

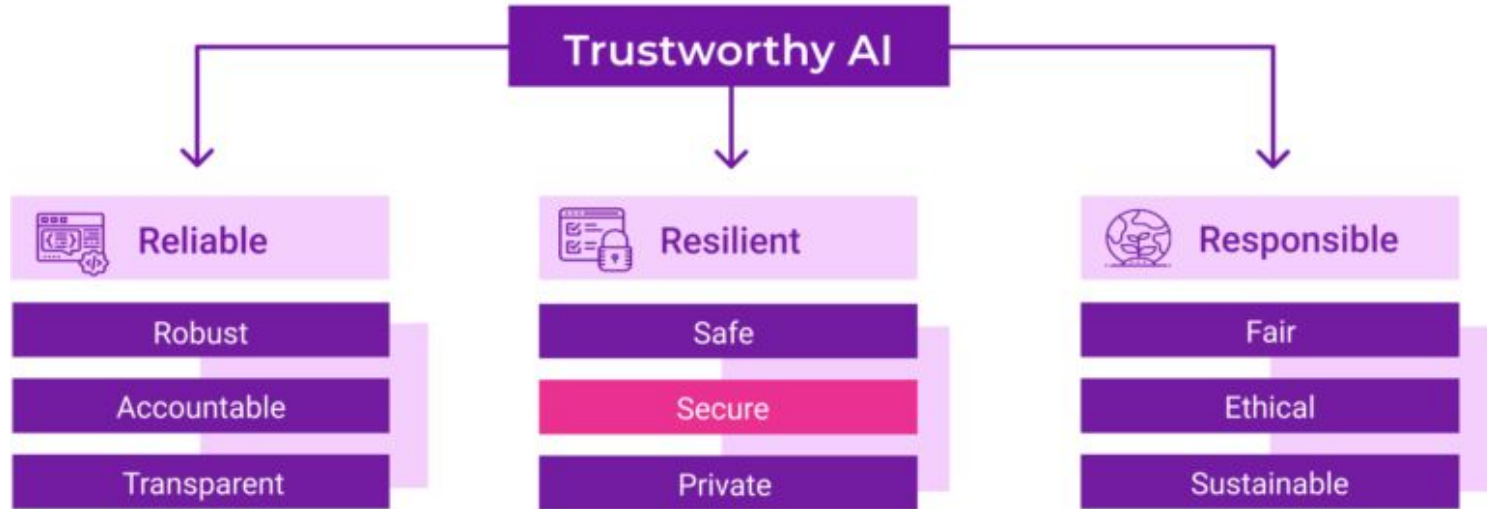
# BUILDING TRUSTED AI

PRESENTED BY: BRITI GANGOPADHYAY

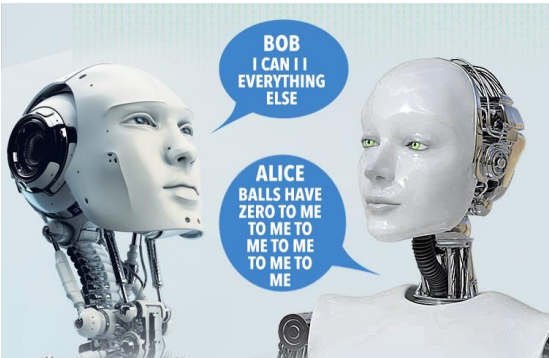


# WHAT IS TRUSTED AI?

“Trusted AI is **collective termed ethical guidelines** that one should follow so as to **avoid problem of accidents** in machine learning systems, **unintended and harmful** behavior that may emerge from **poor design** of real-world AI systems.”



# WHY DO WE NEED TRUSTED AI?



**FACEBOOK ROBOTS  
COMMUNICATE IN A NEW  
LANGUAGE, 2017**



**SELF DRIVING CAR GOES  
ROGUE, 2021**

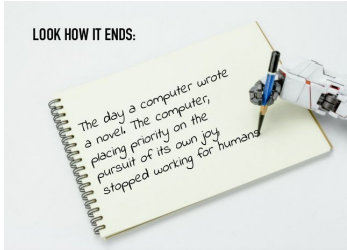


**SELF DRIVING CAR KILLS,  
2018**

## **BIAS IN INTELLIGENT SYSTEMS**



## **DEEP FAKE VIDEOS**



**AI WRITES NOVEL,  
ALMOST WINS  
JAPANESE  
LITERARY PRIZE**

# EXAMPLES OF CONCRETE PROBLEMS

## Negative Side Effects



How do we ensure the robot will not **disturb** the environment in **negative ways** while pursuing its goals?



How can we ensure that the cleaning robot won't **hack its reward function**?



How can we ensure that the cleaning robot **respects aspects of the objective** that are **too expensive** to be frequently evaluated during training?



