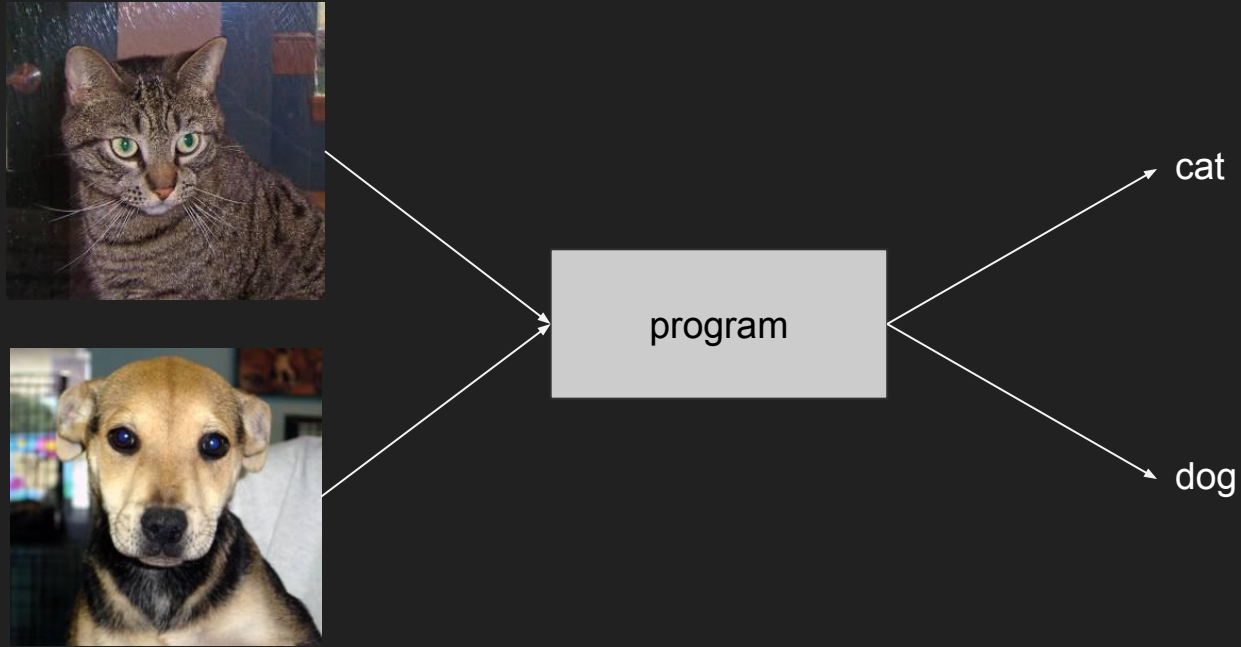


# Ceci n'est pas un chat

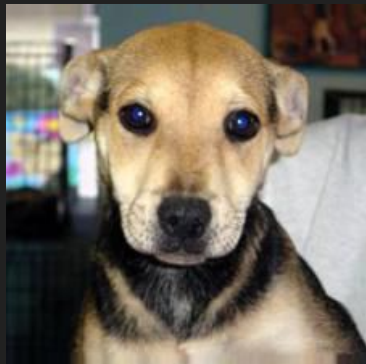
Convolutional neural networks from scratch  
and  
Adversarial examples

# First Goal: create a good cat vs dog classifier



# Second Goal: find adversary images that fool it

dog



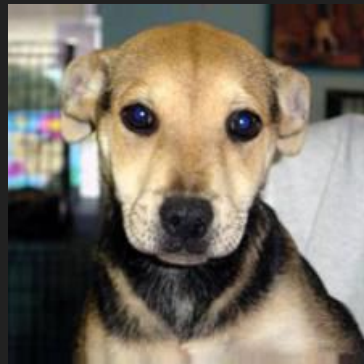
+  $\epsilon^*$

special noise



=

cat

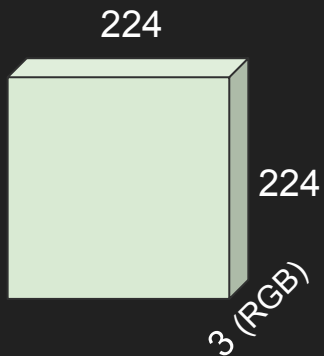


# It's all numbers

Input (150 528)



=



Output (2)

