

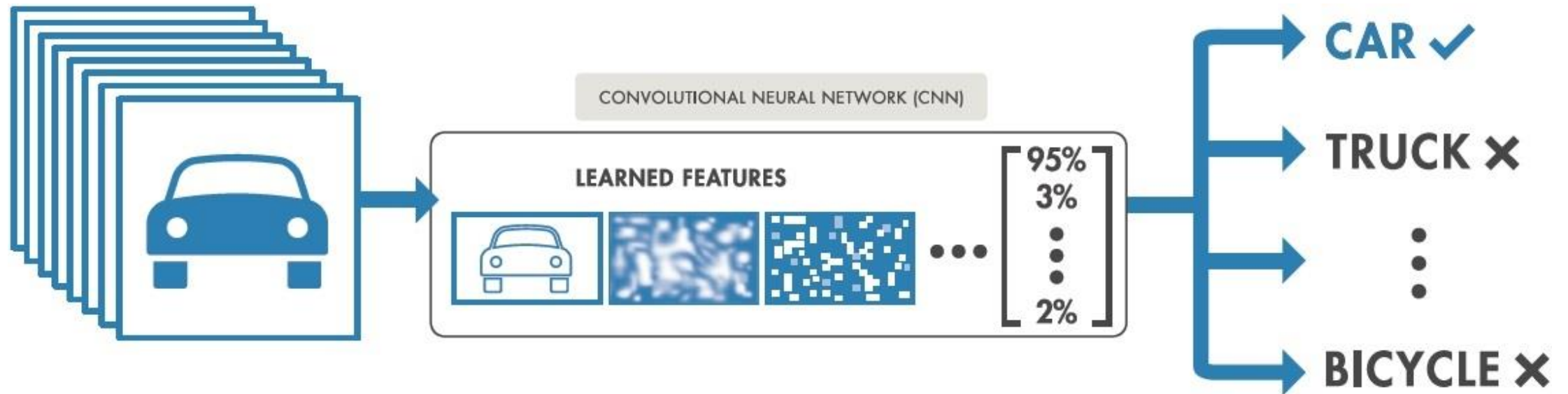
Convolutional Neural Networks

Increasing trust through visualization

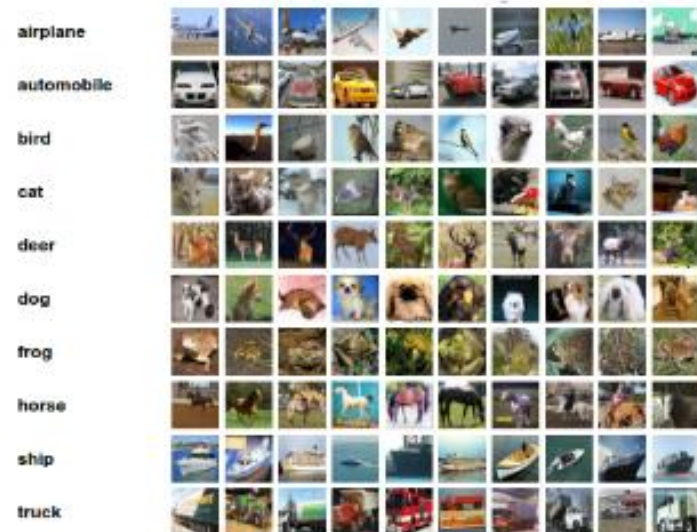
Sam Green, samgreen@gmail.com

<https://twitter.com/teenybiscuit>

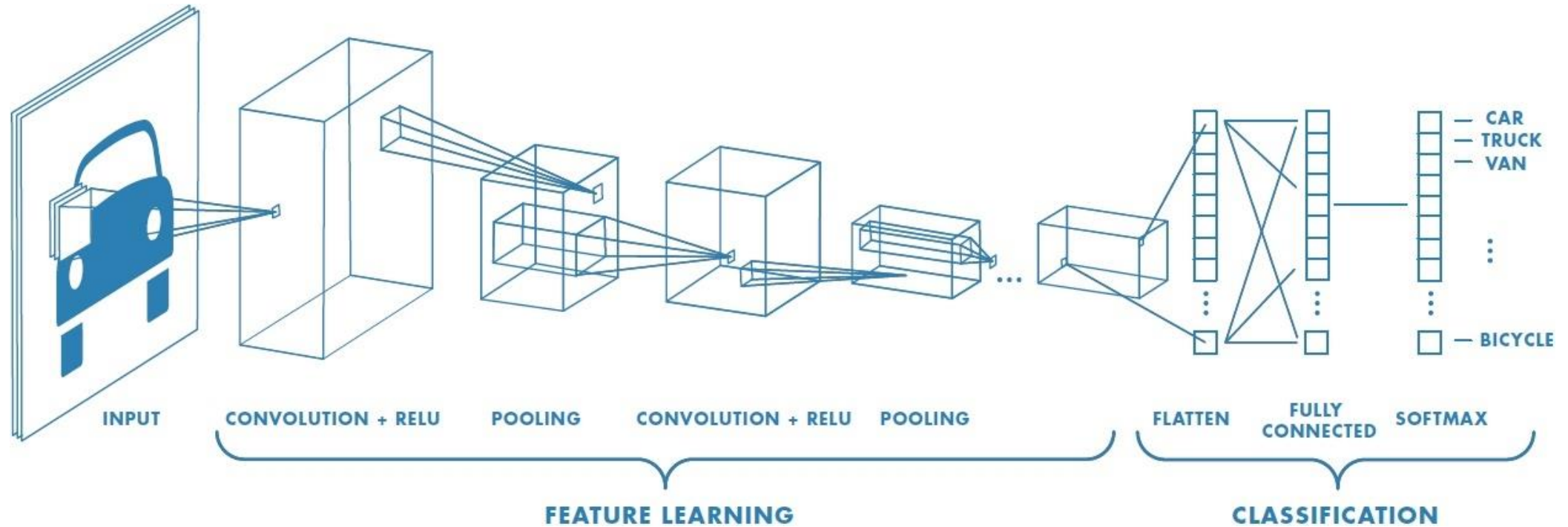
Convolutional Neural Network (CNN)



