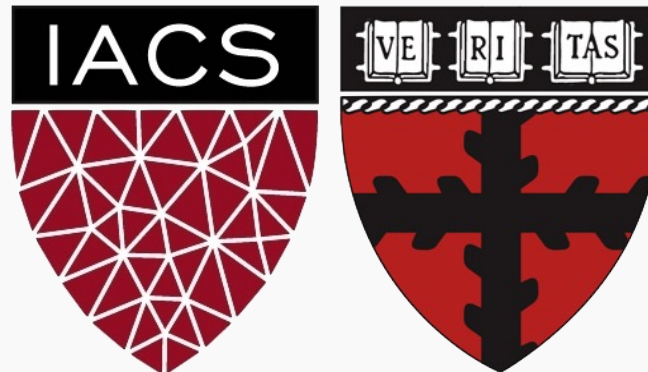# Convolutional Neural Networks III

## CS109B Data Science 2
## Pavlos Protopapas, Mark Glickman

# Outline

1. Regularization for CNN

2. BackProp of MaxPooling layer

3. Layers Receptive Field and dilated convolutions

4. Weights and feature maps visualization

# Outline

1. **Regularization for CNN**

2. BackProp of MaxPooling layer

3. Layers Receptive Field and dilated convolutions

4. Weights and feature maps visualization
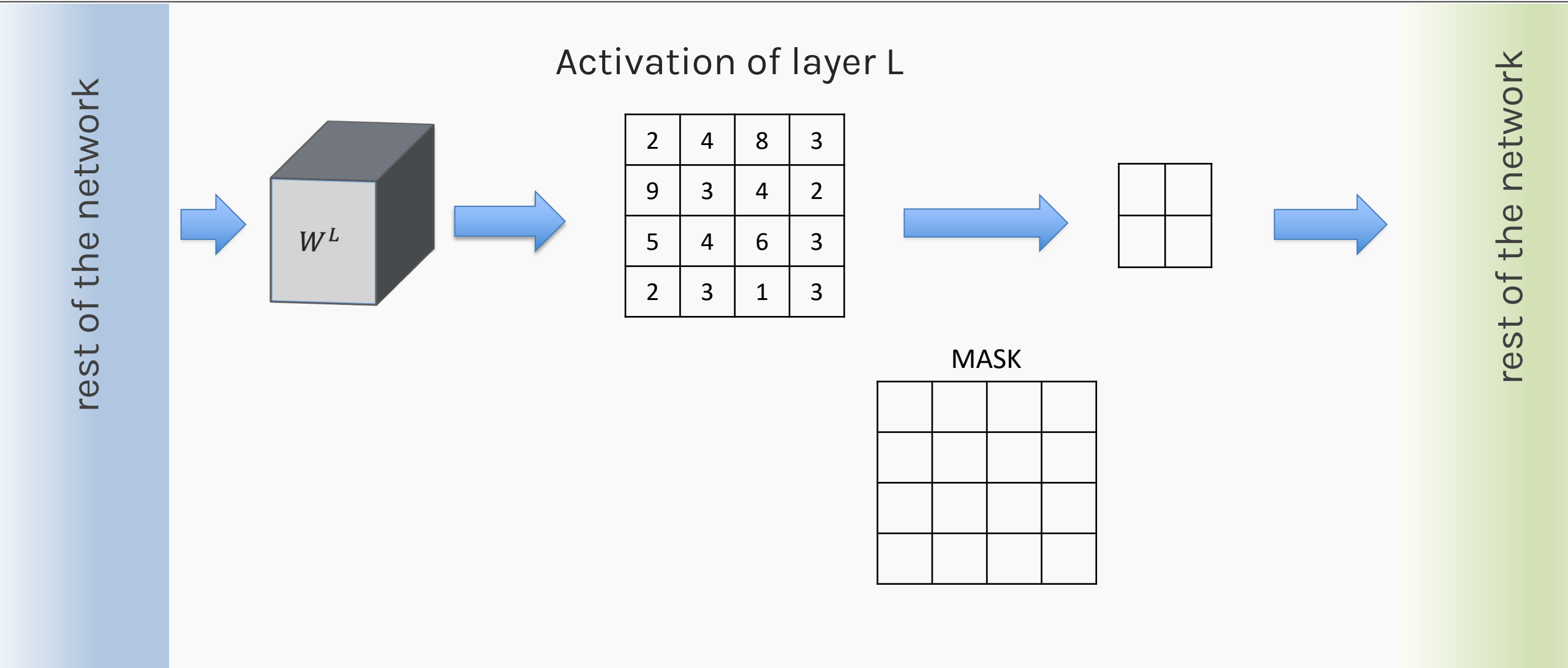
# Regularization for CNN

- L2 and L1 work the same way as in FFNN
- Data Augmentation is the same
- Early Stopping same as in FFNN

- Dropout is slightly different – not the same effect as dropout with FFNN.
  - Dropout in CNN still allows the weights in a kernel to be trained.
  - The name is misleading!
  - The effect of dropout on convolutional layers amounts to multiplying Bernoulli noise into the feature maps of the network.

