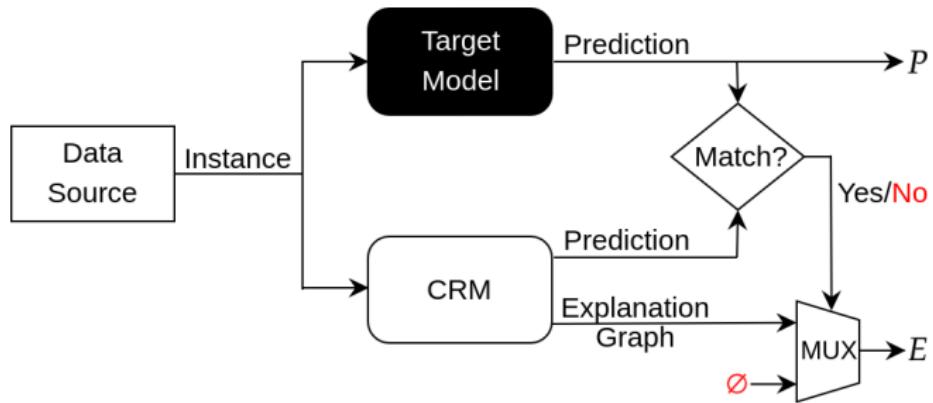


Relational instance 2:



Logically explainable DNNs

Evaluation: (a) Predictive fidelity, (b) Explanatory fidelity



Explanation: Constructed by back-tracing the top activations in each layer of the deep neural network.

Logically explainable DNNs

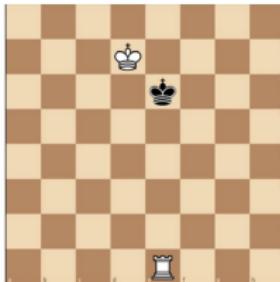
(A) Synthetic datasets (Trains and Chess)

- ▶ Target theory (model) known

Trains



Chess



- ▶ Results

Dataset	Fidelity			
	CRM		Baseline	
	Pred.	Expl.	Pred.	Expl.
Trains	1.0	1.0	0.5	0.4
Chess	1.0	0.9	0.7	0.7

Logically explainable DNNs

- ▶ Some explanations generated by the CRM:



$p(X) \leftarrow$
has_car(X, A), *short*(A),
has_car(X, B), *closed*(B),
 $A = B$

$p(X) \leftarrow$
has_car(X, A'), *short*(A'),
has_car(X, B'), *closed*(B'),

$p(X) \leftarrow$
has_car(X, A''), *short*(A'')
 $p(X) \leftarrow$
has_car(X, A''), *closed*(A'')

Train t_1
(With the substitution $\{X/t1\}$)



Board $(d, 7, e, 1, e, 6)$

$p((A, B, C, D, E, F)) \leftarrow$
lt(A, C), $C = E$, *adj*(B, F)

$p((A, B, C, D, E, F)) \leftarrow$
 $C = E$, *adj*(B, F)

$p((A, B, C, D, E, F)) \leftarrow$
lt(A, C)

$p((A, B, C, D, E, F)) \leftarrow$
lt(A, C), $C = E$

$p((A, B, C, D, E, F)) \leftarrow$
 $C = E$

$p((A, B, C, D, E, F)) \leftarrow$
adj(B, F)

(With the substitution $\{A/d, B/7, \dots, F/6\}$)

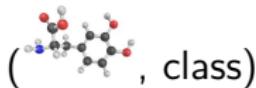
Target theory: Train X has a car Y and
 Y is short and closed.

Target theory: White Rook and Black King are
on the same file (column).

Logically explainable DNNs

(B) Real datasets (drug design: NCI GI50; $n = 10$)

- ▶ Target theory is not known. BotGNNs are used as the target models [Dash et al., MLJ, 2022].