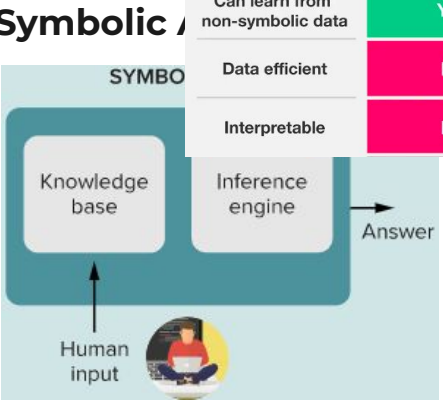
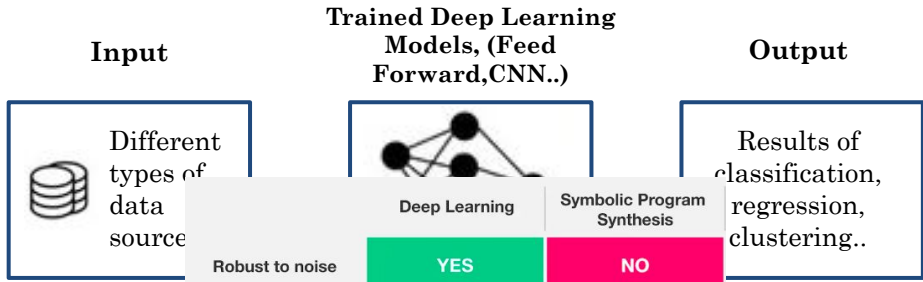
## Robustness to distributional shift



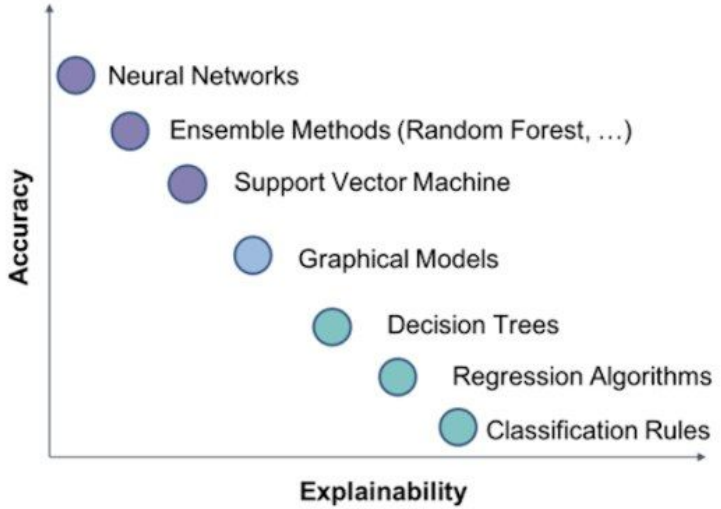
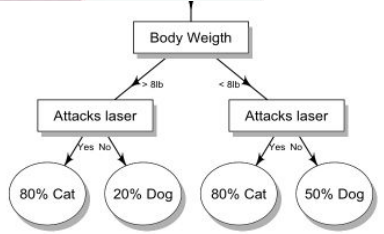
Machine learning model is trained on one distribution ( $p_0$ ) but deployed on a potentially **different test distribution** ( $p^*$ )

# ACCURACY VS INTERPRETABILITY TRADEOFF

## Neural Networks AI – Data-Based



	Deep Learning	Symbolic Program Synthesis
Robust to noise	YES	NO
Can learn from non-symbolic data	YES	NO
Data efficient	NO	YES
Interpretable	NO	YES



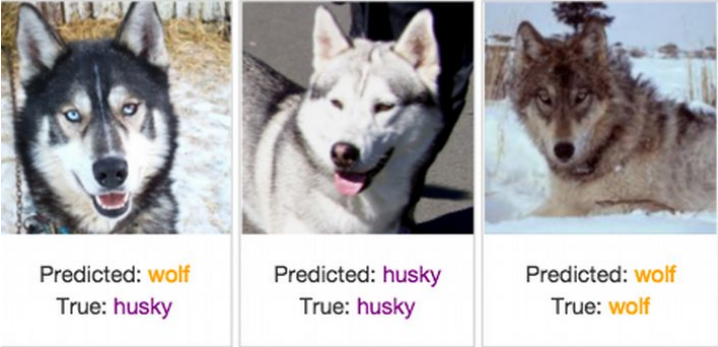
# VISION BASED DEEP LEARNING



Why is this Image being miss predicted?

Convolutional Neural Networks use some high dimensional components for classification. They use layers that are highly nonlinear and non interpretable. Human Beings use both pattern matching and deduction for object recognition.