So, if you try adding dropout after a convolutional layer and get bad results, don't be disappointed! There doesn't appear that there is a good reason it *should* provide good results.
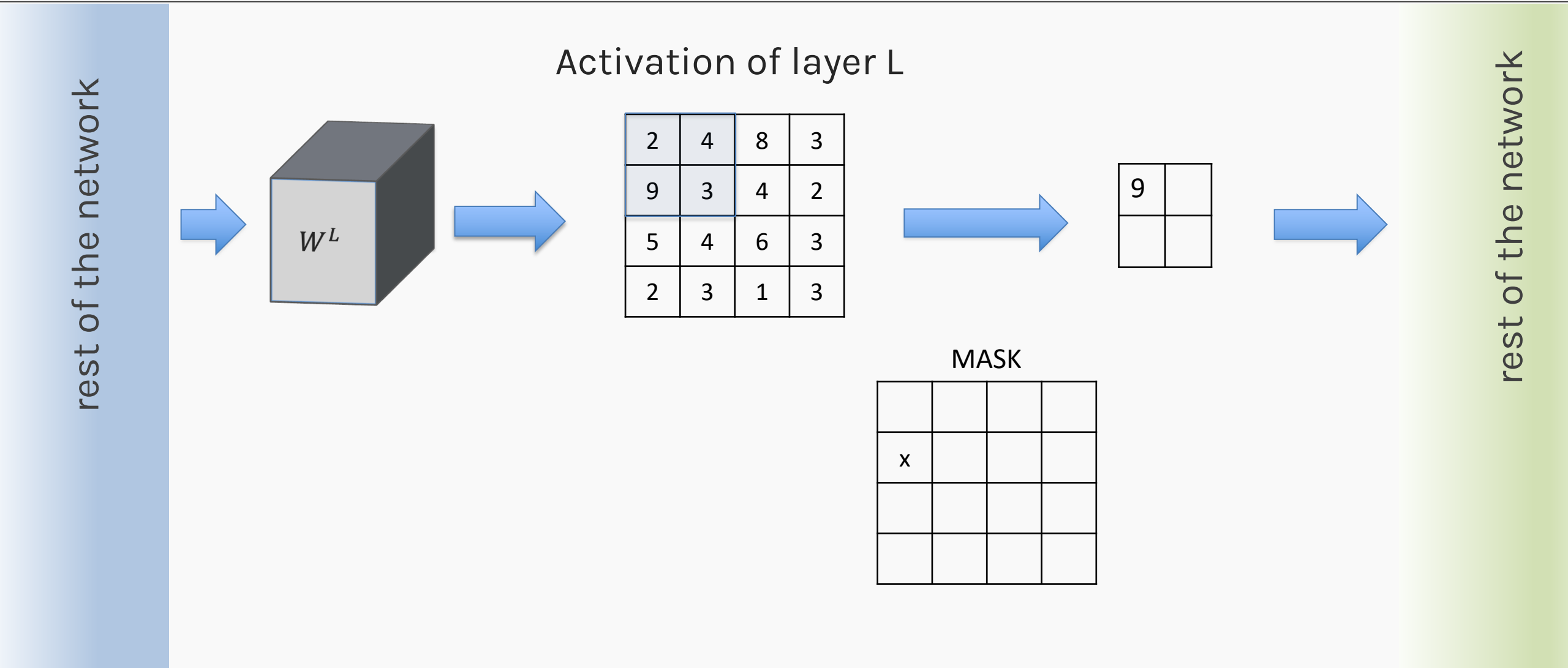
# Outline

1. Regularization for CNN

2. **BackProp of MaxPooling layer**

3. Layers Receptive Field and dilated convolutions

4. Weights and feature maps visualization
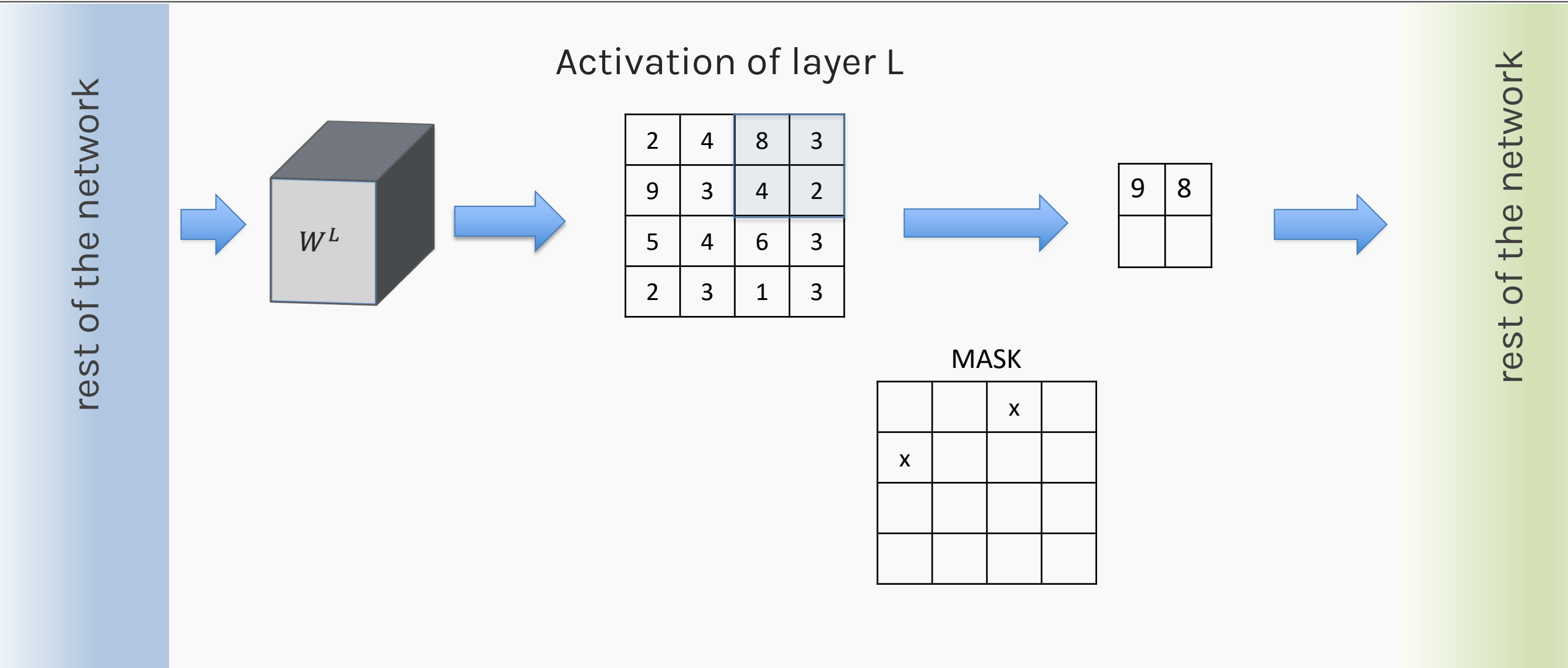
# Backward propagation of Maximum Pooling Layer