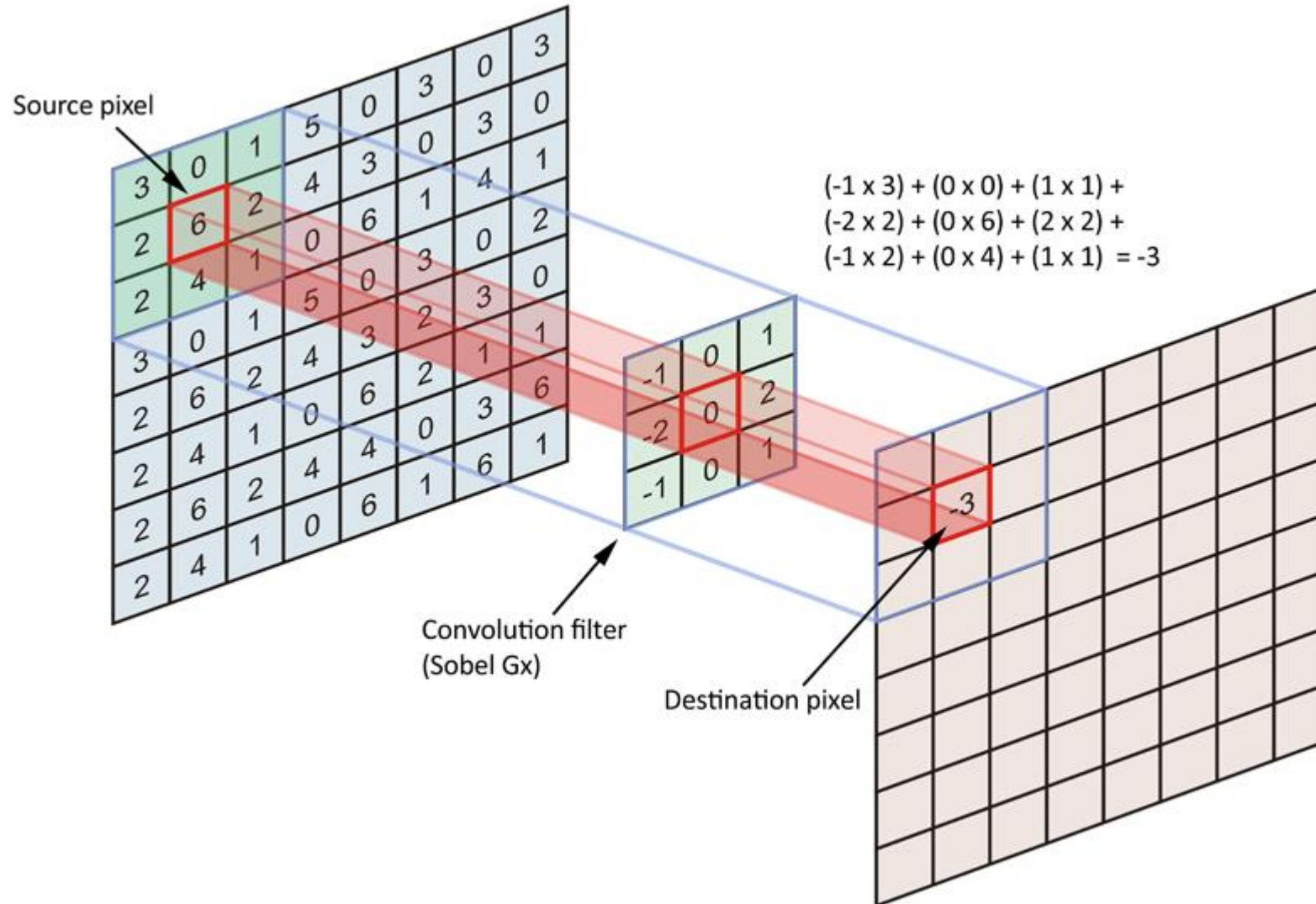
Labeled training images



Convolutional Neural Network (CNN)



