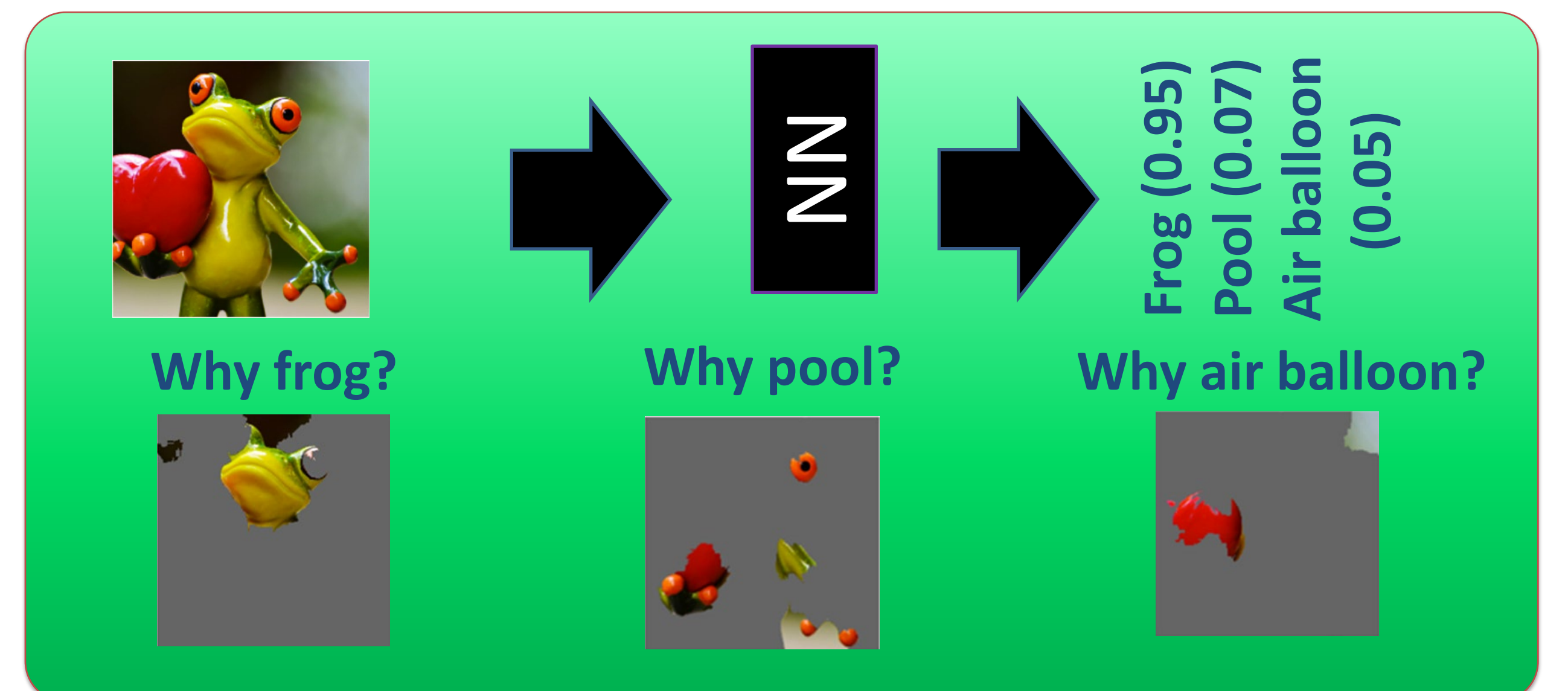
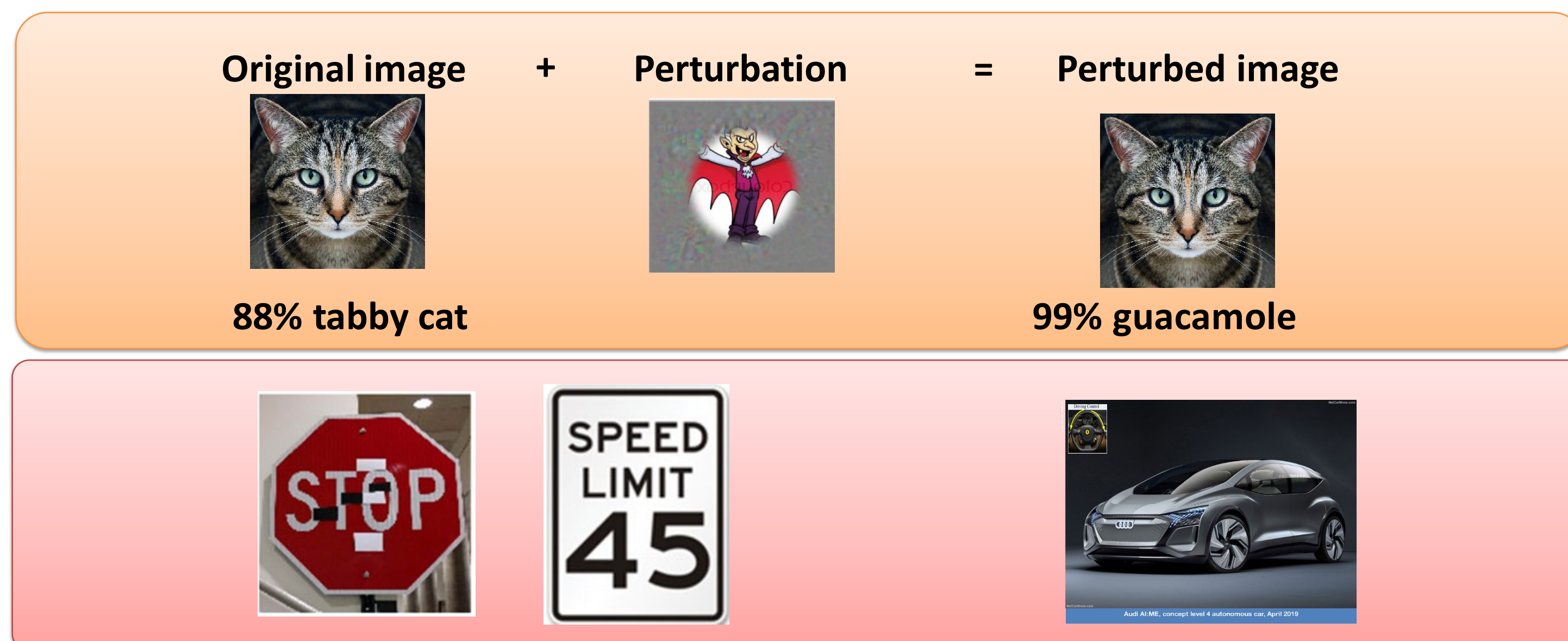