Activation of layer L



| 2 | 4 | 8 | 3 |
|---|---|---|---|
| 9 | 3 | 4 | 2 |
| 5 | 4 | 6 | 3 |
| 2 | 3 | 1 | 3 |

MASK

Forward mode, 2x2 stride 2x2

# Backward propagation of Maximum Pooling Layer

Activation of layer L



| 2 | 4 | 8 | 3 |
|---|---|---|---|
| 9 | 3 | 4 | 2 |
| 5 | 4 | 6 | 3 |
| 2 | 3 | 1 | 3 |

| 9 | |
|---|---|
| | |

MASK

| | | | |
|---|---|---|---|
| x | | | |
| | | | |
| | | | |

$W^L$

Forward mode, 2x2 stride 2x2

# Backward propagation of Maximum Pooling Layer

Activation of layer L

$W^L$

| 2 | 4 | 8 | 3 |
|---|---|---|---|
| 9 | 3 | 4 | 2 |
| 5 | 4 | 6 | 3 |
| 2 | 3 | 1 | 3 |

| 9 | 8 |
|---|---|
|   |   |

MASK

|   |   | x |   |
|---|---|---|---|
| x |   |   |   |
|   |   |   |   |
|   |   |   |   |

rest of the network

rest of the network

Forward mode, 2x2 stride 2x2

# Backward propagation of Maximum Pooling Layer

Activation of layer L



$W^L$

| 2 | 4 | 8 | 3 |
|---|---|---|---|
| 9 | 3 | 4 | 2 |
| 5 | 4 | 6 | 3 |
| 2 | 3 | 1 | 3 |

| 9 | 8 |
|---|---|
| 5 |   |

MASK

|   |   | x |   |
|---|---|---|---|
| x |   |   |   |
| x |   |   |   |
|   |   |   |   |

Forward mode, 2x2 stride 2x2

rest of the network

Activation of layer L



| 2 | 4 | 8 | 3 |
|---|---|---|---|
| 9 | 3 | 4 | 2 |
| 5 | 4 | 6 | 3 |
| 2 | 3 | 1 | 3 |

| 9 | 8 |
|---|---|
| 5 | 6 |

MASK

|   |   | x |   |
|---|---|---|---|
| x |   |   |   |
| x |   | x |   |
|   |   |   |   |

$W^L$

rest of the network

Forward mode, 2x2 stride 2x2

# Backward propagation of Maximum Pooling Layer

Activation of layer L

$W^L$

| 2 | 4 | 8 | 3 |
|---|---|---|---|
| 9 | 3 | 4 | 2 |
| 5 | 4 | 6 | 3 |
| 2 | 3 | 1 | 3 |

| 9 | 8 |
|---|---|
| 5 | 6 |

MASK

|   |   | x |   |
|---|---|---|---|
| x |   |   |   |
| x |   | x |   |
|   |   |   |   |

DERIVATIVES

| 1 | 2 |
|---|---|
| 2 | 5 |

|   |   |   |   |
|---|---|---|---|
|   |   |   |   |
|   |   |   |   |
|   |   |   |   |

Reverse mode

Activation of layer L

| | | | |
|---|---|---|---|
| 2 | 4 | 8 | 3 |
| 9 | 3 | 4 | 2 |
| 5 | 4 | 6 | 3 |
| 2 | 3 | 1 | 3 |

| | |
|---|---|
| 9 | 8 |
| 5 | 6 |

MASK    DERIVATIVES

| | | | |
|---|---|---|---|
| 0 | 0 | | |
| 1 | 0 | | |
| | | | |
| | | | |

| | | | |
|---|---|---|---|
| | | x | |
| x | | | |
| x | | x | |
| | | | |

| | |
|---|---|
| 1 | 2 |
| 2 | 5 |

Reverse mode

# Backward propagation of Maximum Pooling Layer

Activation of layer L

$W^L$

| 2 | 4 | 8 | 3 |
|---|---|---|---|
| 9 | 3 | 4 | 2 |
| 5 | 4 | 6 | 3 |
| 2 | 3 | 1 | 3 |