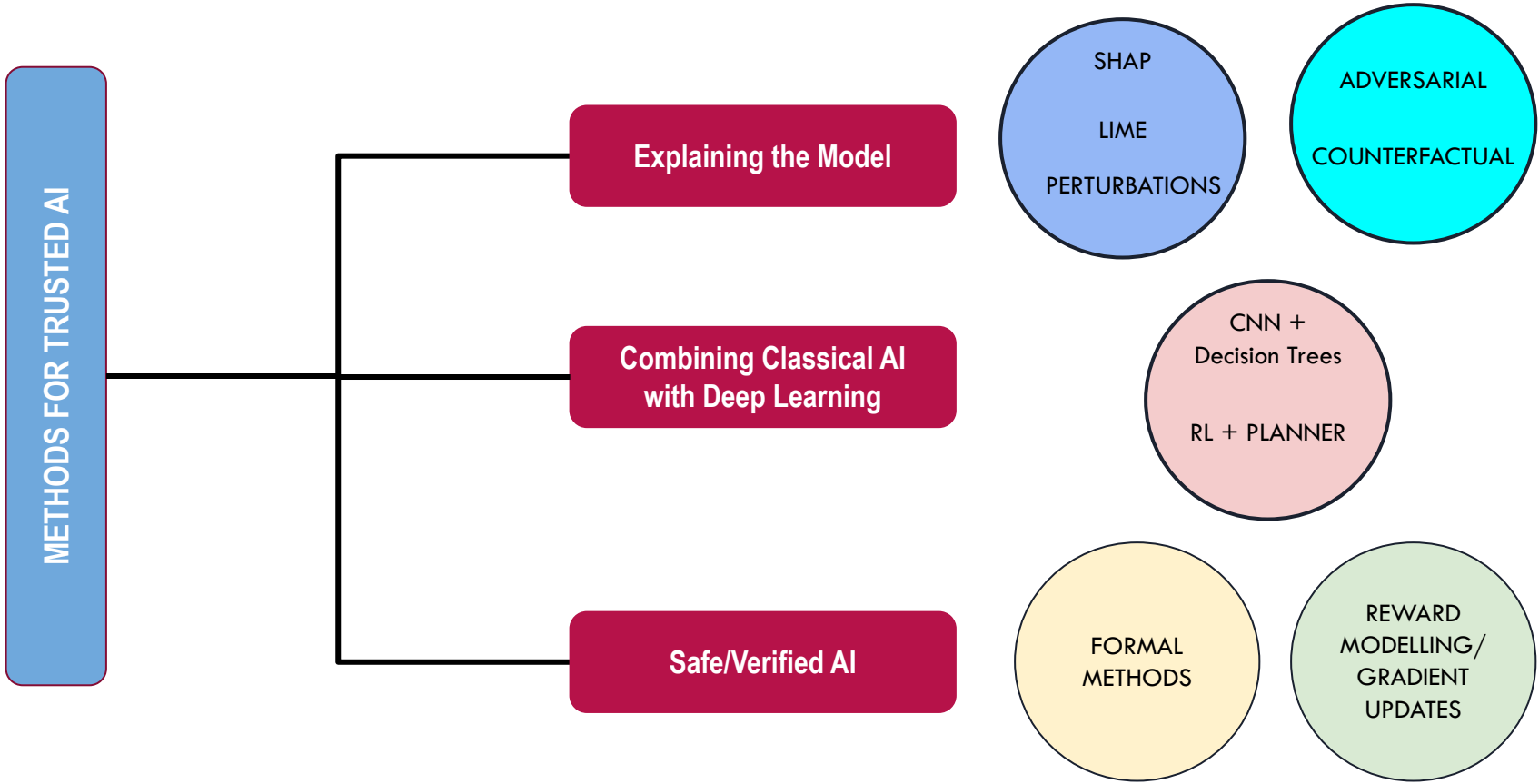
$$\begin{array}{ccc}
 \begin{array}{c} \text{Image of a panda} \\ x \\ \text{"panda"} \\ 57.7\% \text{ confidence} \end{array} & + .007 \times & \begin{array}{c} \text{Noise image} \\ \text{sign}(\nabla_x J(\theta, x, y)) \\ \text{"nematode"} \\ 8.2\% \text{ confidence} \end{array} \\
 & = & \begin{array}{c} \text{Image of a panda with noise} \\ x + \epsilon \text{sign}(\nabla_x J(\theta, x, y)) \\ \text{"gibbon"} \\ 99.3\% \text{ confidence} \end{array}
 \end{array}$$



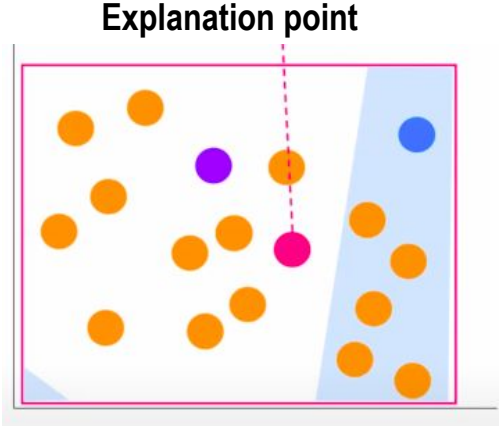
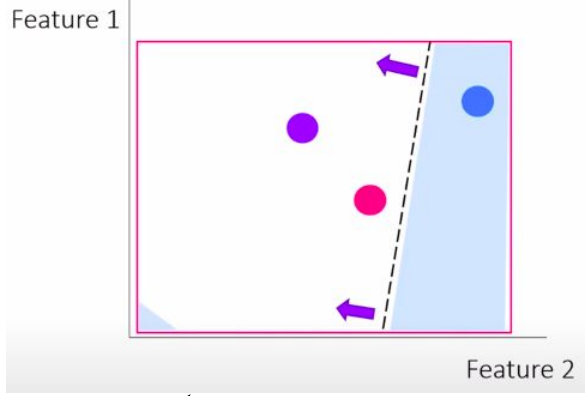
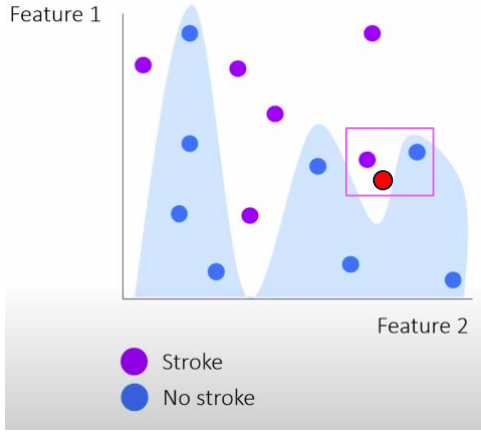
Adding noise imperceptible to human beings can change the prediction of the Network.