Convolution



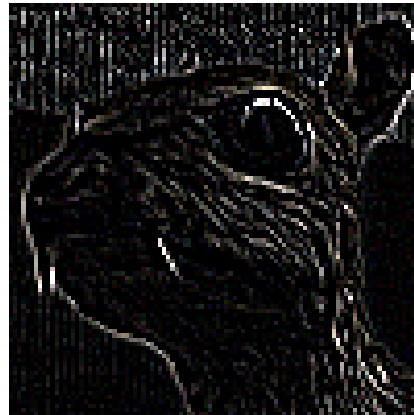
Convolution



$$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$

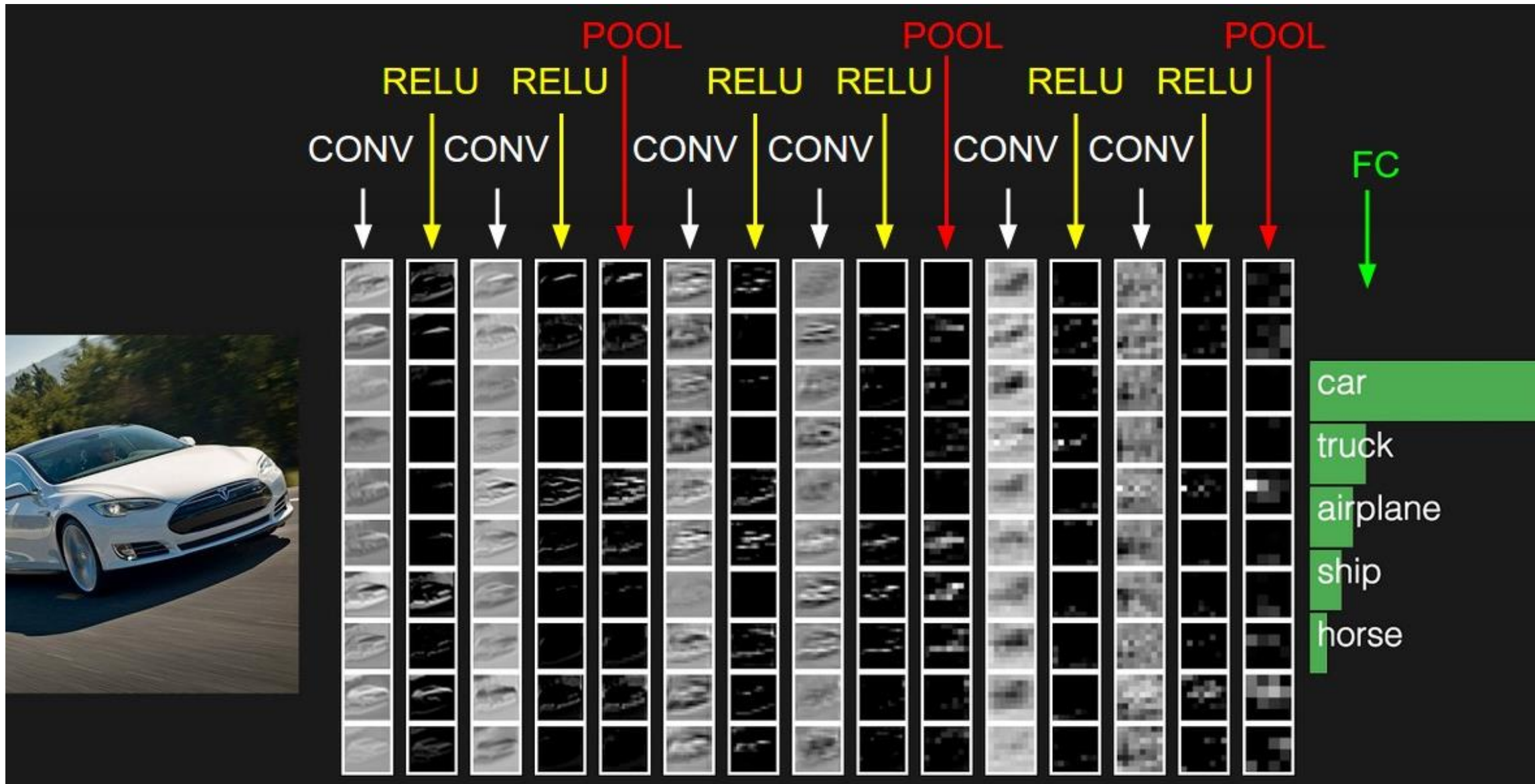


$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$



$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$





Power of CNNs

Beating Go

(and chess, shogi, checkers, backgammon, Dota 2,...)



Face recognition

Image caption generation

Prosthetics control

Colorizing black and white images

ATM deposits

Speech recognition

Breed recognition



Weirdness



A close up of a lush green field

Tags: grass, field, sheep, standing, rainbow, man

Weirdness



A herd of sheep grazing on a lush green hillside
Tags: grazing, sheep, mountain, cattle, horse

Weirdness



Left: A man is holding a dog in his hand
Right: A woman is holding a dog in her hand
Image: @SouperSarah