| 9 | 8 |
|---|---|
| 5 | 6 |

MASK

DERIVATIVES

| 0 | 0 | 2 | 0 |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
|   |   |   |   |
|   |   |   |   |

| | | x | |
|---|---|---|---|
| x | | | |
| x | | x | |
| | | | |

| 1 | 2 |
|---|---|
| 2 | 5 |

Reverse mode

Activation of layer L

| 2 | 4 | 8 | 3 |
|---|---|---|---|
| 9 | 3 | 4 | 2 |
| 5 | 4 | 6 | 3 |
| 2 | 3 | 1 | 3 |

| 9 | 8 |
|---|---|
| 5 | 6 |

MASK

DERIVATIVES

| 0 | 0 | 2 | 0 |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| 2 | 0 |   |   |
| 0 | 0 |   |   |

|   |   | x |   |
|---|---|---|---|
| x |   |   |   |
| x |   | x |   |
|   |   |   |   |

| 1 | 2 |
|---|---|
| 2 | 5 |

rest of the network

rest of the network

$W^L$

Reverse mode

Activation of layer L



| 2 | 4 | 8 | 3 |
|---|---|---|---|
| 9 | 3 | 4 | 2 |
| 5 | 4 | 6 | 3 |
| 2 | 3 | 1 | 3 |

| 9 | 8 |
|---|---|
| 5 | 6 |

MASK

DERIVATIVES

| 0 | 0 | 2 | 0 |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| 2 | 0 | 5 | 0 |
| 0 | 0 | 0 | 0 |

| | | x | |
|---|---|---|---|
| x | | | |
| x | | x | |
| | | | |

| 1 | 2 |
|---|---|
| 2 | 5 |

Reverse mode

rest of the network

rest of the network

Activation of layer L

| 2 | 4 | 8 | 3 |
|---|---|---|---|
| 9 | 3 | 4 | 2 |
| 5 | 4 | 6 | 3 |
| 2 | 3 | 1 | 3 |

| 9 | 8 |
|---|---|
| 5 | 6 |

MASK

| | | x | |
|---|---|---|---|
| x | | | |
| x | | x | |
| | | | |

DERIVATIVES

| 1 | 2 |
|---|---|
| 2 | 5 |

| 0 | 0 | 2 | 0 |
|---|---|---|---|
| 1 | 0 | 0 | 0 |
| 2 | 0 | 5 | 0 |
| 0 | 0 | 0 | 0 |

$W^L$

rest of the network

rest of the network

Reverse mode

# Outline

1. Regularization for CNN

2. BackProp of MaxPooling layer

3. **Layers Receptive Field and dilated convolutions**

4. Weights and feature maps visualization

# Layers Receptive Field

The **receptive field** is defined as the region in the input space that a particular CNN's feature (or activation) is looking at (i.e. be affected by).

The receptive field size is a crucial issue in many visual tasks, as the output must respond to large enough areas in the image to capture information about large objects.

# Layers Receptive Field

The **receptive field** is defined as the region in the input space that a particular CNN's feature (or activation) is looking at (i.e. be affected by).

Let's look at the receptive field in 1D, no padding, stride 1 and kernel 3x1

# Layers Receptive Field

The **receptive field** is defined as the region in the input space that a particular CNN's feature (or activation) is looking at (i.e. be affected by).

Let's look at the receptive field in 1D, no padding, stride 1 and kernel 3x1

# Layers Receptive Field

The **receptive field** is defined as the region in the input space that a particular CNN's feature (or activation) is looking at (i.e. be affected by).

Let's look at the receptive field in 1D, no padding, stride 1 and kernel 3x1

# Layers Receptive Field

The **receptive field** is defined as the region in the input space that a particular CNN's feature (or activation) is looking at (i.e. be affected by).

Let's look at the receptive field in 1D, no padding, stride 1 and kernel 3x1



Layer 1

$a_1^1$   $a_2^1$   $a_3^1$

Input

# Layers Receptive Field

The **receptive field** is defined as the region in the input space that a particular CNN's feature (or activation) is looking at (i.e. be affected by).

Let's look at the receptive field in 1D, no padding, stride 1 and kernel 3x1

# Layers Receptive Field

The receptive field for each element of layer's 2 is shown below.

# Layers Receptive Field

The receptive field for each element of layer's 2 is shown below.

# Layers Receptive Field

The receptive field for each element of layer's 2 is shown below.

# Layers Receptive Field

The receptive field for each element of layers 1 and 2 are shown below.

# Layers Receptive Field

The receptive field for each element of layers 1 and 2 are shown below.

# Layers Receptive Field

In 2D, it works the same way.