% being a cat

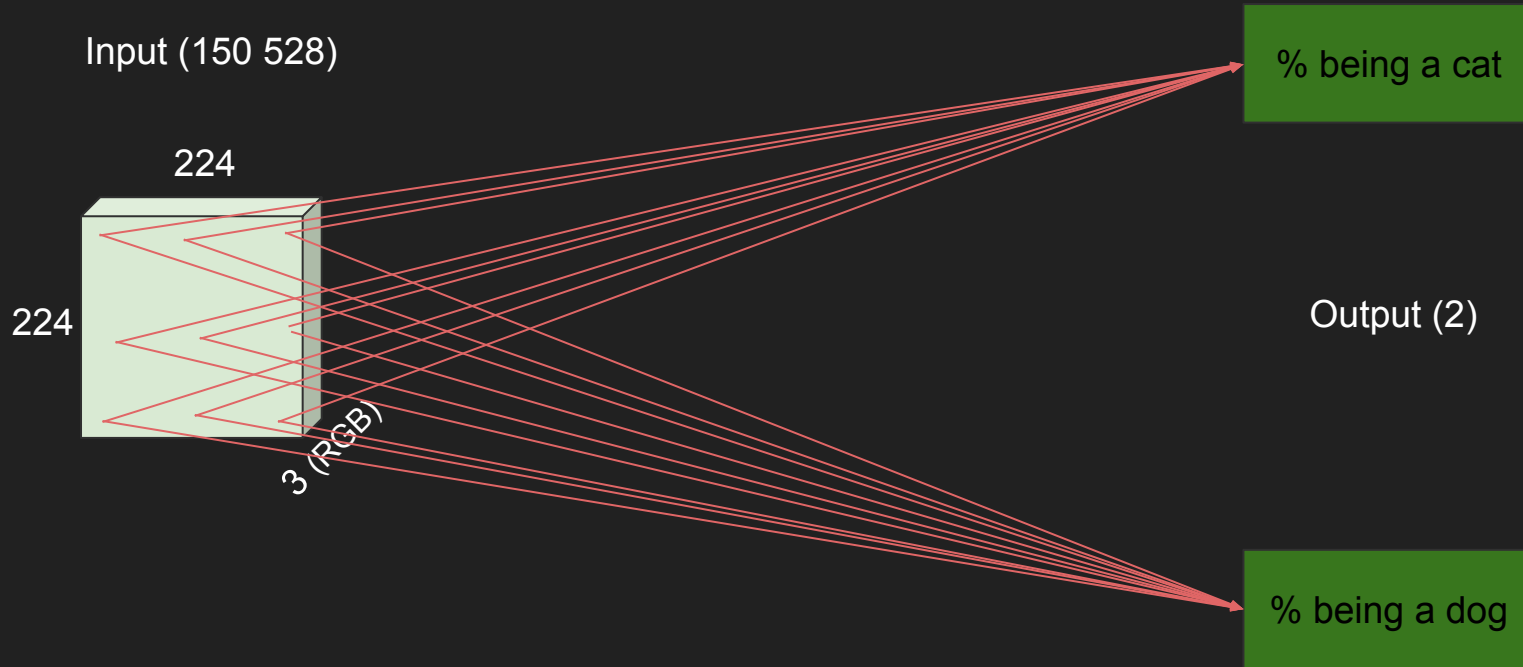
% being a dog

An awesome introduction to neural nets:

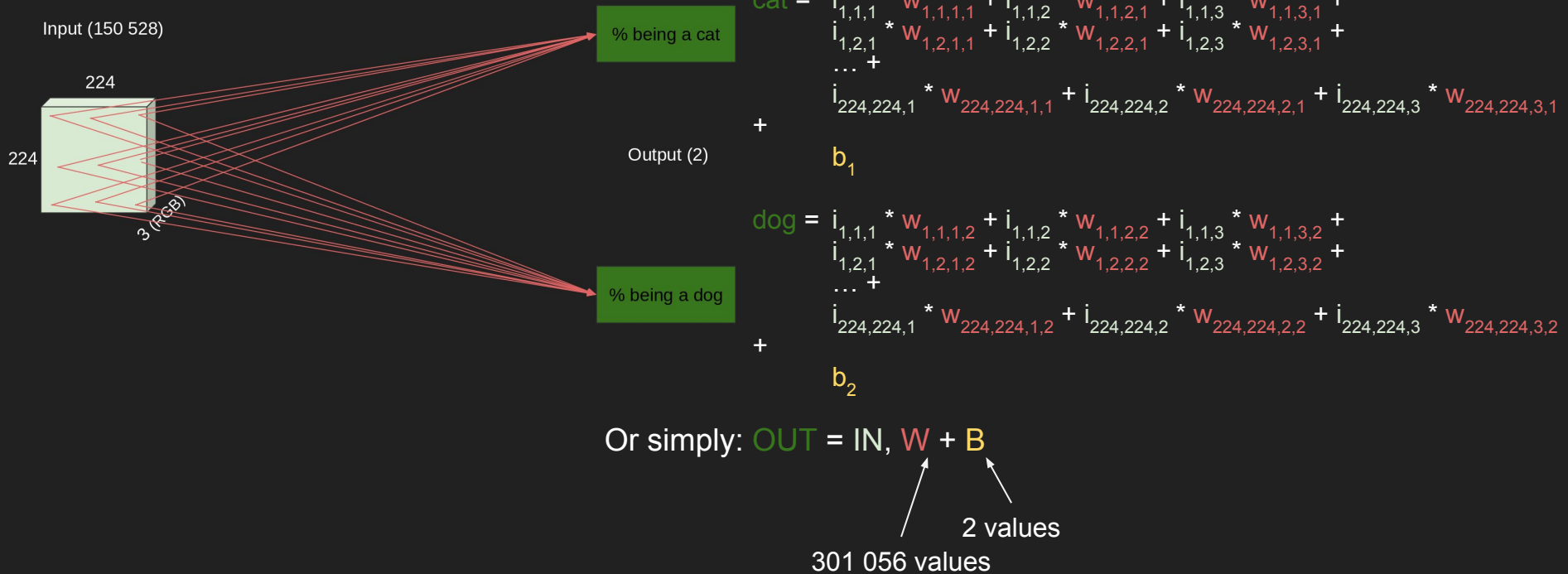
[3Blue1Brown: But what is a Neural Network?](#)