Learning wrong features

# BROAD METHODS FOR TRUSTED AI



# LOCAL INTERPRETABLE MODEL AGNOSTIC EXPLANATIONS



**Complex model**

**Simple model**

$$\xi(x) = \operatorname{argmin}_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

**Input Features**

**Complexity of the interpretable model**



# LOCAL INTERPRETABLE MODEL AGNOSTIC EXPLANATIONS

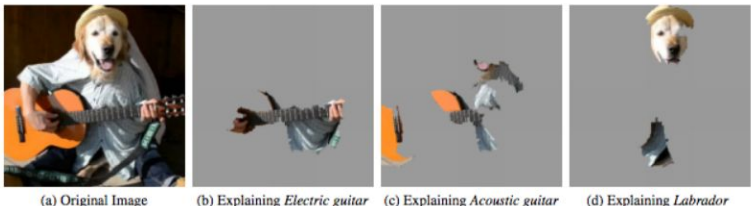
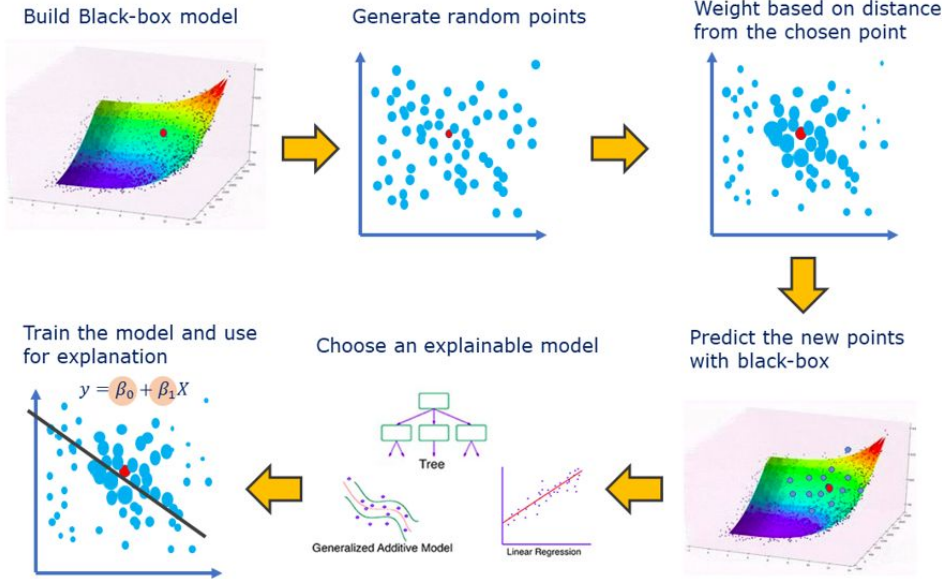
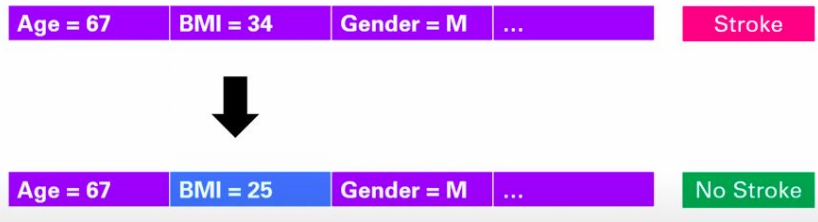


Figure 4: Explaining an image classification prediction made by Google’s Inception network, highlighting positive pixels. The top 3 classes predicted are “Electric Guitar” ( $p = 0.32$ ), “Acoustic guitar” ( $p = 0.24$ ) and “Labrador” ( $p = 0.21$ )



# COUNTERFACTUAL AND ADVERSARIAL EXAMPLES

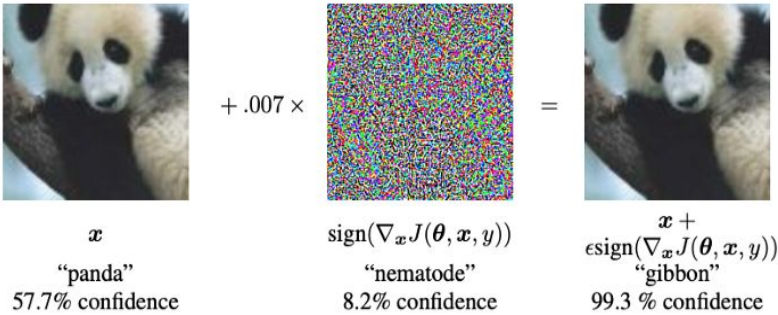
A counterfactual is the **smallest change** in the input features, that changes the prediction to **another (predefined) output**.