# On Relating Explanations and Adversarial Examples

Alexey Ignatiev, Nina Narodytska, and Joao Marques Silva

## Adversarial examples and explanations in neural networks

[Szegedy et al., Goodfellow et al., Ribeiro et al.]

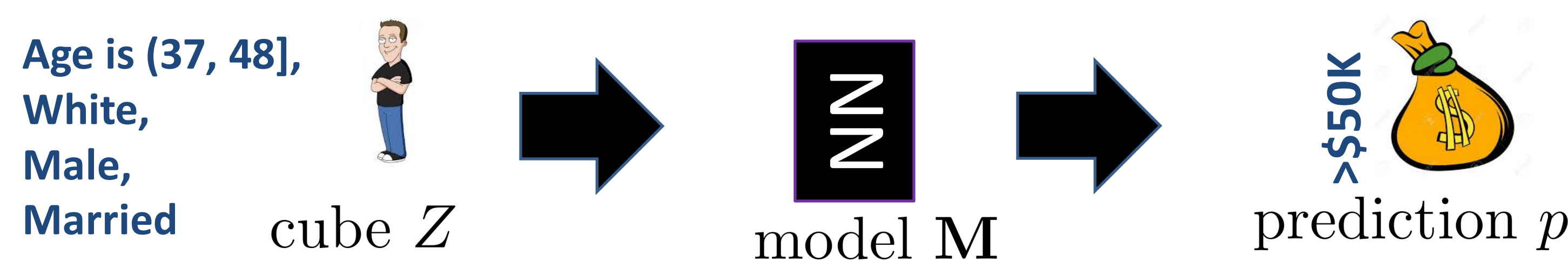


Is there a connection between **adversarial examples** and **explanations**?

