

"Formal XAI": Can We **Formally** Explain ML Models?

Shahaf Bassan

### ML and NNs are used extensively







# NNs are sensitive to adversarial attacks



# NNs are sensitive to adversarial attacks





Cup(16.48%) Soup Bowl(16.74%)

Bassinet(16.59%)
Paper Towel(16.21%)

### **Neural Network** Verification



### The classic property: Adversarial Robustness



**<u>Property:</u>** Verify that there isn't an **adversarial perturbation** around:

$$\mathcal{E}=0.008$$



### How hard is NN-verification?



#### **NN-verification is NP-complete!** (Katz et al., 2017)

### How hard is NN-verification?

- Marabou (Katz et al.)
- Beta- Crown (Wang et al)
- DeepPoly (Singh et al.)
- Al2 (Gehr et al.)
- Prima (Muller et al.)
- Verinet (Henriksen et al.)
- MN-BAB (Ferrari et al.)
- Nnenum (Bak et al.)

verifier

2017: <100 neurons

### How hard is NN-verification?

- Marabou (Katz et al.)
- Beta- Crown (Wang et al)
- DeepPoly (Singh et al.)
- Al2 (Gehr et al.)
- Prima (Muller et al.)
- Verinet (Henriksen et al.)
- MN-BAB (Ferrari et al.)
- Nnenum (Bak et al.)





2023-2025: ~10,000,000 neurons

\*Still not applicable on SOTA\*

### From Formal Verification to Formal Explainability

### Neural Networks are Black-Boxes





### Explainable AI (XAI)

# Tools and frameworks for explaining the decisions made by ML models.













Explainer:



"Labrador"





Explainer:



"Labrador"





#### "Wolf" (and not Husky)







Explainer:



### "Wolf" (and not Husky)

### Can we trust XAI tools?

# XAI tools are often <u>heuristic</u>, and do not provide <u>formal guarantees</u>.

### Can we trust XAI tools?

# If we can't trust the **<u>explainer</u>**, we can't trust the **<u>model</u>**.



### We hence want to produce <u>formal</u> <u>and provable</u> explanations.

### How can we formally define an explanation?



### Sufficient Reason/Abductive Explanation



# $\forall \mathbf{z} \in B_p^{\epsilon_p}(\mathbf{x}). \quad [\operatorname*{arg\,max}_j f(\mathbf{x}_S; \mathbf{z}_{\bar{S}})_j = \operatorname*{arg\,max}_j f(\mathbf{x})_j]$

Sufficient Reason:





### Abductive Explanations/Sufficiency

### **DNN-verification** can be used to verify



if

#### is a <u>sufficient reason</u>.

### Minimal Sufficient Reasons

 Smaller sufficient reasons are more meaningful.
 We are interested in finding <u>subset minimal</u> or preferably <u>cardinally minimal</u> explanations.

Towards Formal XAI: Formally Approximate Minimal Explanations of Neural Networks

**Tacas, 2023** Shahaf Bassan, Guy Katz

Towards Formal XAI: Formally Approximate Minimal Explanations of Neural Networks

**Tacas, 2023** Shahaf Bassan, Guy Katz



Towards Formal XAI: Formally Approximate Minimal Explanations of Neural Networks

**Tacas, 2023** Shahaf Bassan, Guy Katz

Provable subset minimal sufficient reasons

Provable <u>approximation</u> of cardinally minimal sufficient reasons

Towards Formal XAI: Formally Approximate Minimal Explanations of Neural Networks

**Tacas, 2023** Shahaf Bassan, Guy Katz

Provable subset minimal sufficient reasons

Provable <u>approximation</u> of cardinally minimal sufficient reasons














### **Subset** minimal sufficient reasons

Attempt to ``free" features until converging to a subset minimal explanation.