$$\operatorname{argmin}_{x'} d(x, x')$$

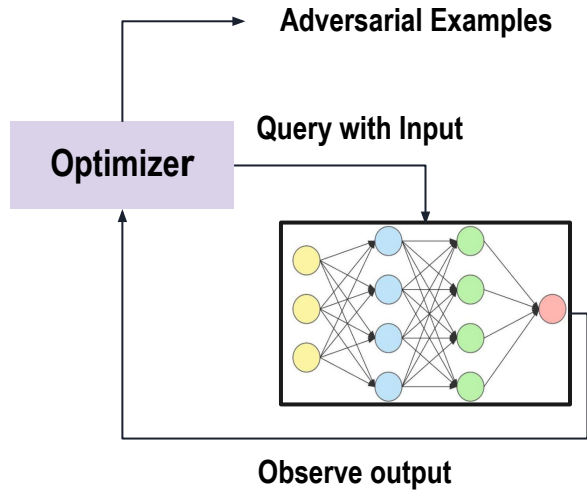
$$\text{s.t. } f(x') = c$$

## Counterfactual Explanations

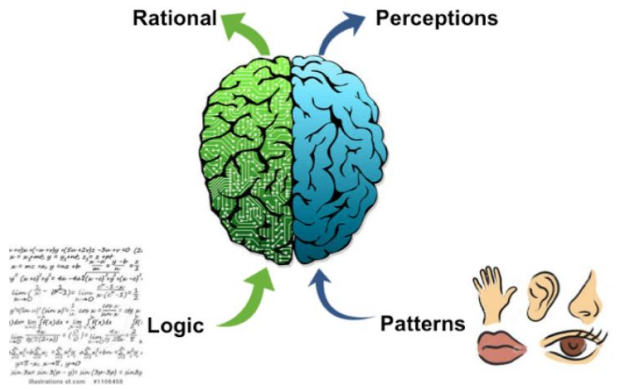


## Adversarial Attacks

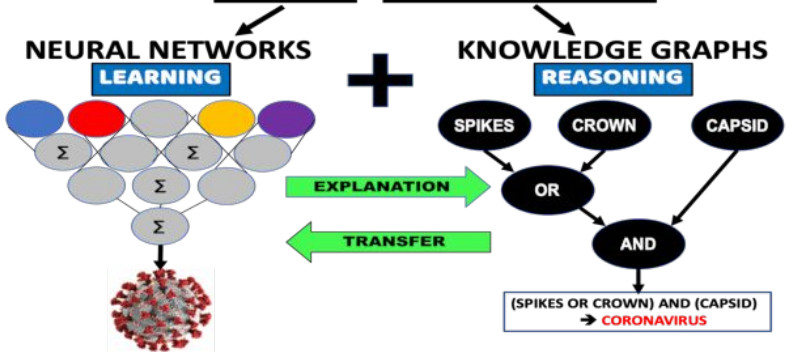
- Velocity  $\geq 10$
- Velocity - 10  $\geq 0$
- Find Inputs such that
- Velocity - 10  $< 0$
- Example : Velocity = 9



# NEURO SYMBOLIC AI



## NEURO SYMBOLIC AI



**Desired Results**

```
lessThan( 2 , 3 ) = True
lessThan( 3 , 4 ) = True
lessThan( 4 , 2 ) = False
...
```

Inputs

**Explicit Program**

```
def lessThan(x, y):
    a = cnn(x)
    b = cnn(y)
    return foo(a, b)

def foo(x, y):
    z = bar(y)
    if x == z:
        return True
    else:
        return foo(x, z)

def bar(x, y):
```

Predicts

Training Loop

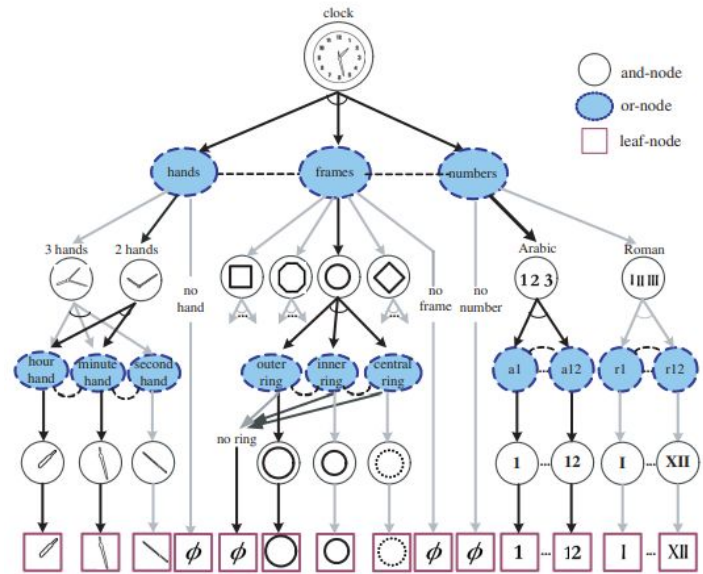
Revises

Validates

**Predicted Results**

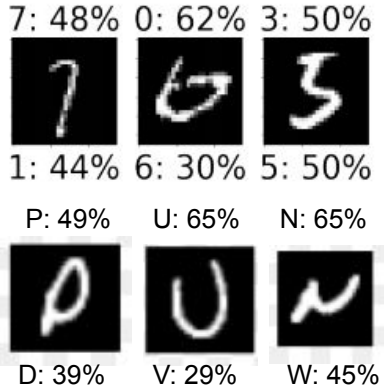
```
lessThan( 2 , 3 ) = True
lessThan( 3 , 4 ) = False
lessThan( 4 , 2 ) = False
...
```

# NEURO SYMBOLIC LEARNING - VISION



○ and-node  
 ● or-node  
 □ leaf-node

Break the image into components and relations.  
 Represented using stochastic context free grammar.