The receptive field can be calculated using the recursive formula:

$$r_0 = 1 + \sum_{l=1}^{L}(k_l - 1)\prod_{i=1}^{l-1} s_i$$

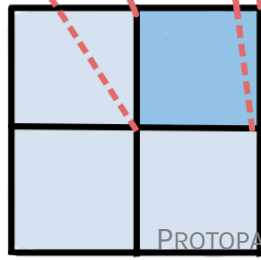- $k_l$ kernel size (positive integer)
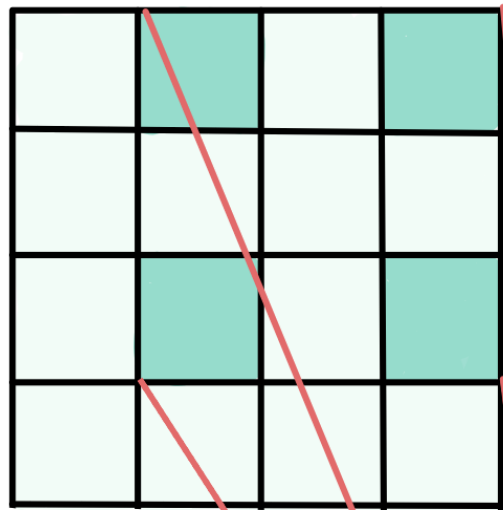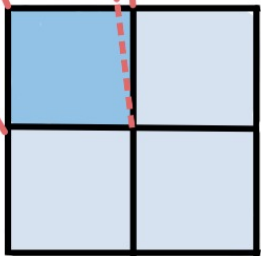- $s_l$ stride (positive integer)

# Dilated CNNs
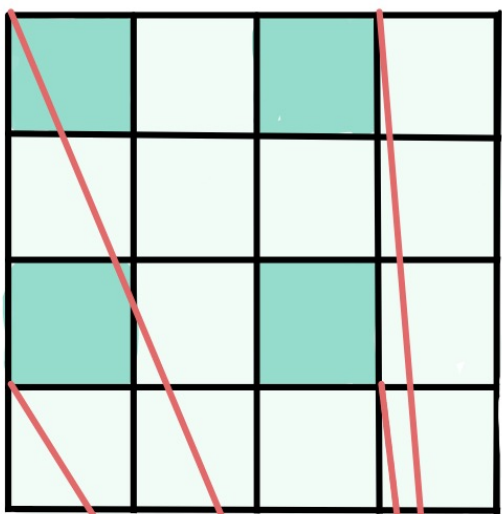
- We can "inflate" the receptive field by inserting holes between the kernel elements.

- These are called **Dilated Convolutions.**

- Dilation rate indicates how much the kernel is widened.

Dilate rate=1



*Original Idea: Algorithme a trous*, an algorithm for wavelet decomposition (Holschneider et al., **1987**; Shensa, 1992)

# Dilated CNNs

- We can "inflate" the receptive field by inserting holes between the kernel elements.

- These are called **Dilated Convolutions.**

- Dilation rate indicates how much the kernel is widened.

Dilate rate=1

# Dilated CNNs

- We can "inflate" the receptive field by inserting holes between the kernel elements.

- These are called **Dilated Convolutions.**

- Dilation rate indicates how much the kernel is widened.

Dilate rate=1

# Dilated CNNs

- We can "inflate" the receptive field by inserting holes between the kernel elements.

- These are called **Dilated Convolutions.**

- Dilation rate indicates how much the kernel is widened.

Dilate rate=1

# Dilated CNNs

**2D Example:** 2x2 kernel, stride=1, dilate rate=1

# Dilated CNNs

There is a relationship between classification accuracy and receptive field size.

Large receptive fields are necessary for high-level recognition tasks, but with diminishing rewards.

Araujo, A., Norris, W., & Sim, J. (2019). Computing receptive fields of convolutional neural networks. *Distill*, *4*(11), e21.

# Outline

1. Regularization for CNN

2. BackProp of MaxPooling layer

3. Layers Receptive Field

4. **Weights and feature maps visualization**

# Lessons for Visualization

Choosing/designing machine learning visualization requires that we think about:

***Why and for whom* to visualize?**

- – are we visualizing to **diagnose** problems with our models?
- – are we visualizing to **interpret** our model's meaningfulness?

***What and how to visualize?***

- – do we visualize decision boundaries, weights of our model, and or distributional differences in the data?

# Why and for whom to visualize?

1. **Interpretability & Explainability:** understand how deep learning models make decisions and what representations they have learned.

For others

2. **Debugging & Improving Models:** help model developers build and debug their models, with the hope of expediting the iterative experimentation process to ultimately improve performance.

For us

3. **Teaching Deep Learning Concepts:** educate non-expert users about AI

For me

From: *Visual Analytics in Deep Learning: An Interrogative Survey for the Next Frontiers*

# What and how to visualize?

**What technical components of neural networks could be visualized?**

- Computational Graph & Network Architecture
- Learned Model Parameters: weights, filters
- Individual Computational Units: activations, gradients
- Aggregate information: performance metrics

**How can they be insightfully visualized?**

How depends on the type of data and model as well as our specific investigative goal.

# What to Visualize for Neural Network Models?

For logistic regression, $p(y = 1|w, x) = \sigma(w^T x)$ we can interrogate the model by printing out the weights of the model.

Recalling from previous lectures, we can visualize the feature importance looking at the coefficients

$$\ln\left(\frac{P(y = 1)}{P(y = 0)}\right) = w^T x$$

log-odds



FEATURE IMPORTANCE

RAINFALL
% IRRIGATED
TOTAL GROWING DAYS
MEAN AIR TEMPERATURE
HUMIDITY
DAILY RADIATION
ORGANIC COMPOUNDS
FERTILIZER APPLIED
SOIL pH
% SILT
% CLAY
CONDUCTIVITY
MEAN COLDEST
ELEVATION
SOIL DENSITY
ERODIBILITY
% SAND
SUMMER SOLSTICE DAY LENGTH
WATER TABLE DEPTH
WIND SPEED
MEAN WARMEST
# OF DAYS OF ABNORMAL TEMP.
# OF WILD FIRES
DISTANCE TO CITY
ENERGY CONSUMPTION
LONGITUDE
LADITUDE
TOTAL POPULATION
COUNTRY GDP

-1   0   1   2   3   4   5   6

# What to Visualize for Neural Network Models?

For a neural network classifier, $p(y = 1 | w, x) = \sigma(\hat{f}(x))$ would it be helpful to print out all the weights?