### A key concern: **high computational complexity**

## A key concern: high computational complexity

### Each verification query is NP-Hard, and we require a linear number of queries only for a subset minimal explanation!

### How can we **<u>speed up</u>** this process?

### How can we **<u>speed up</u>** this process?

Explaining, Fast and Slow: Abstraction and Refinement of Provable Explanations

#### ICML 2025 (To appear)

Shahaf Bassan\*, Yizhak Elboher\*, Tobias Ladner\*, Matthias Althof, Guy Katz

### How can we speed up this process?

We suggest an algorithm composed of **two** main aspects:

### How can we speed up this process?

We suggest an algorithm composed of **two** main aspects:



### How can we speed up this process?

We suggest an algorithm composed of **two** main aspects:



2. Refinement

**Original model** 



**Original model** 







































This requires an additional procedure: refinement!









Gradually **refine** the abstract network (increase its size) and find an explanation



Gradually **refine** the abstract network (increase its size) and find an explanation



Gradually **refine** the abstract network (increase its size) and find an explanation

We "free" features in the abstract model, then refine and "free" more features, and so on...



Gradually **refine** the abstract network (increase its size) and find an explanation

We "free" features in the abstract model, then refine and "free" more features, and so on...



Gradually **refine** the abstract network (increase its size) and find an explanation

We "free" features in the abstract model, then refine and "free" more features, and so on...

Eventually, converges to a **minimal explanation**.

### An example



### More examples

Network size  $\rho = 10\% = 20\% = 30\% = 40\% = 50\% = 60\% = 70\% = 80\% = 90\% = 100\%$ 



### If we **<u>relax</u>** sufficiency, can we <u>scale even more</u>?

## If we **relax** sufficiency, can we **scale even more**?

# Explain Yourself, Briefly! Self-Explaining Neural Networks with Concise Sufficient Reasons

**ICLR 2025** Shahaf Bassan, Ron Eliav, Shlomit Gur
Up until now, we discussed **<u>post-hoc</u>** explanations, that are obtained <u>after</u> training.

Can <u>training time-intervention</u> help with our scalability challenges?

#### A model that **explains itself**

#### A model that **explains itself**



# How do we train a model to give us an explanation that is both **sufficient** and **small**?

How do we train a model to give us an explanation that is both **sufficient** and **small**?

#### With Sufficient Subset Training (SST)

#### The dual propagation in Sufficient Subset Training



#### $L_{\theta}(h) := L_{pred}(h) + \lambda L_{faith}(h) + \xi L_{card}(h)$

$$L_{\theta}(h) := L_{pred}(h) + \lambda L_{faith}(h) + \xi L_{card}(h)$$

(Standard) Prediction loss: Optimize for accuracy  $L_{CE}(h_1(\mathbf{x}), t)$ 

#### The dual propagation in Sufficient Subset Training



#### $L_{\theta}(h) := L_{pred}(h) + \lambda L_{faith}(h) + \xi L_{card}(h)$

$$L_{\theta}(h) := L_{pred}(h) + \lambda L_{faith}(h) + \xi L_{card}(h)$$

Faithfulness loss: Optimize for sufficiency  $L_{CE}(h_1(\mathbf{x}_S; \mathbf{z}_{\bar{S}}), \arg \max_j h_1(\mathbf{x})_j),$ 

#### The dual propagation in Sufficient Subset Training



#### $L_{\theta}(h) := L_{pred}(h) + \lambda L_{faith}(h) + \xi L_{card}(h)$

$$L_{\theta}(h) := L_{pred}(h) + \lambda L_{faith}(h) + \xi L_{card}(h)$$

Cardinality loss: Optimize for minimal cardinality  $||h_2(\mathbf{x})||_1$ 

$$L_{\theta}(h) := L_{pred}(h) + \lambda L_{faith}(h) + \xi L_{card}(h)$$

Cardinality loss: Optimize for minimal cardinality  $||h_2(\mathbf{x})||_1$ 

#### How do we perform the <u>masking</u>?



#### What about the masking?

#### How do we perform the <u>masking</u>?



1. **<u>Baseline masking</u>** - fix some baseline to the complementary.

- 1. **<u>Baseline masking</u>** fix some baseline to the complementary.
- 2. **Probabilistic masking** sample values to the complementary from some distribution.