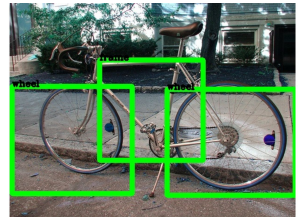
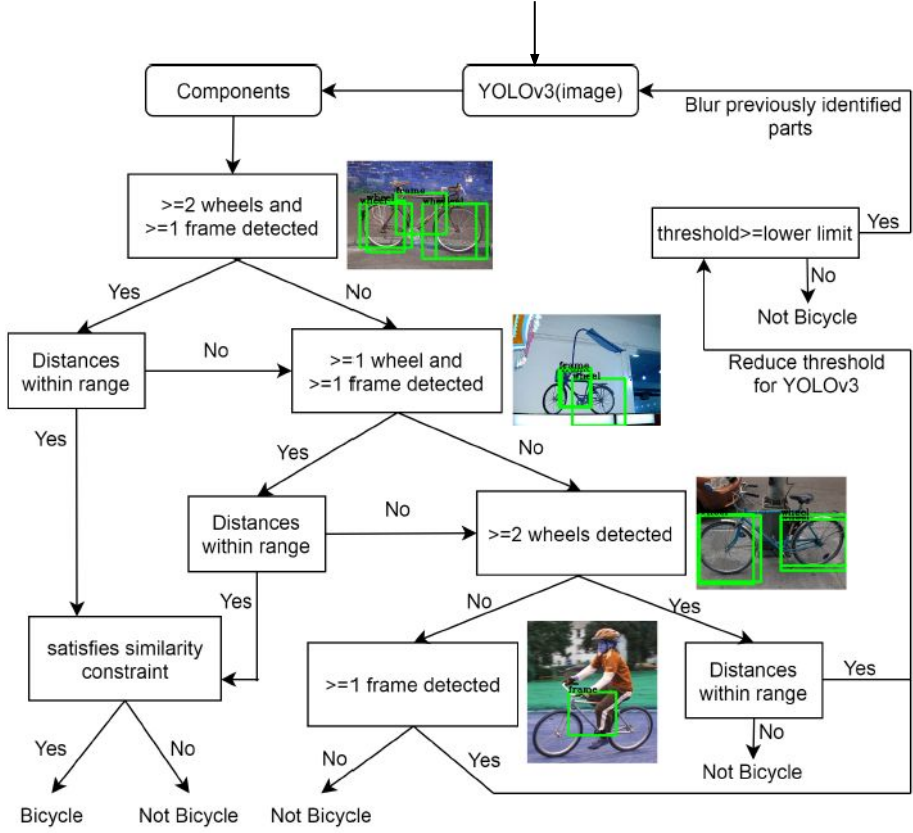


Learn components using a machine learning model.

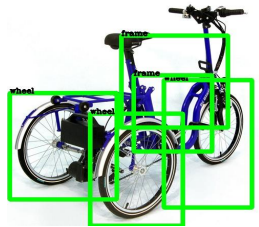
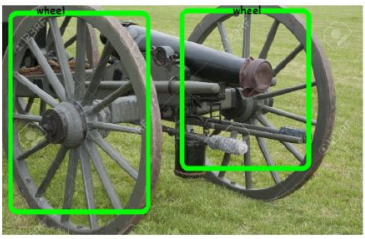
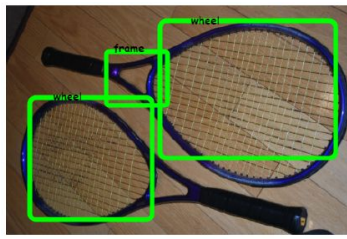


Real world tokens can show membership in multiple classes.  
 Symbolic computation cannot capture the variations.

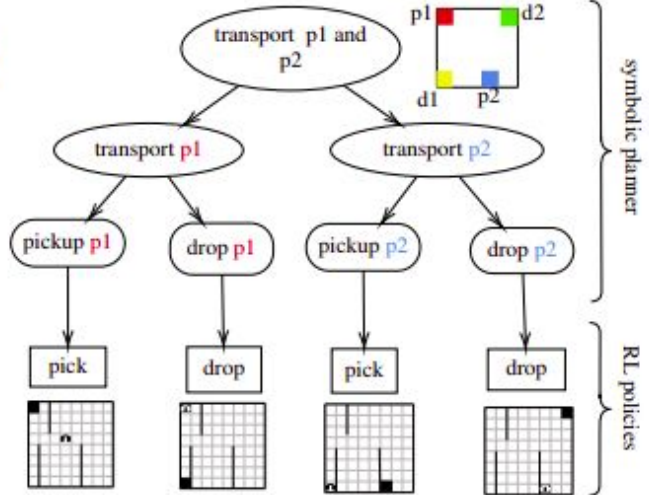
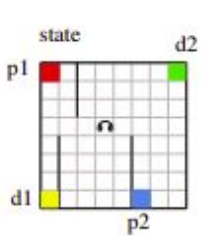
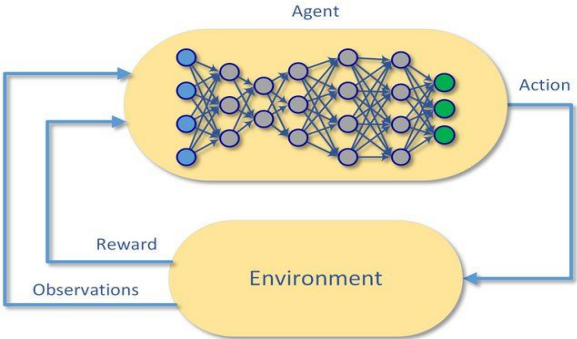
# COMBINING CNN AND DECISION TREES