Weirdness

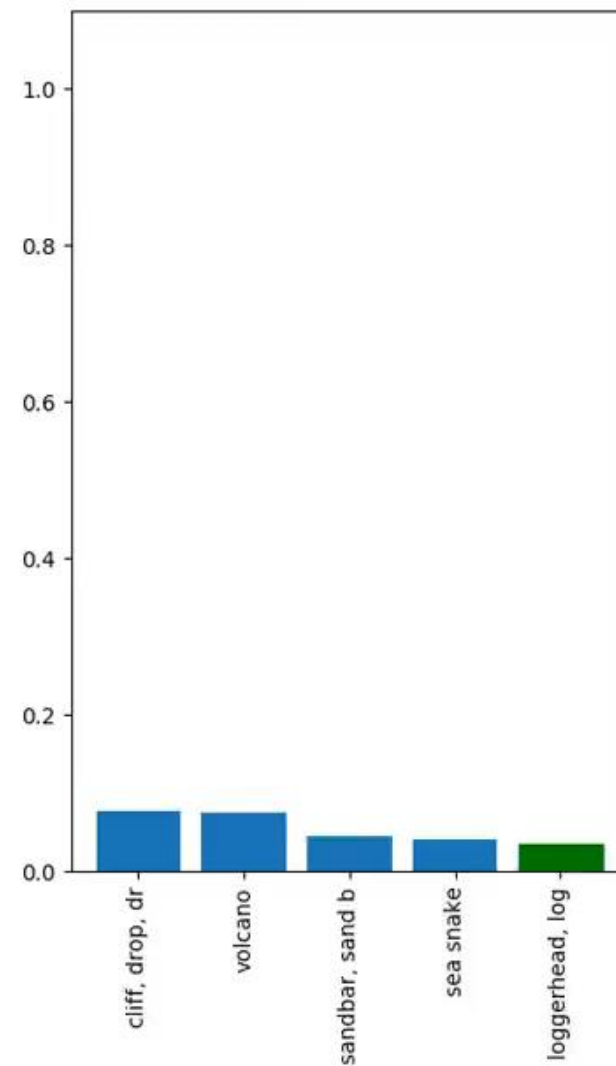


NeuralTalk2: A flock of birds flying in the air

Microsoft Azure: A group of giraffe standing next to a tree

Image: Fred Dunn, <https://www.flickr.com/photos/gratapictures> - CC-BY-NC

Attacks



t-SNE

t-Distributed Stochastic Neighbor Embedding

t-SNE results with AlexNet

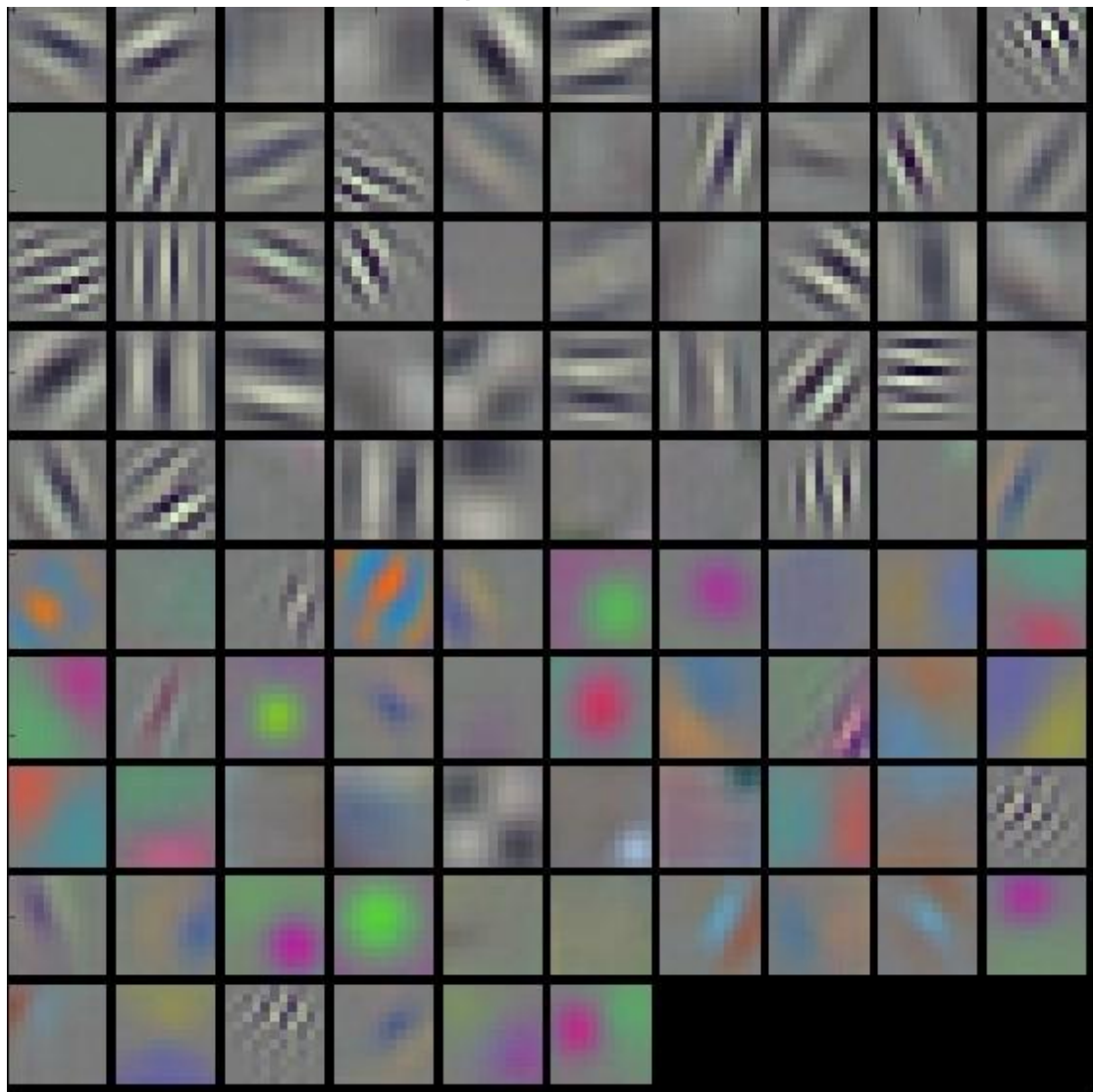
Maps 4096-dimensional CNN output
to 2-dimensions