# Neural networks = MANY! simple operations

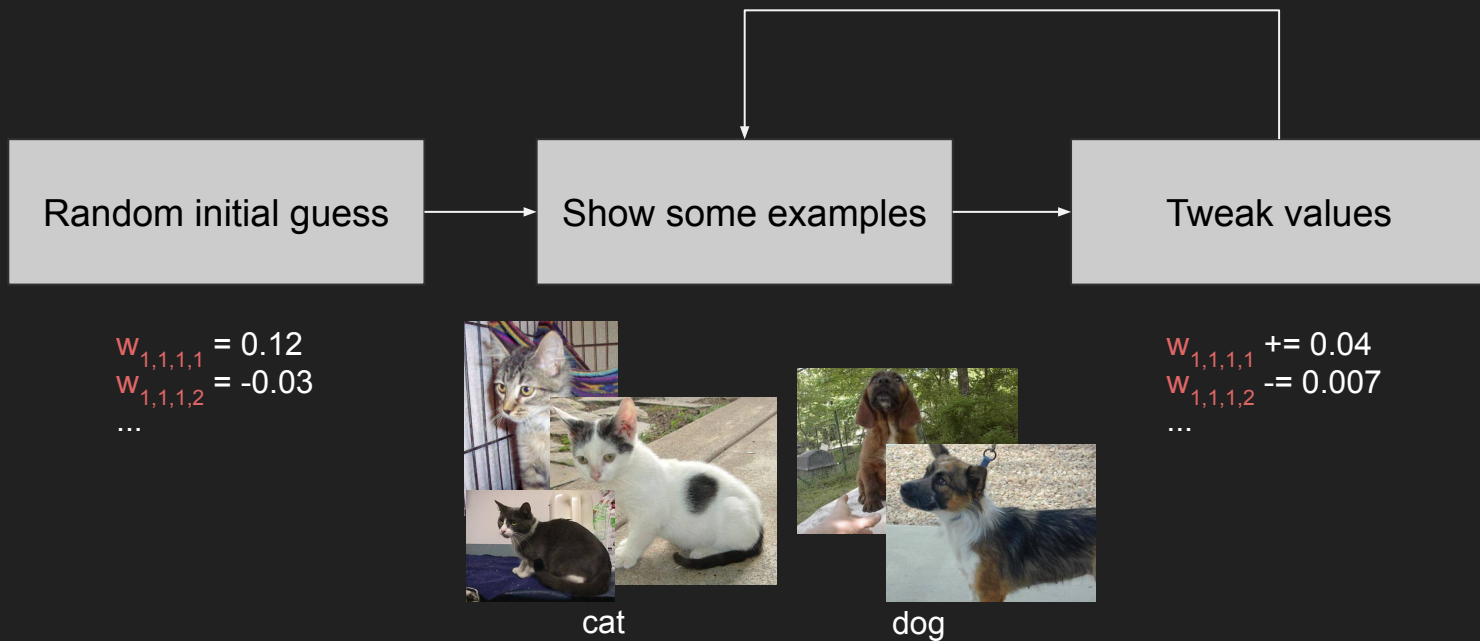
Connect each input pixel to each output



# Neural networks = MANY! simple operations



# Train a NN = find those values (**W** and **B**)



Keywords: stochastic gradient descent (SGD), back-propagation

Nice takes on the subjects:

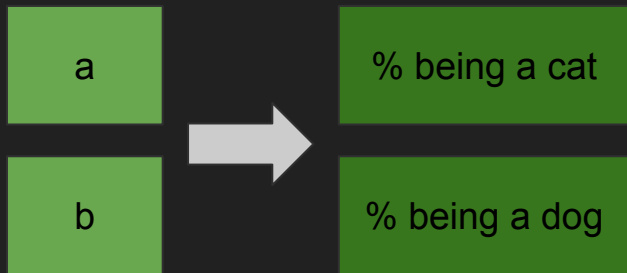
[3Blue1Brown: Gradient descent, how neural networks learn](#)

[colah's blog: Calculus on Computational Graphs: Backpropagation](#)

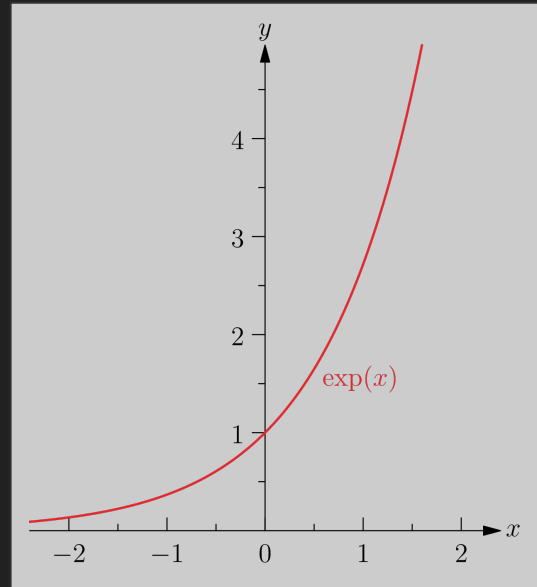
# Softmax: keep the output under control

We must guarantee:

- $\text{cat} + \text{dog} = 1$
- $0 \leq \text{cat}, \text{dog} \leq 1$



$$\text{cat} = e^a / (e^a + e^b)$$
$$\text{dog} = e^b / (e^a + e^b)$$

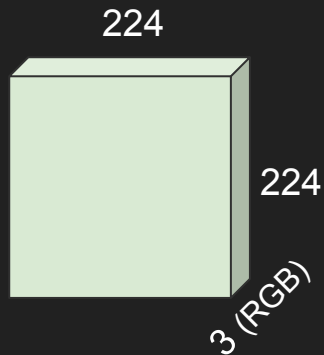


Or simply:  $\text{OUT} = \text{softmax}(\text{IN}, \text{W} + \text{B})$

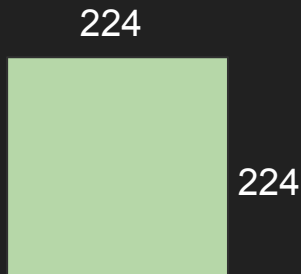


# Let's go deeper

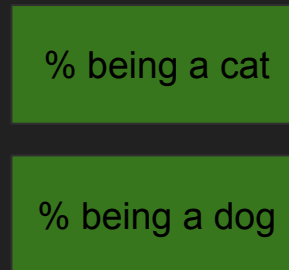
Input (150 528)



Hidden (50 176)



Output (2)



$$\text{HIDDEN} = \text{IN} \cdot W_1 + B_1$$
$$\text{OUT} = \text{softmax}(\text{HIDDEN} \cdot W_2 + B_2)$$