# Weight Space Versus Function Space

While it's convenient to build up a complex function by composing simple ones  -as in neural networks- understanding the impact of each weight on the outcome is difficult.

In fact, the relationship between weights of a neural network and the function the network represents is extremely complicated:

1.  the same function may be represented by two very different set of weights for the same architecture.

2.  the architecture may be overly expressive - it can express the function $\hat{f}$ using a subset of the weights and hidden nodes (i.e. the trained model can have weights that are zero or nodes that contribute little to the computation).

# Debugging & Improving Models : Tensorboard

What happens if we want to know the outputs of a specific hidden layer?

By visualizing the network weights and activations as we train, we can diagnose issues that ultimately impact model performance.

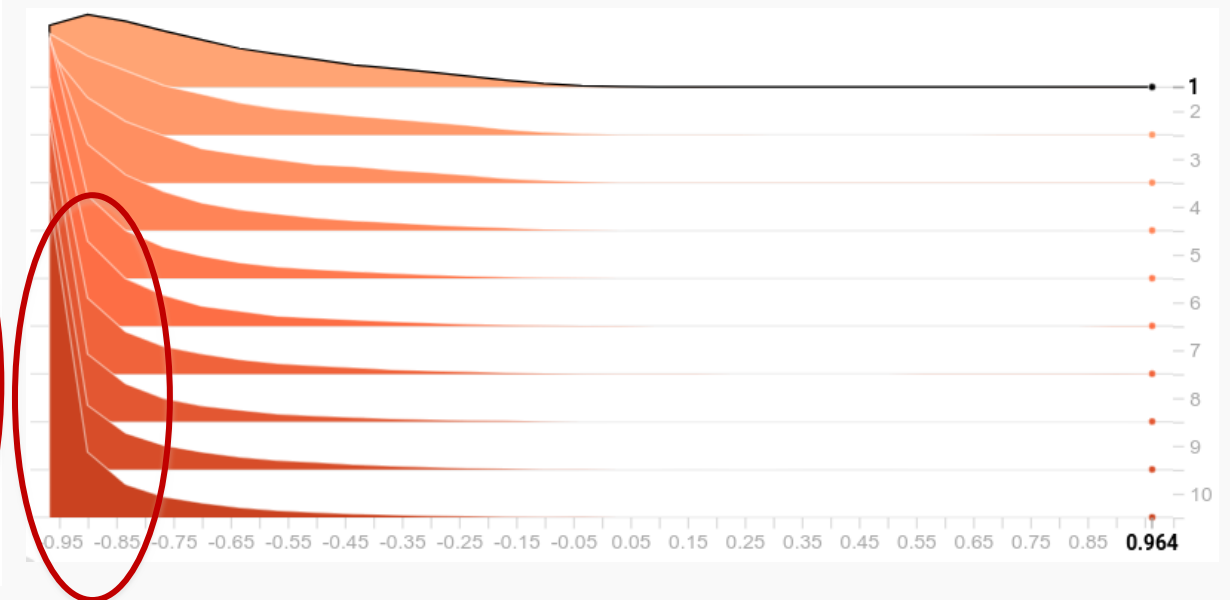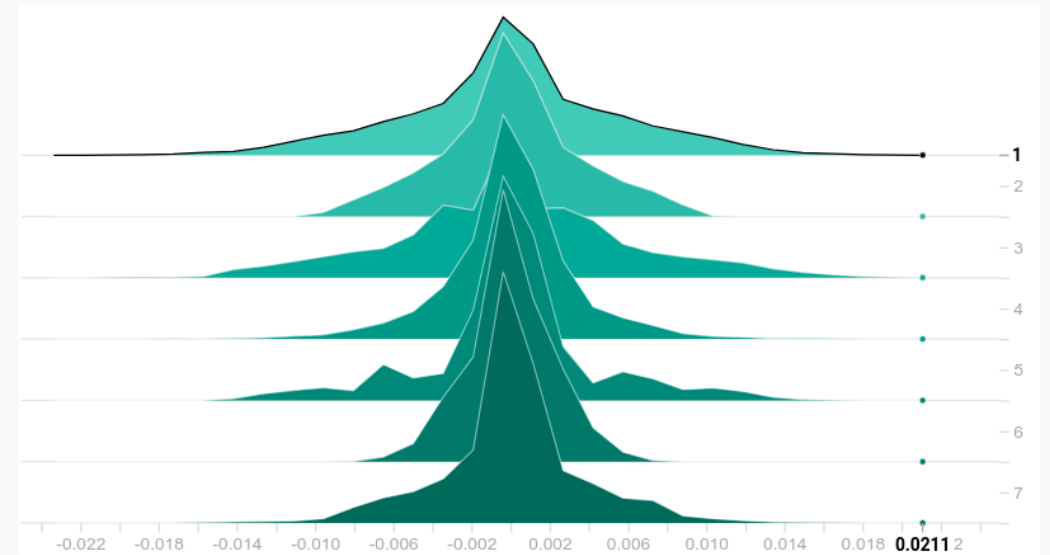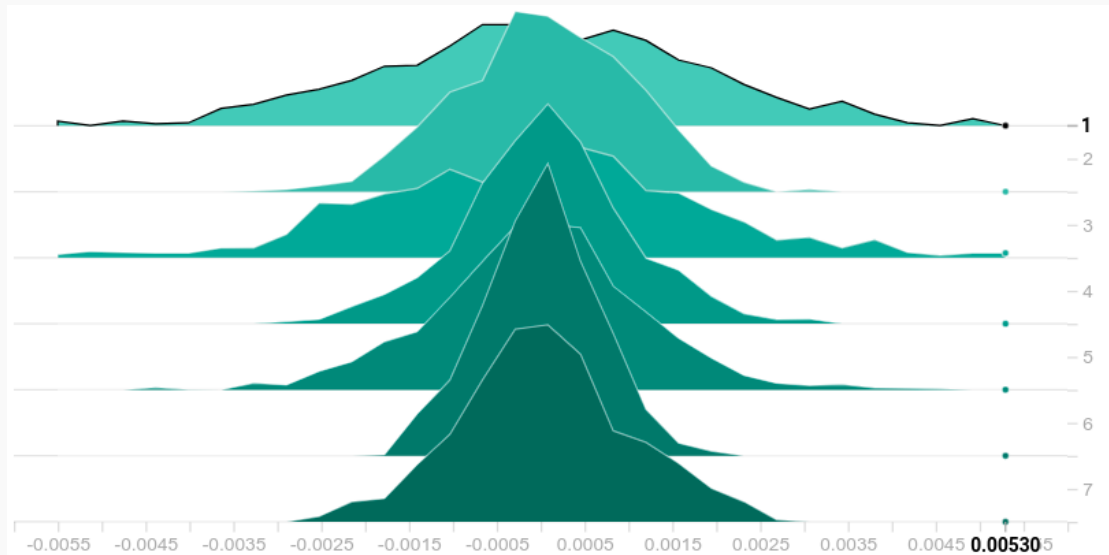TensorFlow provides a functionality to explore the inner workings of the network.