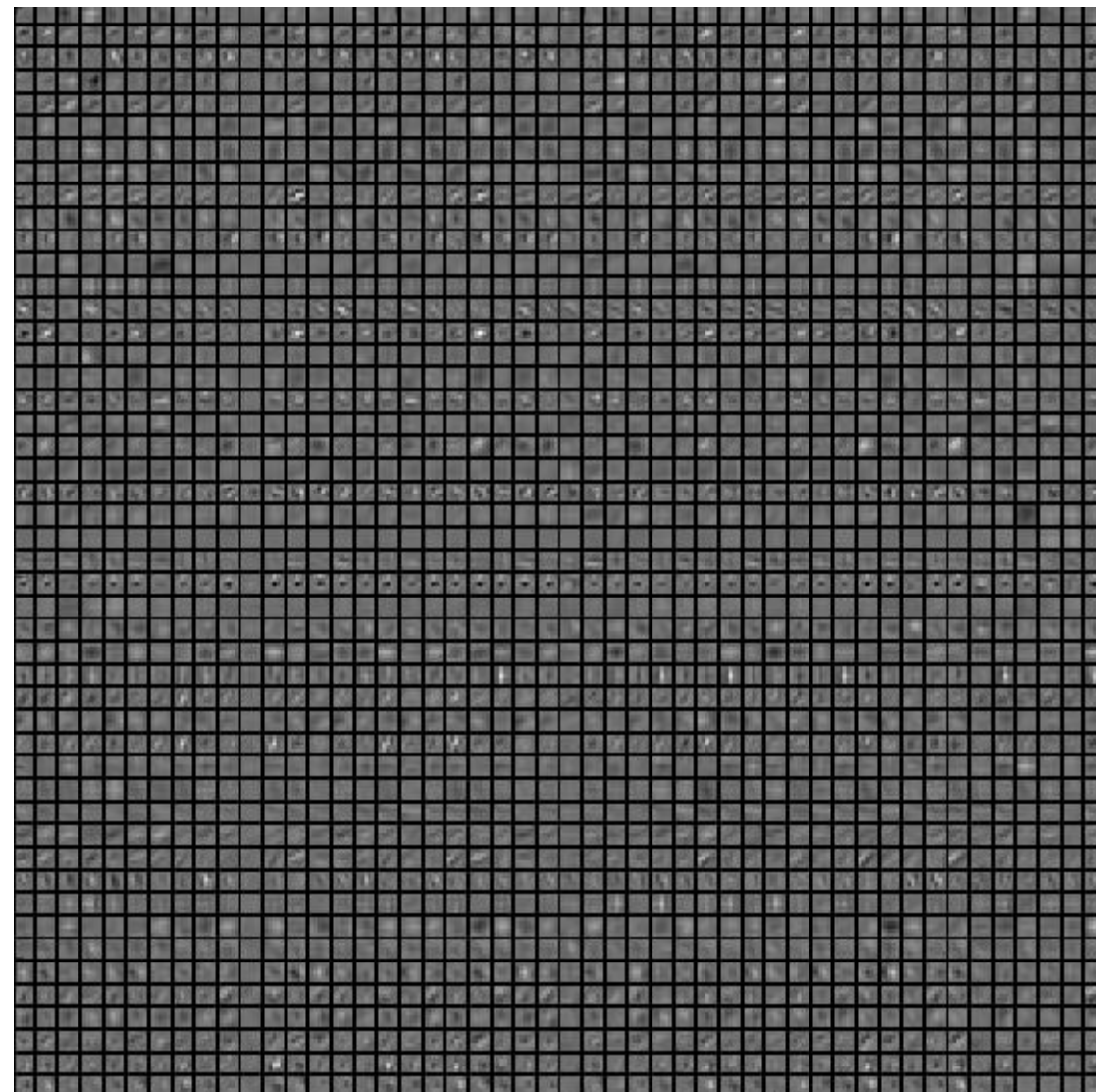
Filter Visualization

Inspecting what the CNN learned

AlexNet | conv1



AlexNet | conv2



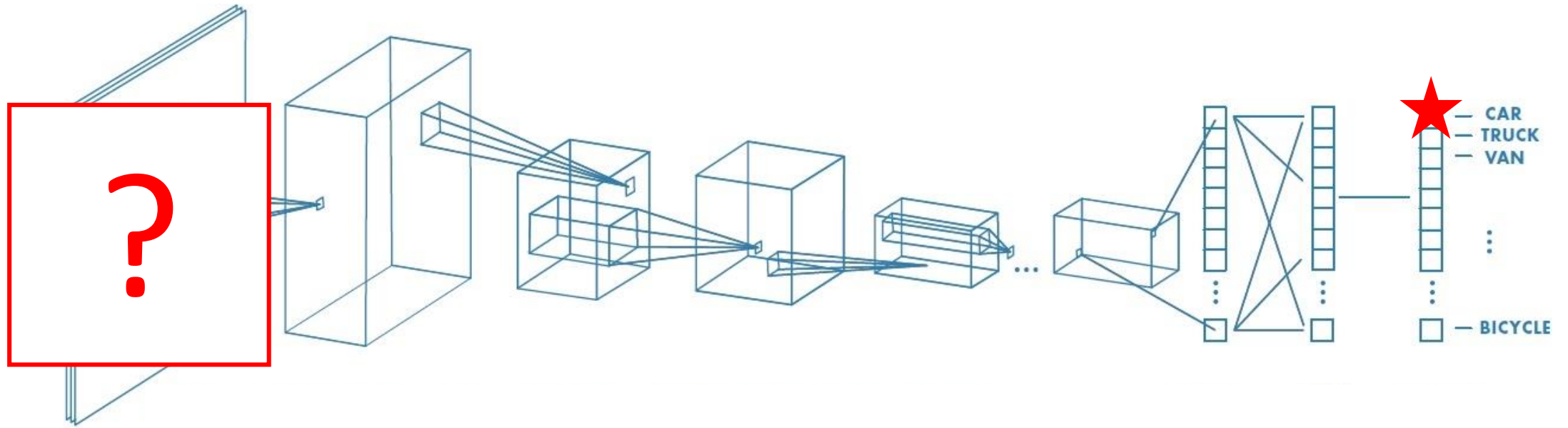


MNIST Filters

Class Visualization

Generating inputs to activate classes

What will trigger an output?



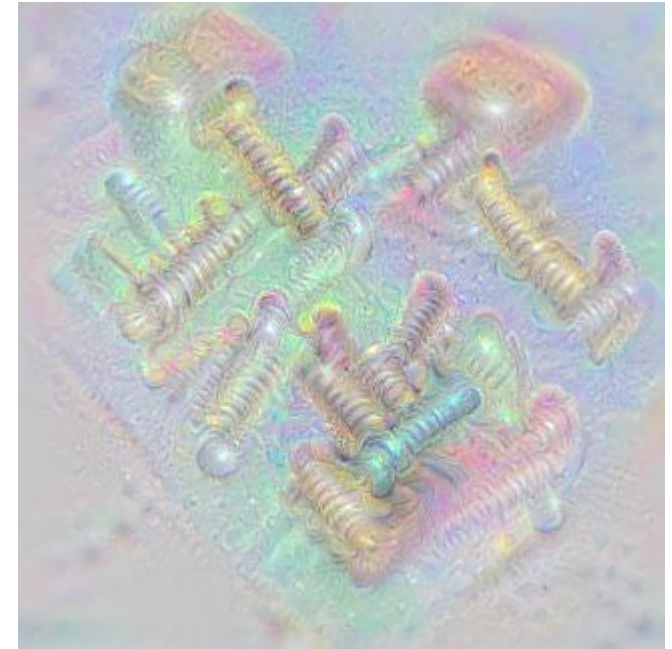
Visualizing GoogLeNet Classes



Loggerhead turtle



Saxophone



Screws

Visualizing GoogLeNet Classes



Shetland sheepdog

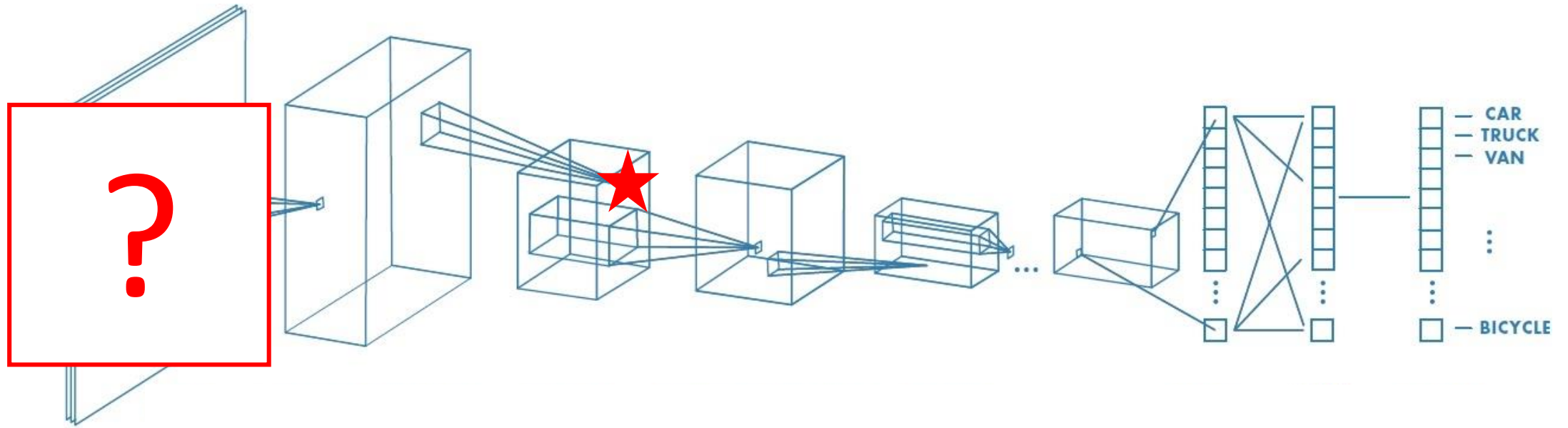


Basset hound

Feature Visualization

Generating inputs to activate neurons

What image triggers a neuron?