# Problem 1: linear collapse

$$\begin{aligned}\text{HIDDEN} &= \text{IN}, W_1 + B_1 \\ \text{OUT} &= \text{softmax}(\text{HIDDEN}, W_2 + B_2)\end{aligned}$$



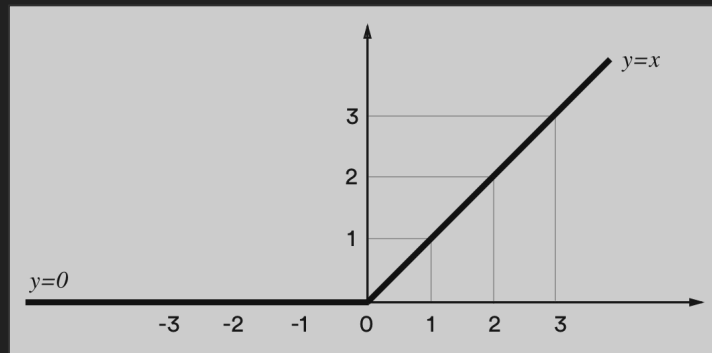
$$\text{OUT} = \text{softmax}(\text{IN}, \underbrace{(W_1, W_2)}_{\text{}} + \underbrace{(B_1, W_2 + B_2)}_{\text{}})$$

Exactly the same equation form without  
the hidden layer

# Solution: toss some non-linearity

Activation function

$\text{HIDDEN} = \text{relu}(\text{IN}, W_1 + B_1)$   
 $\text{OUT} = \text{softmax}(\text{HIDDEN}, W_2 + B_2)$



$$\text{relu}(x) = \max(0, x)$$

A mind-blowing dive into the effects of non-linearity:

[colah's blog: Neural Networks, Manifolds, and Topology](#)

# Problem 2: parameter explosion

Input (150 528)  Output (2)

$$\begin{aligned}\|W\| &= 150\,528 * 2 = 301\,056 \\ \|B\| &= 2\end{aligned}$$

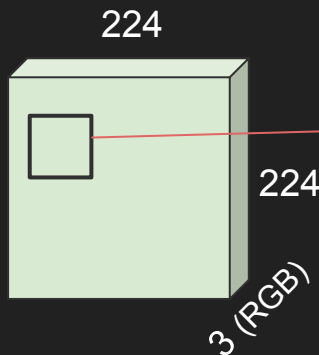
Input (150 528)  Hidden (50 176)  Output (2)

$$\begin{aligned}\|W\| &= 150\,528 * 50\,176 + 50\,176 * 2 = \underline{7\,552\,993\,280} \\ \|B\| &= 50\,176 + 2 = 50\,178\end{aligned}$$

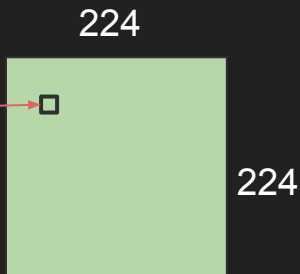
 Too complex to solve

# Solution - part 1: local connection

Input (150 528)



Hidden (50 176)



Connect each unit to only a neighboring input region, instead of all of it

Typical size: 3x3

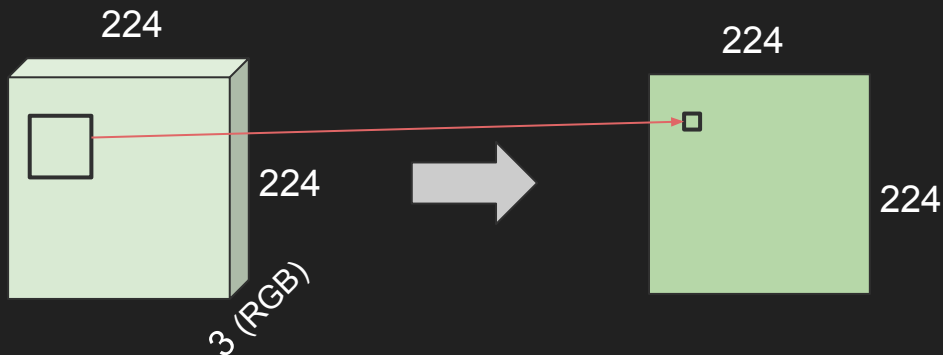
$$\|W\| = 50\,176 * 9 = 451\,584$$

$$\|B\| = 50\,176$$

# Solution - part 2: shared kernel

Input (150 528)

Hidden (50 176)



Use the same set of values for each unit (invariant to translation)

Typical size: 3x3

$$\|W\| = 9$$

$$\|B\| = 1$$

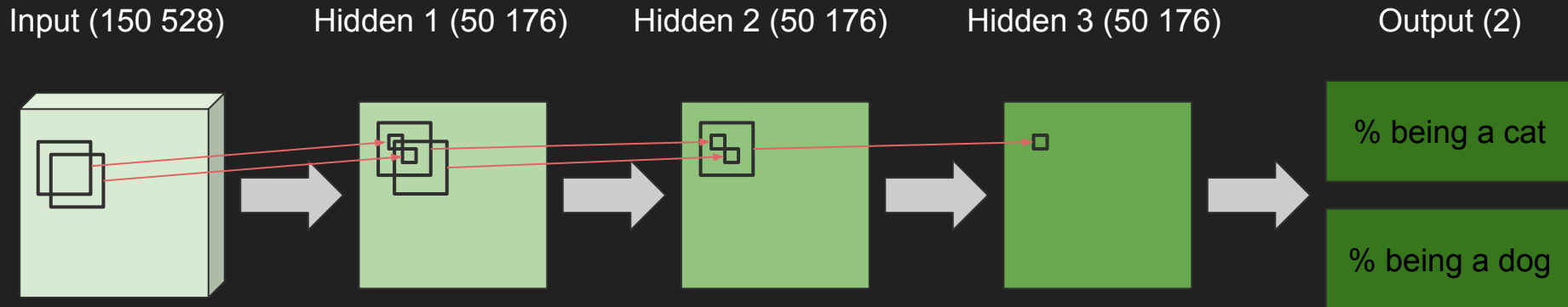
This is called  
**convolution**

Nice extras:

[Computerphile: How Blurs & Filters Work](#)

[Computerphile: Separable Filters and a Bauble](#)

# Let's go deeper



Each new layer has a more abstract representation of the input:

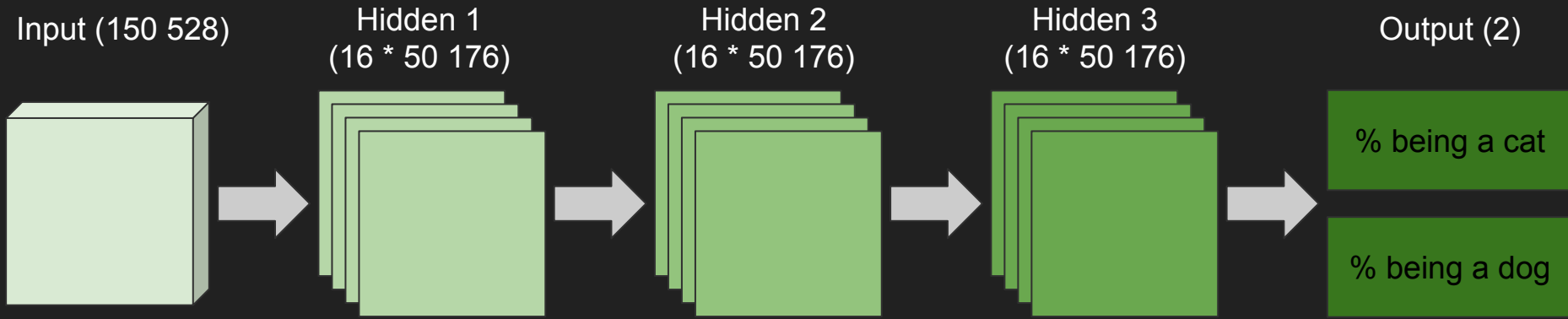
pixels

edge

corner

eye

# Let's go wider



Each new layer has **multiple** abstract representations of the input:

pixels

vertical edge  
horizontal edge  
brightness  
blueness

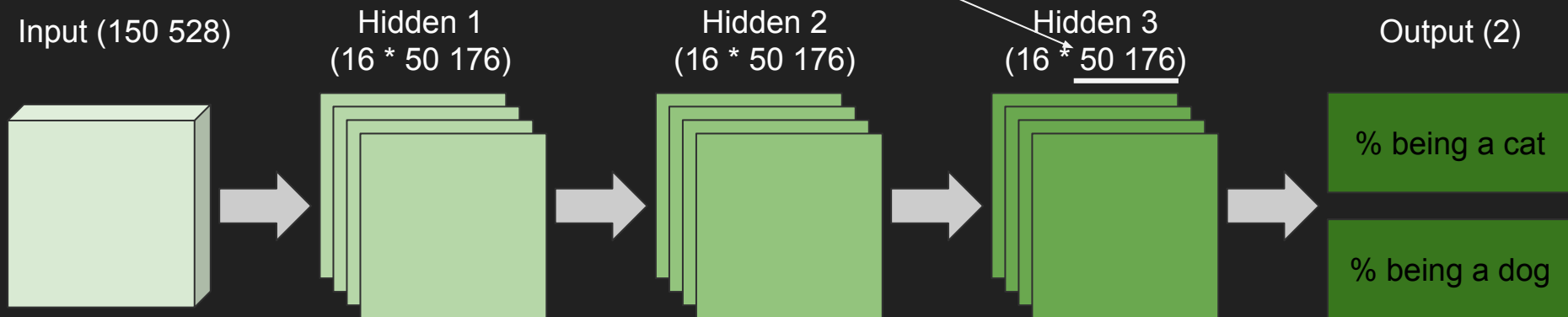
square  
circle  
triangle

eye  
mouth  
ear



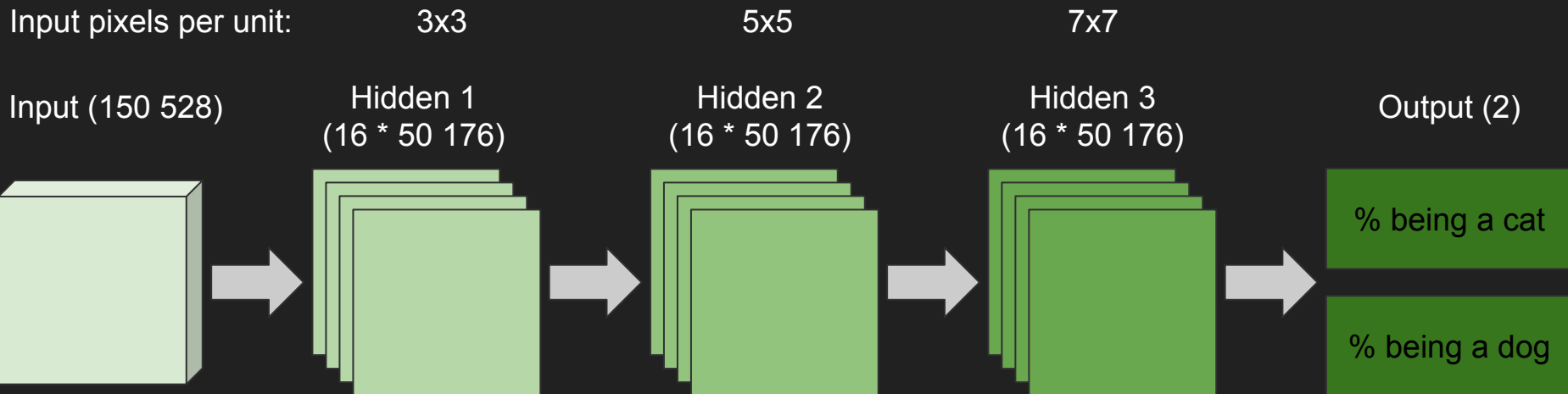
# Two problems

1: Too much detail at higher layers:  
no need to know that many times that an eye was detected