# Debugging & Improving Models : Tensorboard

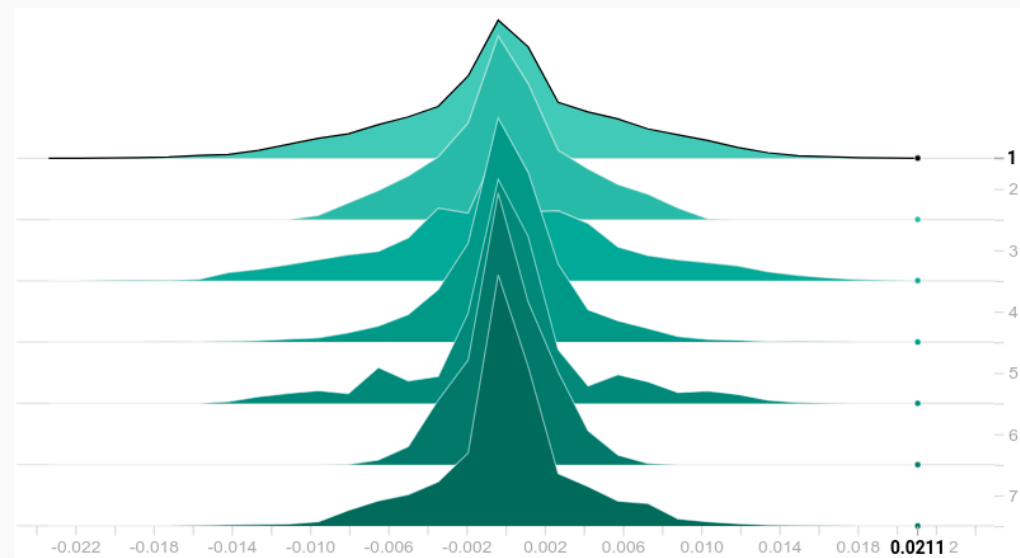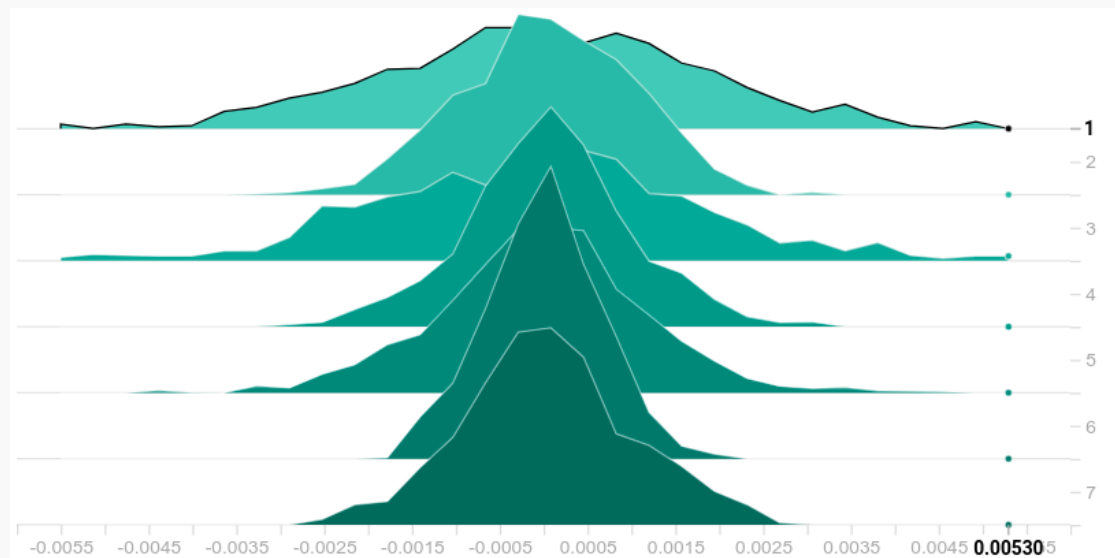The following visualizes the distribution of activations in two hidden layers over the course of training.