Images that activate neurons in GoogLeNet

Mixed4a, neuron 6

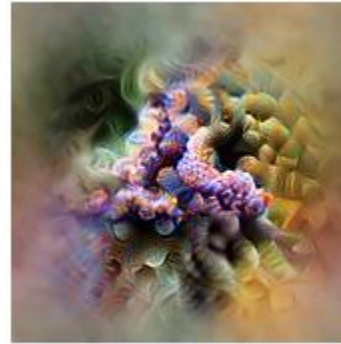
Mixed4a, neuron 240

Mixed4a, neuron 492

Triggering images
in training set



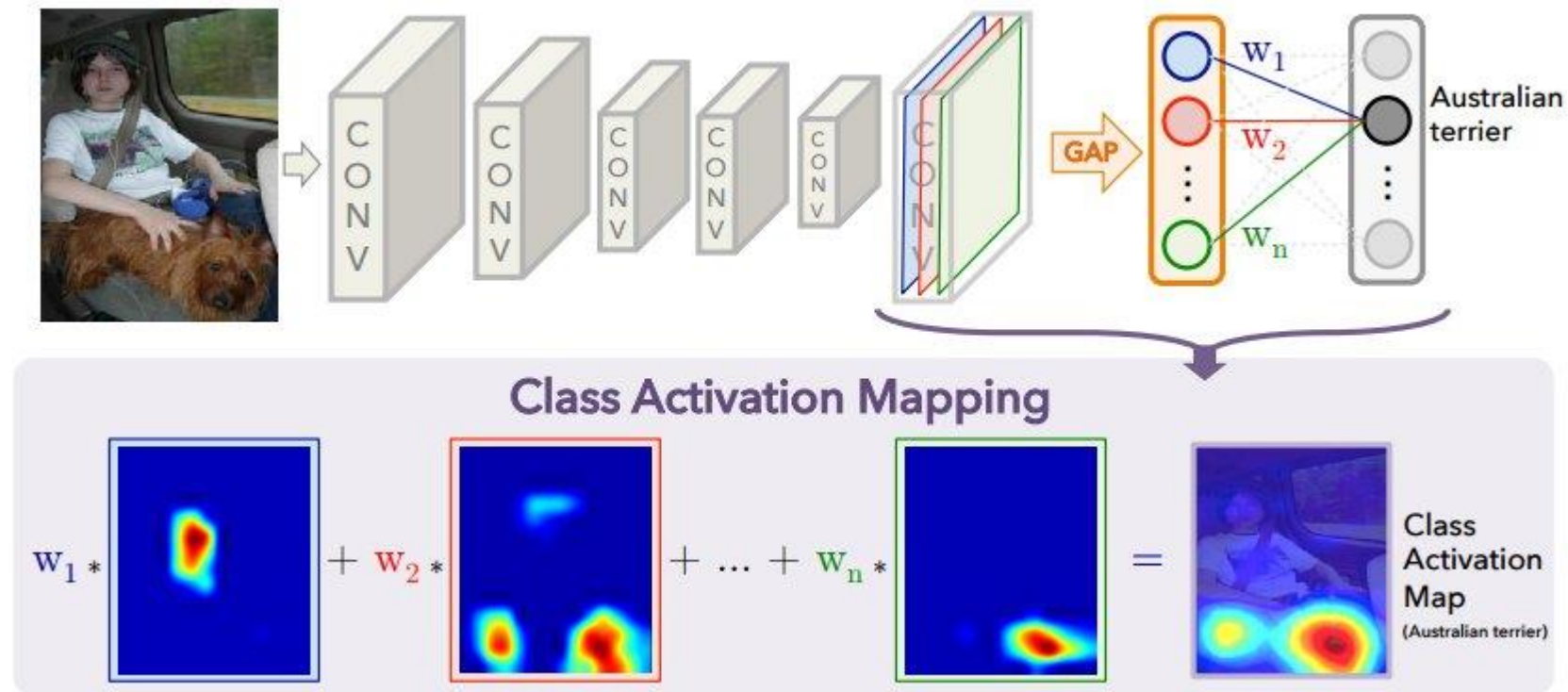
Feature visualization



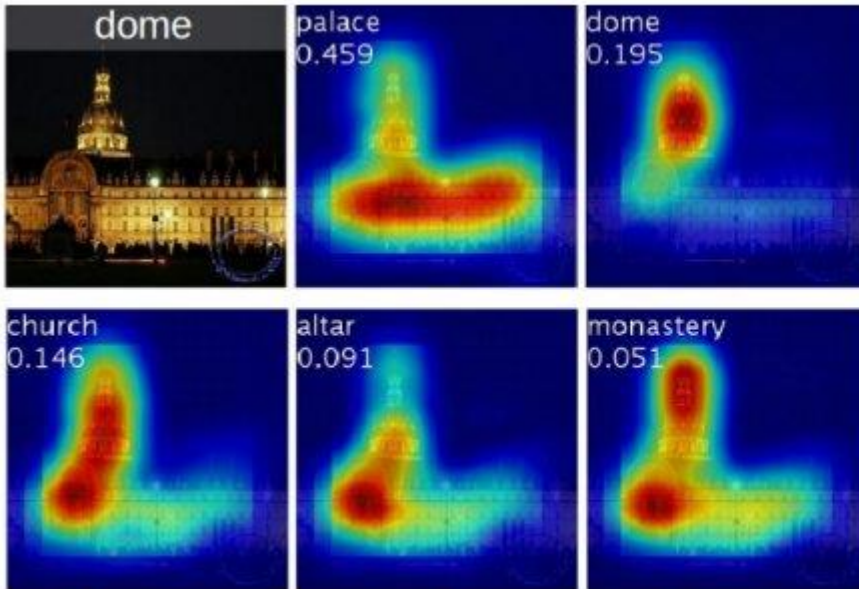
Attribution Visualization

Understanding what triggers a class selection

Class Activation Mapping



Class Activation Mapping



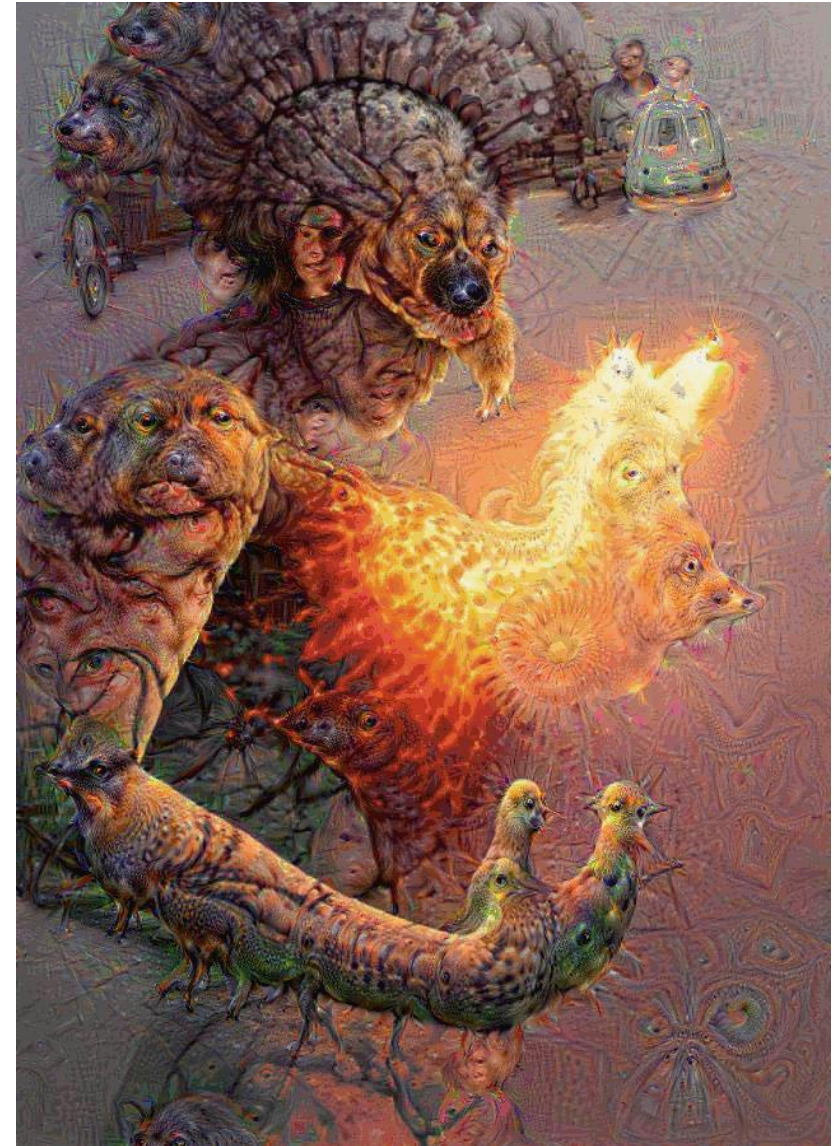
Class activation maps of top 5 predictions