# Two problems

2: Too slow composition: it doesn't see the "big picture"



# Solution: group pixels (pooling)

Typical: that the max of each group of 2x2 units

Layer:	conv	pool	conv	pool	conv	pool	out
Size:	224x224x16	<b>112x112x16</b>	112x112x16	<b>56x56x16</b>	56x56x16	<b>28x28x16</b>	2
Input pixels per unit	3x3	<b>4x4</b>	8x8	<b>10x10</b>	18x18	<b>22x22</b>	224x224



# Let's code it



```
model = Sequential([
    Lambda(lambda x: x / 255 - 0.5, input_shape=(224, 224, 3)),
    Conv2D(16, (3, 3), padding='same', activation='relu'),
    MaxPool2D(),
    Conv2D(16, (3, 3), padding='same', activation='relu'),
    MaxPool2D(),
    Conv2D(16, (3, 3), padding='same', activation='relu'),
    MaxPool2D(),
    Flatten(),
    Dense(2, activation='softmax')
])
```

Model: "sequential"

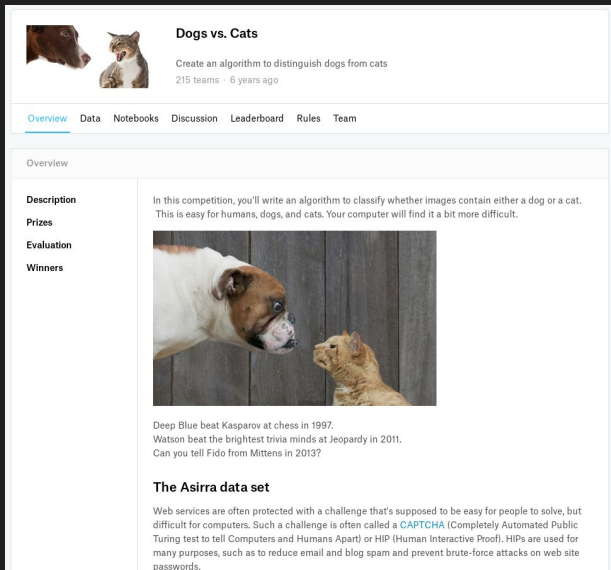
Layer (type)	Output Shape	Param #
lambda (Lambda)	(None, 224, 224, 3)	0
conv2d (Conv2D)	(None, 224, 224, 16)	448
max_pooling2d (MaxPooling2D)	(None, 112, 112, 16)	0
conv2d_1 (Conv2D)	(None, 112, 112, 16)	2320
max_pooling2d_1 (MaxPooling2D)	(None, 56, 56, 16)	0
conv2d_2 (Conv2D)	(None, 56, 56, 16)	2320
max_pooling2d_2 (MaxPooling2D)	(None, 28, 28, 16)	0
flatten (Flatten)	(None, 12544)	0
dense (Dense)	(None, 2)	25090

Total params: 30,178  
Trainable params: 30,178  
Non-trainable params: 0

Documentation: [tf.keras.Sequential](https://keras.io/api/sequential/)

Repository: [cecil-nest-pas-un-chat](https://github.com/cecil-nest-pas-un-chat)

# Let's prepare our data




**Dogs vs. Cats**  
Create an algorithm to distinguish dogs from cats  
215 teams · 6 years ago

[Overview](#) [Data](#) [Notebooks](#) [Discussion](#) [Leaderboard](#) [Rules](#) [Team](#)

Overview

**Description**  
**Prizes**  
**Evaluation**  
**Winners**

In this competition, you'll write an algorithm to classify whether images contain either a dog or a cat. This is easy for humans, dogs, and cats. Your computer will find it a bit more difficult.

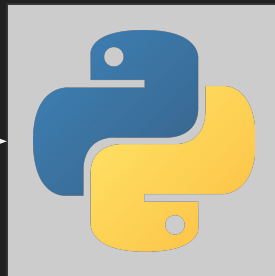


Deep Blue beat Kasparov at chess in 1997.  
Watson beat the brightest trivia minds at Jeopardy in 2011.  
Can you tell Fido from Mittens in 2013?

**The Asirra data set**

Web services are often protected with a challenge that's supposed to be easy for people to solve, but difficult for computers. Such a challenge is often called a **CAPTCHA** (Completely Automated Public Turing test to tell Computers and Humans Apart) or HIP (Human Interactive Proof). HIPs are used for many purposes, such as to reduce email and blog spam and prevent brute-force attacks on web site passwords.