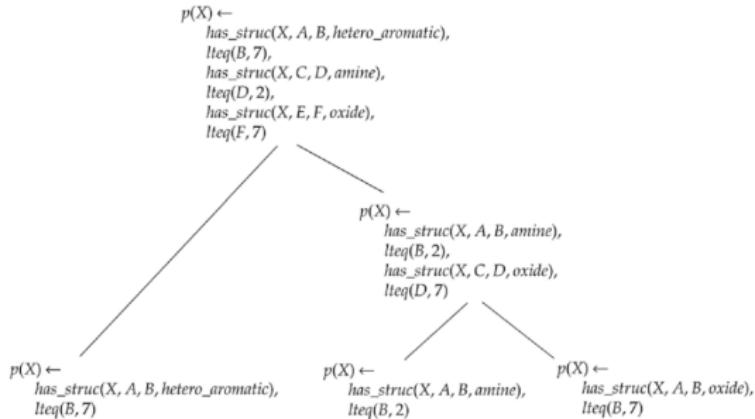
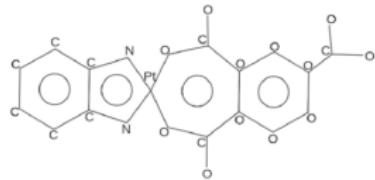


- ▶ Results: Predictive fidelity

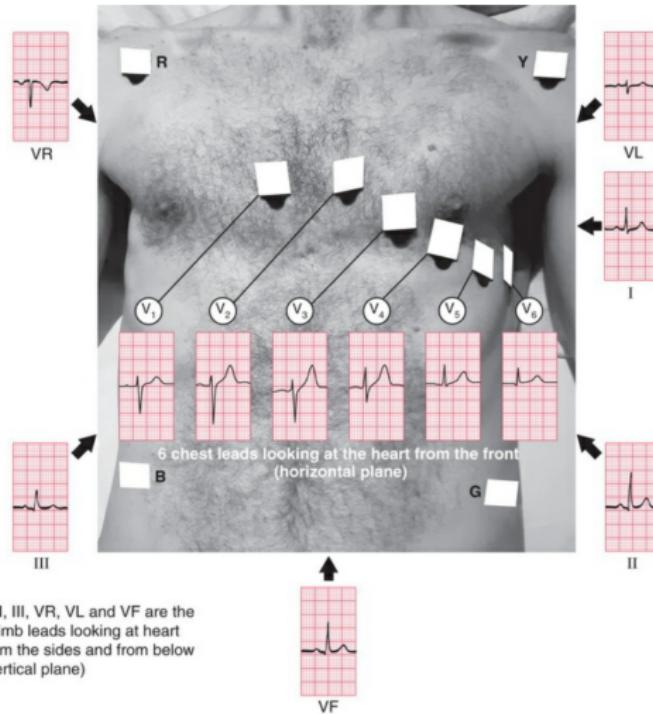
Dataset	CRM	Baseline
786_0	0.77	0.53
A498	0.79	0.59
A549_ATCC	0.85	0.63
ACHN	0.73	0.58
BT_549	0.78	0.51
CAKI_1	0.81	0.69
CCRF_CEM	0.82	0.68
COLO_205	0.77	0.53
DLD_1	0.90	1.00
DMS_114	0.89	0.91
Avg.	0.81 (0.05)	0.66 (0.17)

Logically explainable DNNs

- ▶ An explanation generated by CRM:

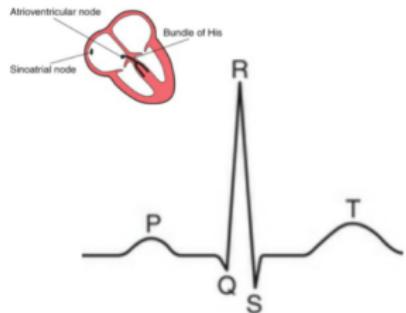


XDL in Healthcare TSF



Lead positions for a 12-lead ECG with 12 views of the heart

XDL in Healthcare TSF

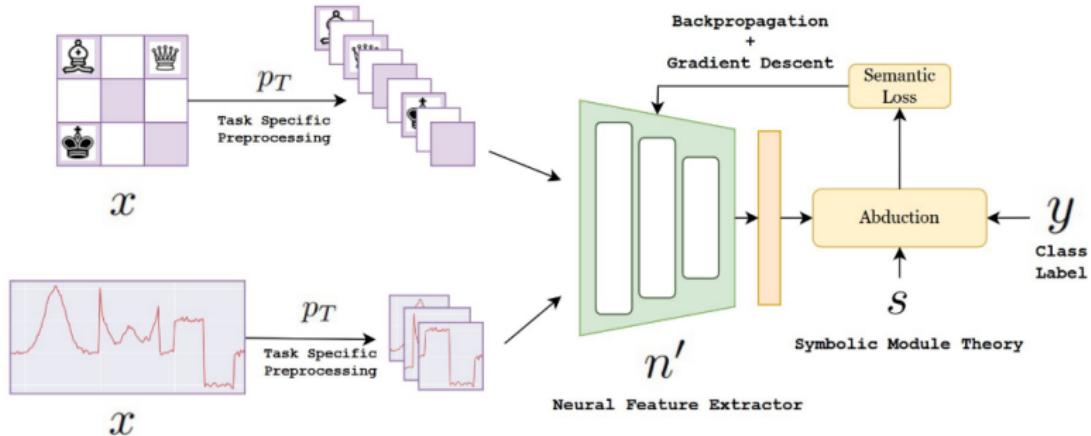


R	Rate	What is the rate (measured in beats per minute [bpm])?
R	Rhythm	What is the rhythm?
P	P wave	Is there one P wave before every QRS complex?
W	Width	Is the width of the QRS complex normal (< 3 small squares)?
Q	Q wave	Are there any deep Q waves present?
S	ST segment	Is there ST segment depression or elevation?
T	T wave	Are there any abnormal inverted (upside down) T waves?

ECG complex - 1 heartbeat

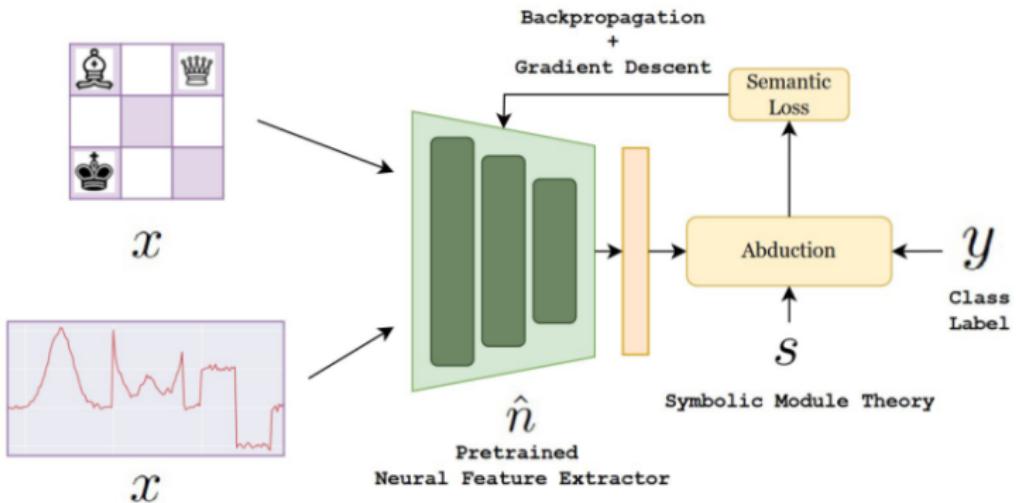
How doctors see it?

XDL in Healthcare TSF



Shah et al.: NEUROLOG, arXiv, 2022.

XDL in Healthcare TSF



Thank you!

tirtharajdash.github.io