Class activation maps for one object class

DeepDream

Modifying image to generate class activations



DeepDream



Horizon



Towers & Pagodas



Trees



Buildings



Leaves

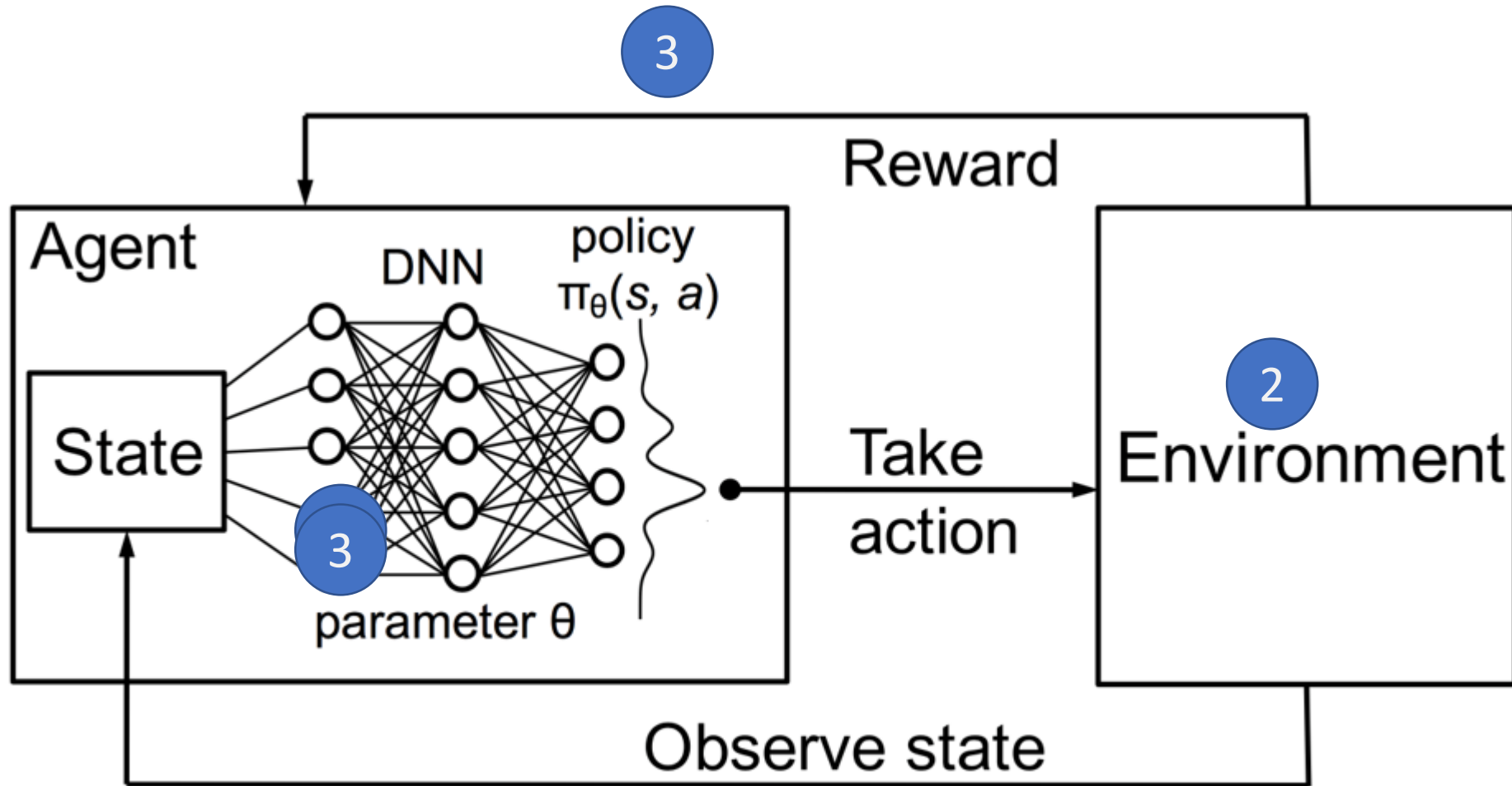


Birds & Insects

Applications to Reinforcement Learning

Visualizing the policy

Reinforcement Learning Paradigm



Reinforcement Learning Cube Example

State observation is camera

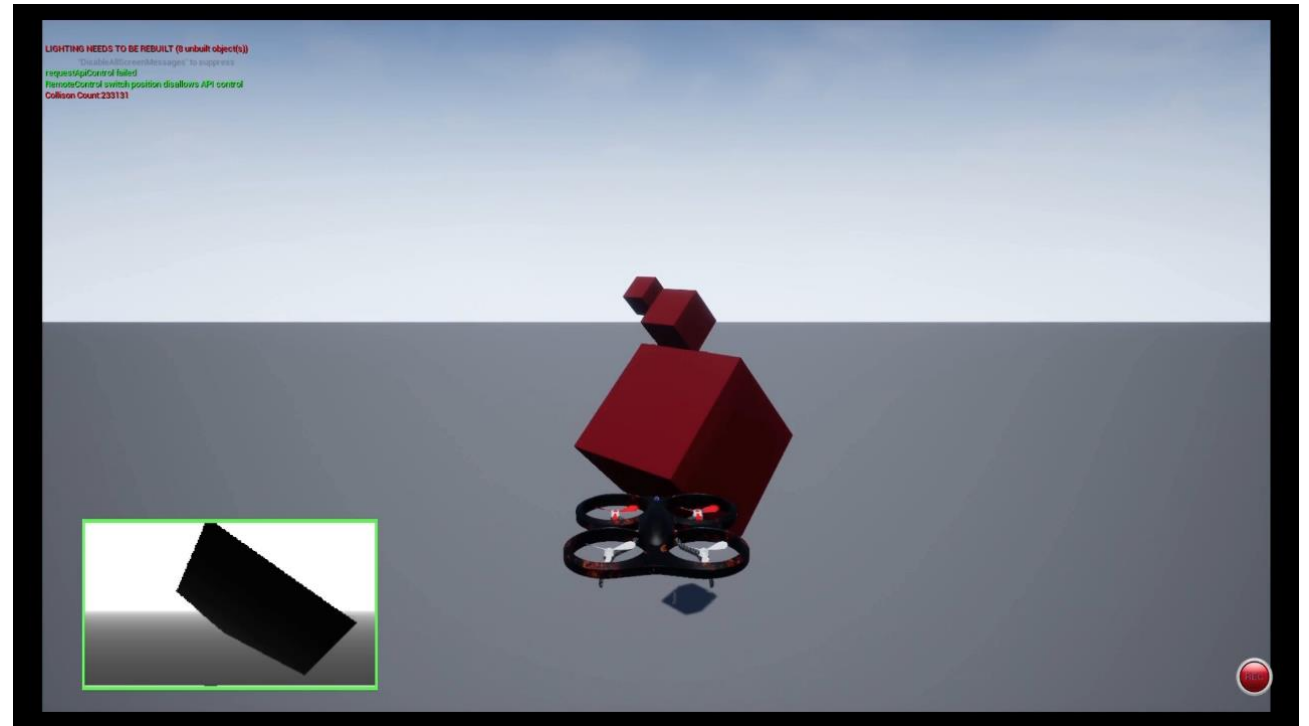
CNN for policy π_θ

Actions: left, forward, right

Reward +1 for hitting box

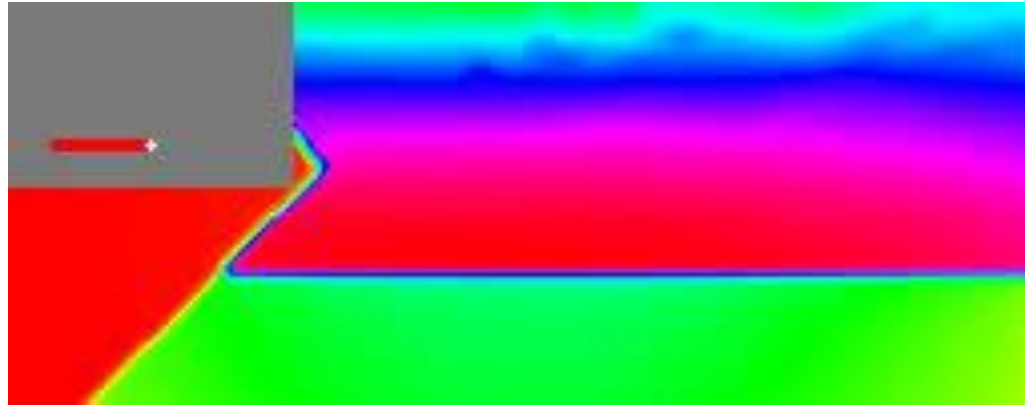