Split, resize and crop



Train (15 000)

Validation (7 500)

Test (2 500)



# Let's train it

```
base_gen = ImageDataGenerator()
```

```
train_generator = base_gen.flow_from_directory('data/training', (224, 224))
```

```
validation_generator = base_gen.flow_from_directory('data/validation', (224, 224))
```

```
test_generator = base_gen.flow_from_directory('data/test', (224, 224))
```

```
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])
```

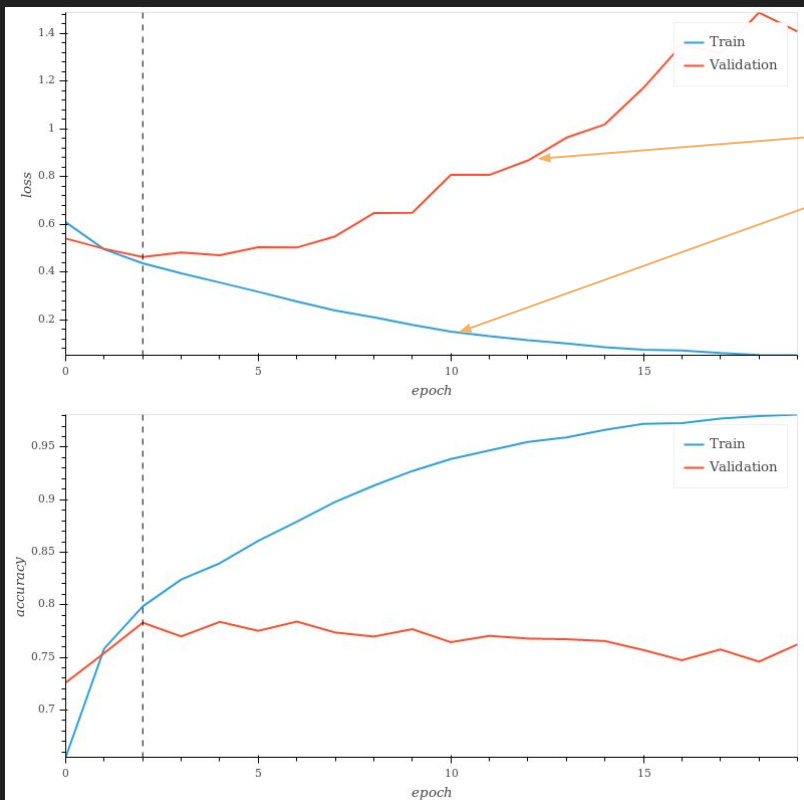
```
model.fit_generator(train_generator, epochs=20, validation_data=validation_generator)
```

pass through the  
training set 20x

=  $-\log(\text{right\_class})$

=  $\text{right\_answers} / \text{total}$

# Training stats



Overfit: model “memorizes” answers for test set, and generalization is compromised

Use validation to know when to stop:

```
model.fit_generator(train_generator, epochs=1000,  
validation_data=validation_generator, callbacks=[  
    EarlyStopping(patience=1, restore_best_weights=True)  
])
```

More robust techniques: dropout, batch normalization, residual convolution

# Results over test set



dog at 100%



dog at 100%



dog at 100%



cat at 100%



dog at 99%



cat at 99%



dog at 99%



dog at 99%



	predicted		
real		cat	dog
	cat	1013	237
	dog	270	980



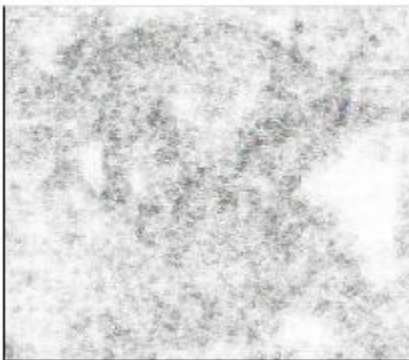
# Cool! Let's break it

Goal: modify an image to manipulate the result

**Initial dog at 83%**



**+ 0.032x**



**= cat at 95% - end**



# In the real world



Mistaken by 45 mi/h speed limit  
<https://arxiv.org/abs/1707.08945>

Active research area

Open-source collection of attacks and defences:  
<https://github.com/tensorflow/cleverhans>