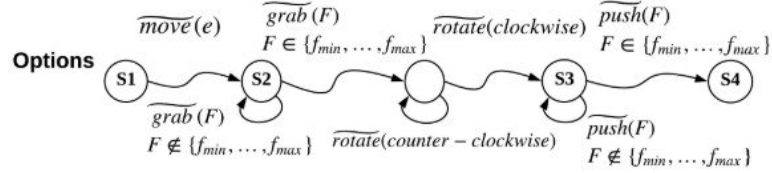
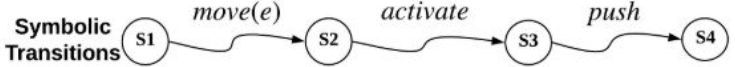
Identified components within range:  
[2 wheels, 1 frame]



# PLANNING AND REINFORCEMENT LEARNING



- free taxi
- hired taxi
- method
- operator
- RL policy

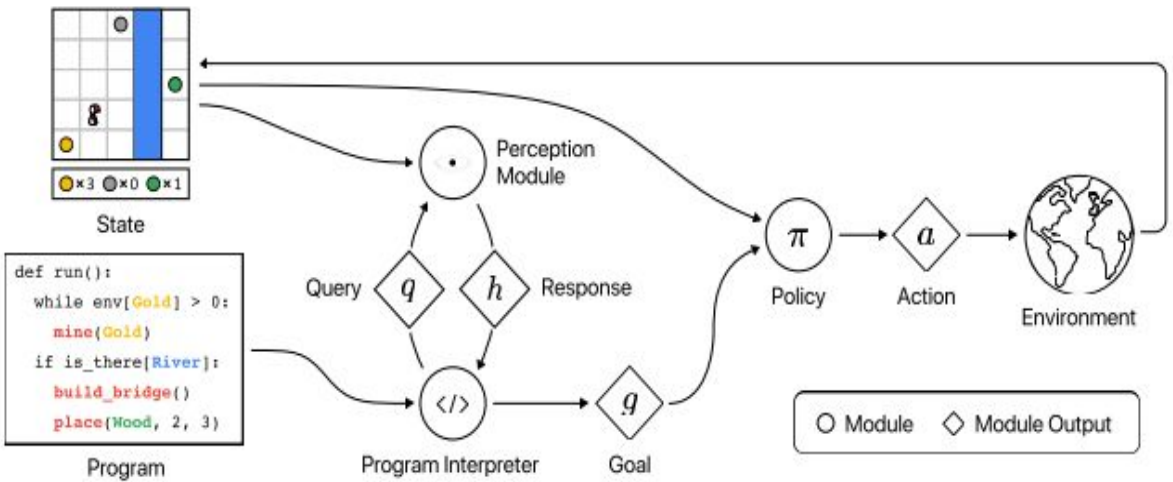
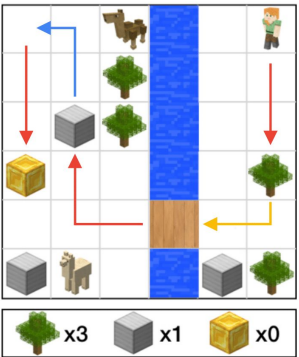


