

# Explaining, Fast and Slow

## Abstraction and Refinement of *Provable* Explanations

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

Hebrew University of Jerusalem (HUJI),  
Technical University of Munich (TUM)

An efficient algorithm to find *minimal* and *sufficient* explanations in neural networks with *provable guarantees*. It *abstracts* the model to a smaller version, then *refines* it by growing its size, ensuring *sufficiency* and *minimality*.

What does *provably sufficient* mean?

$$f\left(\text{Image of dog}\right) = \text{Dog} \quad f\left(\text{Image of dog with UK flag}\right) = \text{Dog}$$

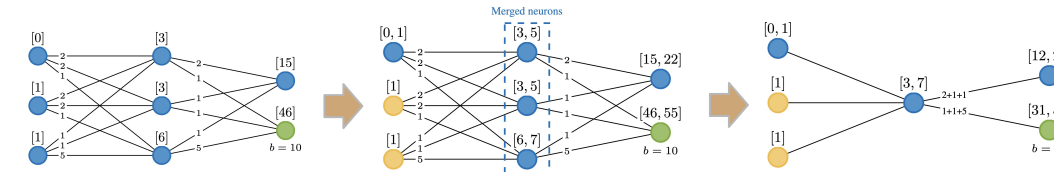
$$f\left(\text{Image of dog}\right) \stackrel{?}{=} \text{Dog}$$

What does *provably minimal* mean?

Any subset of  is not sufficient!

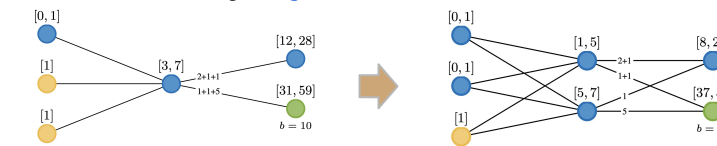
Previous algorithms struggle with **high compute** (more than **NP-hard**) especially for **large models**.

**Our method:** *Abstract* the model (merge neurons and build a **much smaller one**).



**We prove:** A **sufficient explanation** for the **abstract (small)** model is **sufficient** for the **original**!  
But a **minimal** one **might not stay minimal**.

This means we need to iteratively *refine* the model (increase its size) to reach **minimality**.



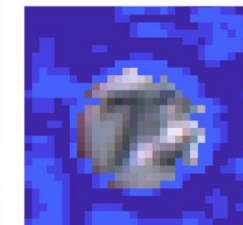
Network size  $\rho =$  10% 20% 30% 40% 50% 60% 70% 80% 90% 100%



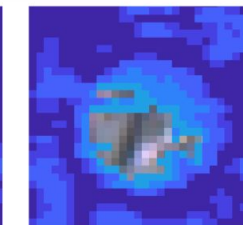
(a) Original image



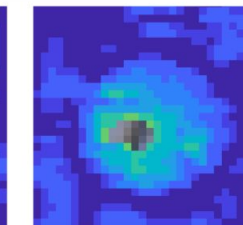
(b)  $\rho \leq 10\%$



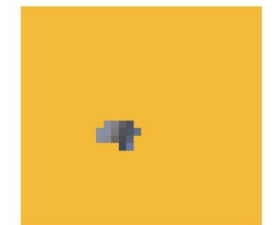
(c)  $\rho \leq 30\%$



(d)  $\rho \leq 50\%$



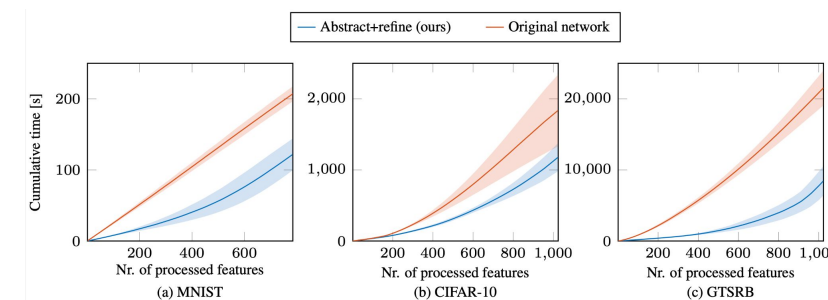
(e)  $\rho \leq 80\%$



(f) Original network

Explanation size increases

Network size  $\rho$  and computation time increases



- Significantly *faster* than previously suggested algorithms.
- Significantly *smaller* explanations compared to previous approaches.
- Scales to *larger* models.