## Logic-based approach to attacks and explanations

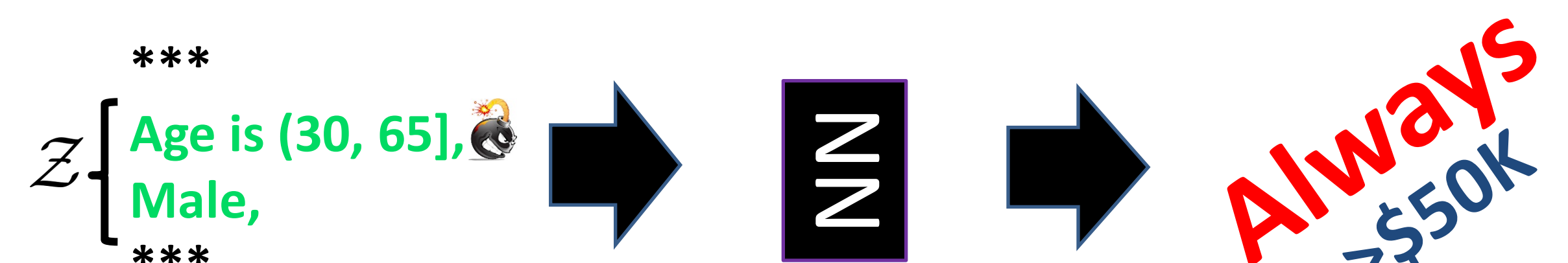
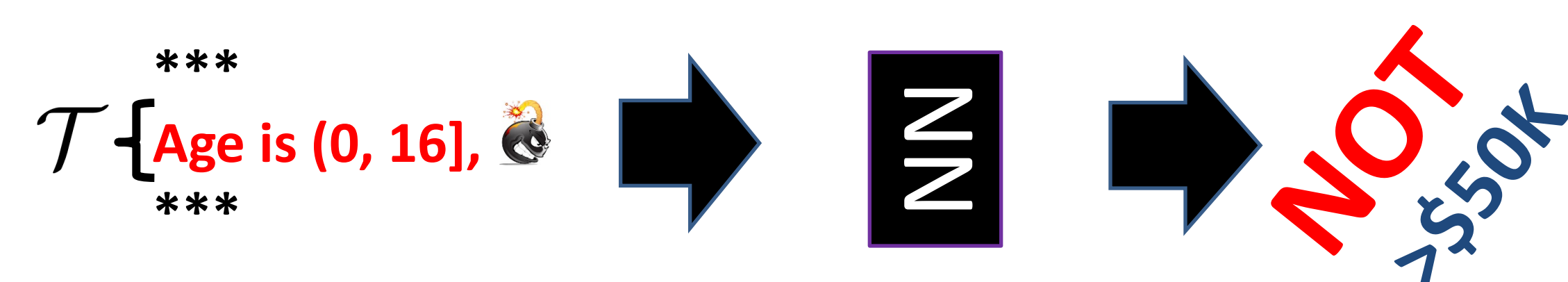
[Ignatiev et al., Shih et al.]