# How?

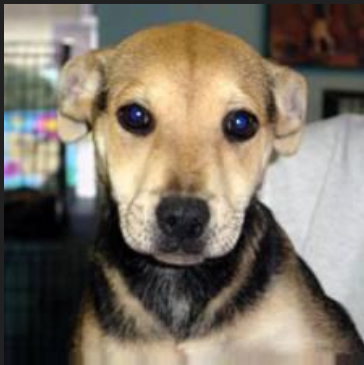


image + 255 \* noise



base\_model

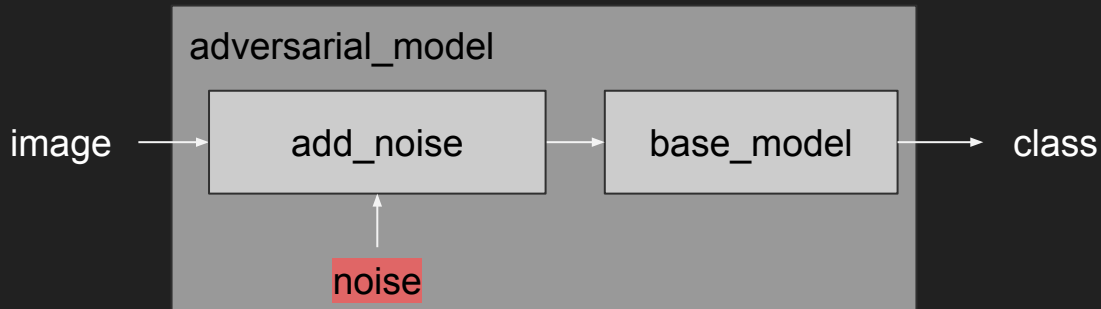


target wrong classification

cat

$$\text{loss} = 0.1 * l_2(\text{noise}) - \log(\text{target\_wrong\_category})$$

# How?

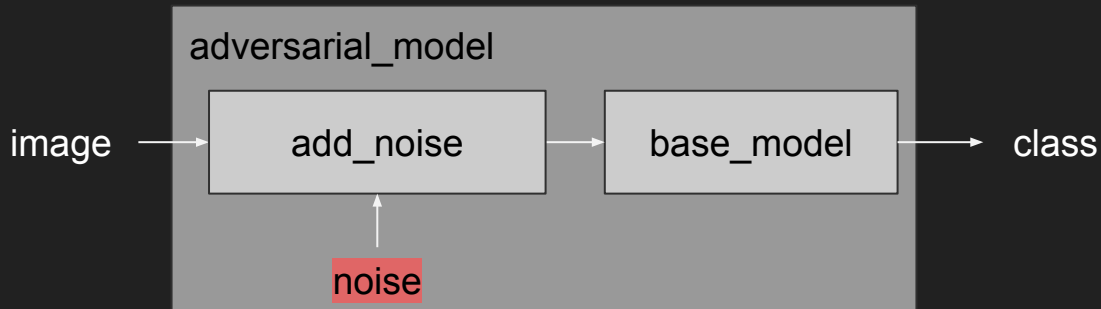


```
base_model = load_model('models/simple')
base_model.trainable = False
```

```
base_img_input = Input((224, 224, 3))
adversarial_img_layer = AddNoise(12(0.1), 255)
adversarial_img = adversarial_img_layer(base_img_input)
model_output = base_model(adversarial_img)
adversarial_model = Model(inputs=base_img_input, outputs=model_output)
```

control the tradeoff  
between image similarity  
and fooling power

# How?



```
class AddNoise(Layer):  
    def __init__(self, regularizer=None, scale=1, **kwargs):  
        self.scale = scale  
        self.regularizer = regularizer  
        super().__init__(**kwargs)  
    def build(self, input_shape):  
        self.noise = self.add_weight('noise', input_shape[1:], regularizer=self.regularizer)  
        super().build(input_shape)  
    def call(self, inputs):  
        return tf.clip_by_value(inputs + self.scale * self.noise, 0, 255)
```

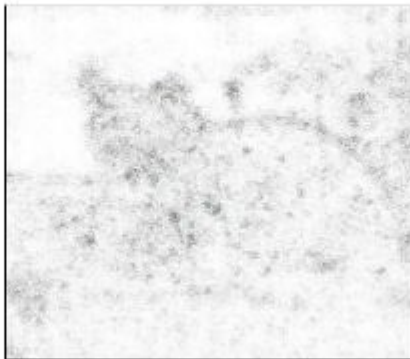
# How?

```
adversarial_model.compile(optimizer='adam', loss='categorical_crossentropy')
adversarial_model.fit(x=base_img, y=target_wrong_class, epochs=1000, callbacks=[
    EarlyStopping(monitor='loss', patience=10, restore_best_weights=True)
])
```

**Initial cat at 94%**



**+ 0.071x**



**= dog at 90% - end**



# Extra: transfer learning with VGG16

```
model_input = Input((224, 224, 3))
base_model_input = Lambda(vgg16.preprocess_input)(model_input)

base_model = VGG16(include_top=False,
                    input_tensor=base_model_input,
                    pooling='avg')
base_model.trainable = False

x = base_model.output
x = Dense(128, activation='relu', name='fc1')(x)
x = Dense(16, activation='relu', name='fc2')(x)
x = Dense(2, activation='softmax', name='predictions')(x)

model = Model(model_input, x)
```

Model: "model"		
Layer (type)	Output Shape	Param #
=====		
input_1 (InputLayer)	[(None, 224, 224, 3)]	0
lambda (Lambda)	(None, 224, 224, 3)	0
block1_conv1 (Conv2D)	(None, 224, 224, 64)	1792
block1_conv2 (Conv2D)	(None, 224, 224, 64)	36928
block1_pool (MaxPooling2D)	(None, 112, 112, 64)	0
block2_conv1 (Conv2D)	(None, 112, 112, 128)	73856
block2_conv2 (Conv2D)	(None, 112, 112, 128)	147584
block2_pool (MaxPooling2D)	(None, 56, 56, 128)	0
block3_conv1 (Conv2D)	(None, 56, 56, 256)	295168
block3_conv2 (Conv2D)	(None, 56, 56, 256)	590080
block3_conv3 (Conv2D)	(None, 56, 56, 256)	590080
block3_pool (MaxPooling2D)	(None, 28, 28, 256)	0
block4_conv1 (Conv2D)	(None, 28, 28, 512)	1180160
block4_conv2 (Conv2D)	(None, 28, 28, 512)	2359808
block4_conv3 (Conv2D)	(None, 28, 28, 512)	2359808
block4_pool (MaxPooling2D)	(None, 14, 14, 512)	0
block5_conv1 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv2 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv3 (Conv2D)	(None, 14, 14, 512)	2359808
block5_pool (MaxPooling2D)	(None, 7, 7, 512)	0
global_average_pooling2d (Gl	(None, 512)	0
fc1 (Dense)	(None, 128)	65664
fc2 (Dense)	(None, 16)	2064
predictions (Dense)	(None, 2)	34
=====		
Total params: 14,782,450		
Trainable params: 67,762		
Non-trainable params: 14,714,688		

# Quite of an improvement!

previous model  
accuracy = 79.7%

	predicted		
real		cat	dog
	cat	1013	237
	dog	270	980

VGG16  
accuracy = 97.8%

	predicted		
real		cat	dog
	cat	1223	27
	dog	28	1222

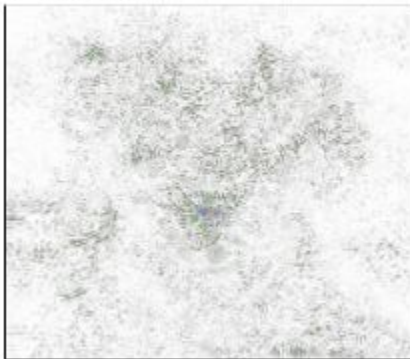


# Let's break it

**Initial dog at 100%**



**+ 0.046x**



**= cat at 99% - end**