- 1. **<u>Baseline masking</u>** fix some baseline to the complementary.
- 2. **<u>Probabilistic masking</u>** sample values to the complementary from some distribution.
- 3. **<u>Robust masking</u>** perform a gradient attack over the complementary features.

# We ran experiments on both image domains (IMAGENET, CIFAR10, MNIST) and language domains (IMDB, SNLI)

#### We ran comparisons to post-hoc methods



#### We ran comparisons to post-hoc methods



# We ran comparisons to post-hoc methods, and ablations



#### Probabilistic-Sufficiency (negative):

I *found* this movie <u>really hard</u> ... <u>wandering</u> off the tv ...

#### **Baseline-Sufficiency (negative):**

I <u>found this</u> movie really hard ... wandering <u>off</u> the tv ... <u>Don't</u> bother with it.

1. Improved faithfulness (more subsets verified as sufficient)

1. Improved faithfulness (more subsets verified as sufficient)

2. Improved conciseness (smaller subsets)

1. Improved faithfulness (more subsets verified as sufficient)

- 2. Improved conciseness (smaller subsets)
- 3. More efficient (inherent sufficient reasons)

1. Improved faithfulness (more subsets verified as sufficient)

- 2. Improved conciseness (smaller subsets)
- 3. More efficient (<u>inherent</u> sufficient reasons)
  - 4. Comparable predictive performance

# classification





#### **<u>Reactive Systems</u>** are more complex





#### Formal Explanations in **Reactive Systems**

Formally Explaining Neural Networks within Reactive Systems

**FMCAD, 2023** Shahaf Bassan\*, Guy Amir\*, Davide Corsi, Idan Refaeli, Guy Katz **Best paper runnerup** 

#### A naive solution: Encode everything together


### The Problem with the naive solution

#### Growth of neural network: <u>exponential</u> <u>blow up</u> in computation time.

### Our solution

# We suggest algorithms that avoid the exponential blow-up, but are still optimal.

# Formal Explanations in **<u>Reactive Systems</u>**



### Let's move to some more theoretical stuff

# The computational complexity of finding explanations

# The computational complexity of finding explanations

- Local vs. Global Interpretability: A Computational Complexity Perspective ICML 2024 (Spotlight) Shahaf Bassan, Guy Amir, Guy Katz
- What makes an Ensemble (Un) Interpretable?
  ICML 2025 (To appear)
  Shahaf Bassan, Guy Amir, Meirav Zehavi, Guy Katz
- On the Computational Tractability of the (Many) Shapley Values Al'STATS 2025
   Reda Marzouk\*, Shahaf Bassan\*, Guy Katz, Colin De La Higuera
- Hard to Explain: On the Computational Hardness of In-Distribution Model Interpretation ECAI 2024 Guy Amir\*, Shahaf Bassan\*, Guy Katz

Factors that influence the complexity:

#### Factors that influence the complexity:



#### Factors that influence the complexity:







1. Formal XAI lets us compute explanations that are provably correct.

- 1. Formal XAI lets us compute explanations that are provably correct.
- 2. We can do this using NN verification

- 1. Formal XAI lets us compute explanations that are provably correct.
- 2. We can do this using NN verification
- 3. Abstraction-refinement can improve efficiency

- 1. Formal XAI lets us compute explanations that are provably correct.
- 2. We can do this using NN verification
- 3. Abstraction-refinement can improve efficiency
- 4. Self-explaining neural networks can scale even more

- 1. Formal XAI lets us compute explanations that are provably correct.
- 2. We can do this using NN verification
- 3. Abstraction-refinement can improve efficiency
- 4. Self-explaining neural networks can scale even more
- 5. Understanding the complexity of finding explanations is an important theoretical aspect of this area.

1. Verifying alternative forms of explanations



- 1. Verifying alternative forms of explanations
- 2. Methods for improving scalability



- 1. Verifying alternative forms of explanations
- 2. Methods for improving scalability



- 1. Verifying alternative forms of explanations
- 2. Methods for improving scalability
- 3. Additional work on computational complexity



- 1. Verifying alternative forms of explanations
- 2. Methods for improving scalability
- 3. Additional work on computational complexity
- 4. Formal certification during training