Given a classifier  $M$  (in a logic encoding) and a prediction  $p$ ,

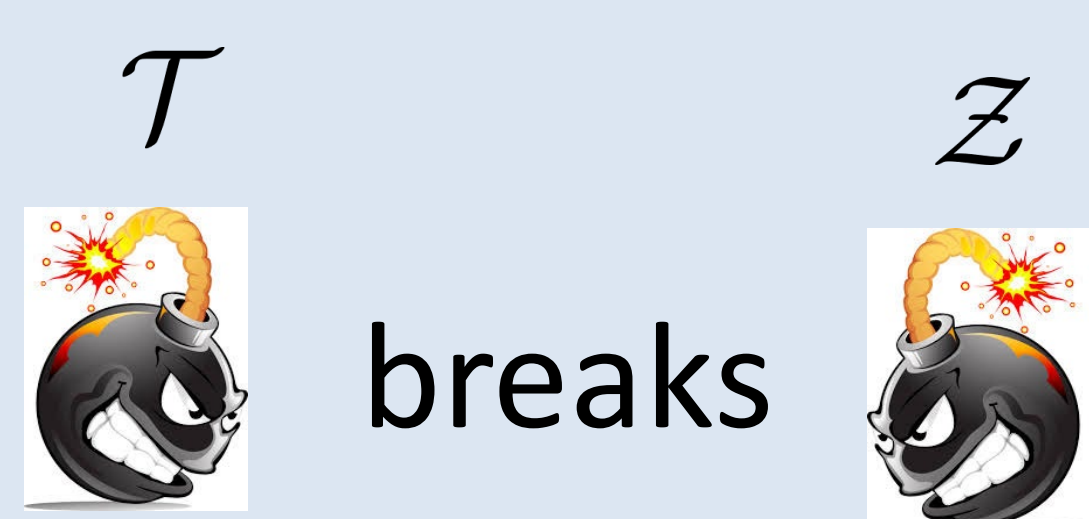


a **counterexample** to  $p$  is  
a subset-minimal  $\mathcal{T}$ , s.t.  $\mathcal{T} \models \forall_{t, t \neq p} (M \rightarrow t)$

an **explanation** of  $p$  is  
a subset-minimal  $\mathcal{Z}$  s.t.  $\mathcal{Z} \models (M \rightarrow p)$



## Duality between counterexamples and explanations



- every explanation  $\mathcal{Z}$  of  $p$  *breaks every* counterexample of  $p$ , and
- every counterexample  $\mathcal{T}$  of  $p$  *breaks every* explanation of  $p$ .

## Duality-based computation of all explanations ( or all counterexamples).

```

Input: formula  $M$  and prediction  $p$ 
Output: set  $\mathbb{E}$  of all absolute explanations of prediction  $p$ 
1  $(\mathbb{T}, \mathbb{E}, \mathcal{Z}) \leftarrow (\emptyset, \emptyset, \emptyset)$ 
2 do:
3   if  $\mathcal{Z} \models (M \rightarrow p)$  :
4      $\mathbb{E} \leftarrow \mathbb{E} \cup \{\mathcal{Z}\}$       #  $\mathcal{Z}$  is an explanation; save it
5   else:
6      $(\mathcal{T}, t) \leftarrow \text{ExtractInstance}()$   # get an instance  $\mathcal{T}$  with a
7     for  $l \in \mathcal{T}$  :
8       if  $(\mathcal{T} \setminus \{l\}) \models (M \rightarrow t)$  :
9          $\mathcal{T} \leftarrow \mathcal{T} \setminus \{l\}$ 
10     $\mathbb{T} \leftarrow \mathbb{T} \cup \{\mathcal{T}\}$       # update  $\mathbb{T}$  with a new counterexample  $\mathcal{T}$ 
11     $\mathcal{E} \leftarrow \text{MinimumHS}(\mathbb{T})$       # get a new hitting set of  $\mathbb{T}$ 
12  while  $\mathcal{Z} \neq \emptyset$ 
13 return  $\mathbb{E}$ 
Algorithm 1: Duality-based computation of all absolute explanations
    
```