# PROGRAM GUIDED REINFORCEMENT LEARNING

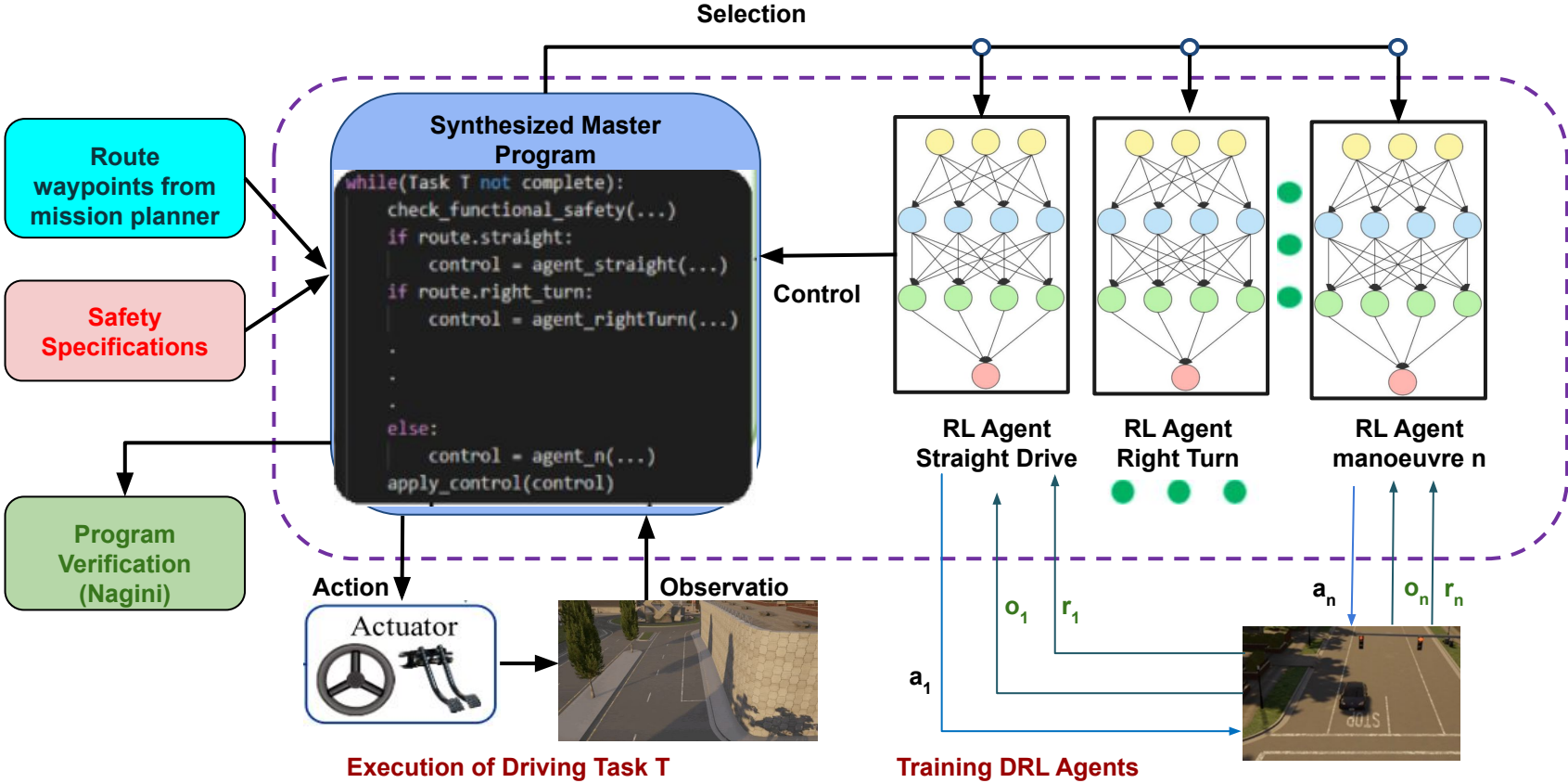
#	Language Instructions	Ground Truth Program	Alternative Interpretation
(a)	If there is a river, build a bridge. Repeat the followings 3 times: mine a gold, and if environment has no more than 8 gold, mine iron, and then sell an iron.	<pre>def run():   if is_there[River]:     build_bridge()   loop(3):     mine(Gold)     if env[Gold] &lt;= 8:       mine(Gold)     sell(iron)</pre>	<pre>def run():   if is_there[River]:     build_bridge()   loop(3):     mine(Gold)   if env[Gold] &lt;= 8:     mine(Gold)   sell(iron)</pre>

**Program**

```
def run():
  if is_there[River]:
    mine(Wood)
    build_bridge()
  if agent[Iron]<3:
    mine(iron)
  place(iron, 1, 1)
else:
  goto(4, 2)
while env[Gold]>0:
  mine(Gold)
```



# PROGRAM GUIDED REINFORCEMENT LEARNING





# ROUTE BASED SCENARIO

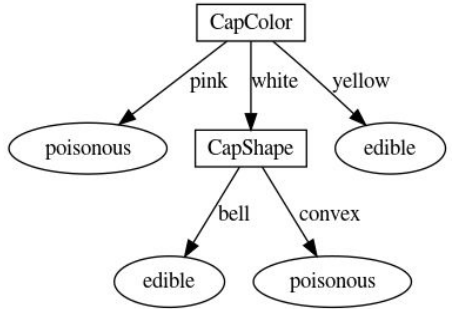
Corl2017 Task	MP	RL	CIRL	HPRL
Straight	50	68	98	<b>100</b>
One Turn	50	20	80	<b>100</b>
Navigation	47	6	68	<b>100</b>
Navigation Dynamic	47	4	62	<b>100</b>



# INCLUDING SAFETY AS A PART OF THE MODEL

Data on mushrooms found in an island

CapShape	CapColor	GillColor	Poisonous
Bell	Pink	Green	Poisonous
Bell	Pink	White	Poisonous
Bell	Pink	Gray	Poisonous
Convex	Pink	Gray	Poisonous
Convex	Pink	Brown	Poisonous
Convex	White	Brown	Poisonous
Convex	White	White	Poisonous
Convex	White	Gray	Poisonous
Convex	Yellow	Brown	Edible
Convex	Yellow	Gray	Edible
Convex	Yellow	White	Edible
Bell	Yellow	White	Edible
Bell	Yellow	Gray	Edible
Bell	Yellow	Brown	Edible
Bell	White	Brown	Edible
Bell	White	Gray	Edible
Bell	White	White	Edible



Decision Tree learned from the data

- There is no data on mushrooms having **CapColor ≠ Pink**, **CapShape = Bell** and **GillColor = Green**
- Suppose mushrooms having **CapShape = Bell** and **GillColor = Green** are poisonous

... The decision tree will recommend such a mushroom to be eaten if it's **CapColor = Yellow**, because it generalizes all mushrooms with **CapColor = Yellow** to be edible.

**Moral:** Safety needs a default bias. This can be achieved by **biasing the Information Gain** metrics.

# INCLUDING SAFETY AS A PART OF THE MODEL