Policy is stochastic:

$$\pi_\theta \left(\text{img} \right) = \begin{bmatrix} .4 \\ .5 \\ .1 \end{bmatrix}$$



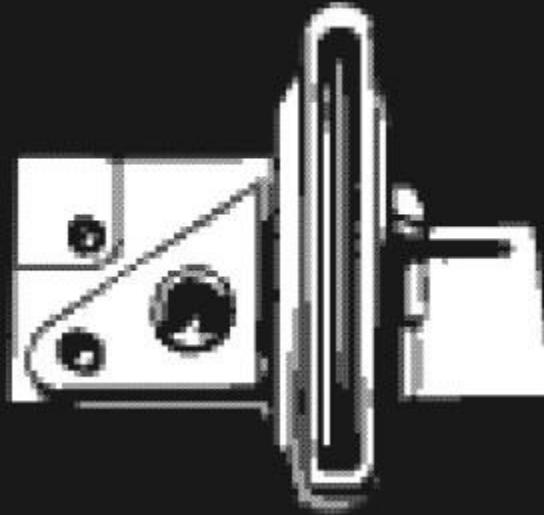
Reinforcement Learning Cube Example

- Using attribution visualization to understand decision making



Learning resources

- PyTorch
<http://pytorch.org/>
- Linear algebra for deep learning
<https://goo.gl/RJiNQJ>
- Calculus for deep learning
<https://goo.gl/zyQp7k>



THE INTERROGATOR ASKS YOU:

DO YOU DREAM ABOUT BEING
A UNICORN?

WHAT DO YOU SAY?

> **INTERLINKED.**
CELLS.