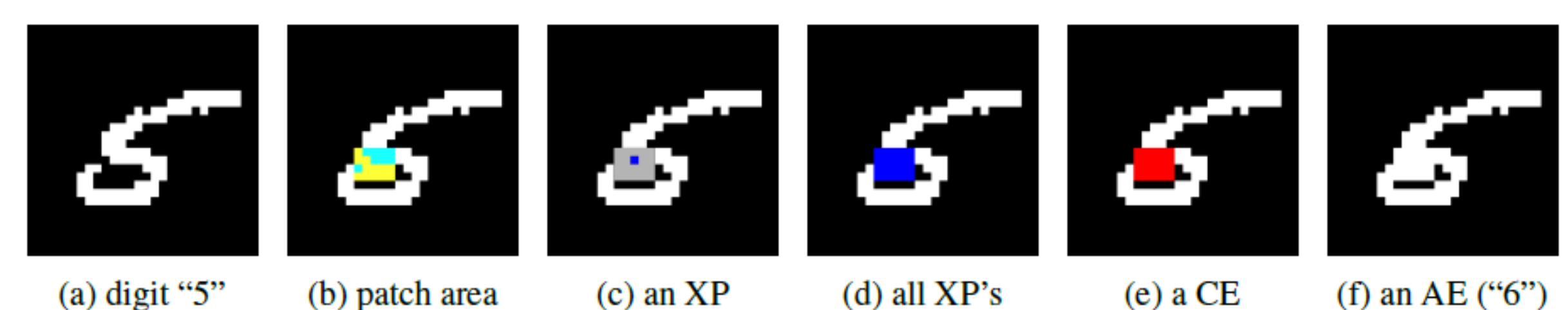


Figure 1: An example of digit five.

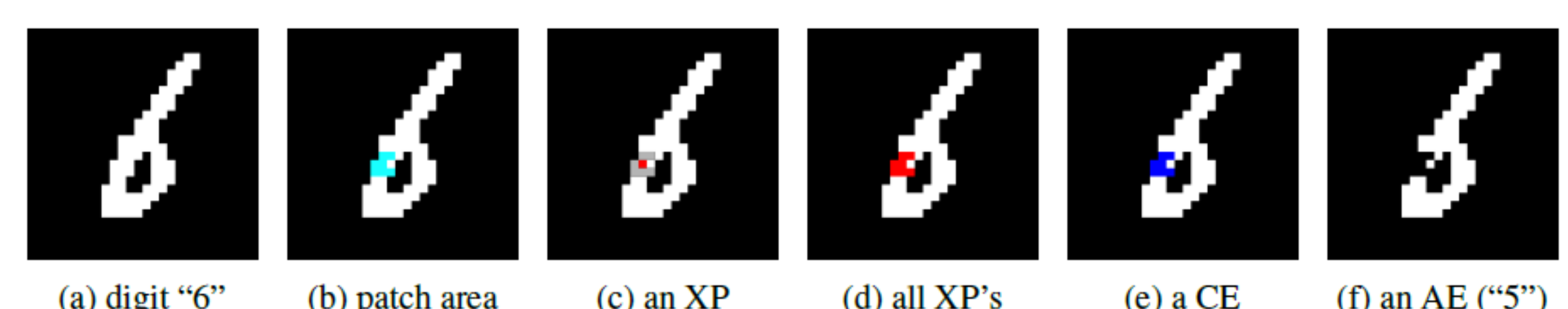


Figure 2: An example of digit six.

[Goodfellow et al.] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In ICLR, 2015.

[Ribeiro et al.] Marco Tulio Ribeiro, Sameer Singh, Carlos Guestrin. Why Should I Trust You? Explaining the Predictions of Any Classifier. In KDD, 2016.

[Shih et al.] Andy Shih and Arthur Choi and Adnan Darwiche. Abduction-Based Explanations for Machine Learning Models. In AAAI, 2019.

[Szegedy et al.] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In ICLR, 2014.

[Ignatiev et al.] Alexey Ignatiev, Nina Narodytska and Joao Marques-Silva. Abduction-Based Explanations for Machine Learning Models. In AAAI, 2019.

[Narodytska et al.] Nina Narodytska, Shiva Prasad Kasiviswanathan, Leonid Ryzhyk, Mooly Sagiv, Toby Walsh. Verifying Properties of Binarized Deep Neural Networks. In AAAI, 2018.