What problems do we see?



Tanh activation

Tanh activation

# Debugging & Improving Models : Tensorboard

The following visualizes the distribution of activations in two hidden layers over the course of training.

The activations are starting to saturate, reducing the learning speed.



Tanh activation

Sigmoid activation

# Debugging & Improving Models : Tensorboard

The following visualizes the distribution of gradients in two hidden layers over the course of training.

What problems do we see?

# Debugging & Improving Models : Tensorboard

The following visualizes the distribution of gradients in two hidden layers over the course of training.

In both layers, the gradients start to became zero. Smaller gradients imply smaller weight updates and slow training speed.

# CNN Feature Extraction Visualization

We know that CNNs extract features that best helps us to perform our downstream task (e.g. classification).

**Idea:** We train a CNN for feature extraction and a model (e.g. MLP, decision tree, logistic regression) for classification, *simultaneously* and *end-to-end*.

# CNN Feature Extraction Visualization

We know that CNNs extract features that best helps us to perform our downstream task (e.g. classification).

**Idea:** We train a CNN for feature extraction and a model (e.g. MLP, decision tree, logistic regression) for classification, *simultaneously* and *end-to-end*.

The resulting feature maps are matrices, that we can interpret as images. As such, we can analyze them a look for relevant patterns.
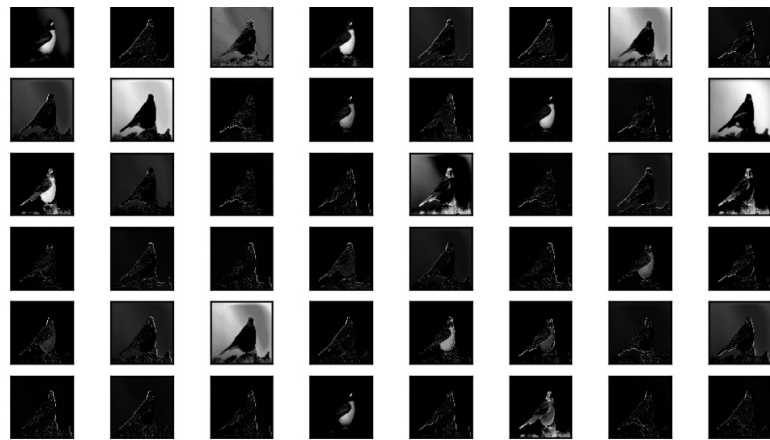


Convolution + Nonlinearity    Max Pooling

$p_{bird}$

$p_{chair}$

$p_{cat}$

$p_{dog}$

CONVOLUTION + POOLING LAYERS

FULLY CONNECTED LAYERS

52

# What to Visualize for CNNs?

The first things to try are:

1. visualize the result of applying a learned filter to an image
2. visualize the filters themselves

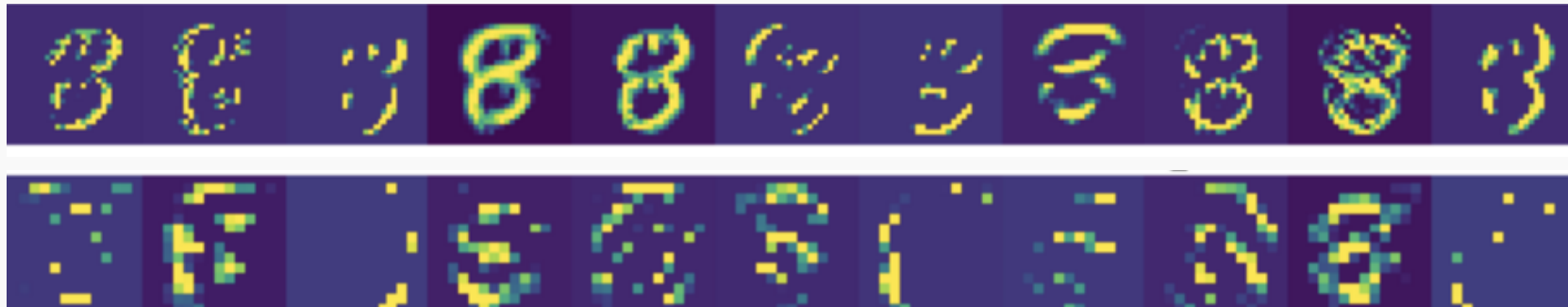**HARD TO UNDERSTAND OR INTERPRETATE**

Input image

Feature maps

Kernels

# Occlusion methods

If we want to interpret what part of the image the network is paying more attention, these visualizations might not be the best solution.



Activations maps for layer1

Activations maps for layer2

We have no guarantees that the feature maps will provide meaningful information. Their interpretation can be even more difficult than the original problem.
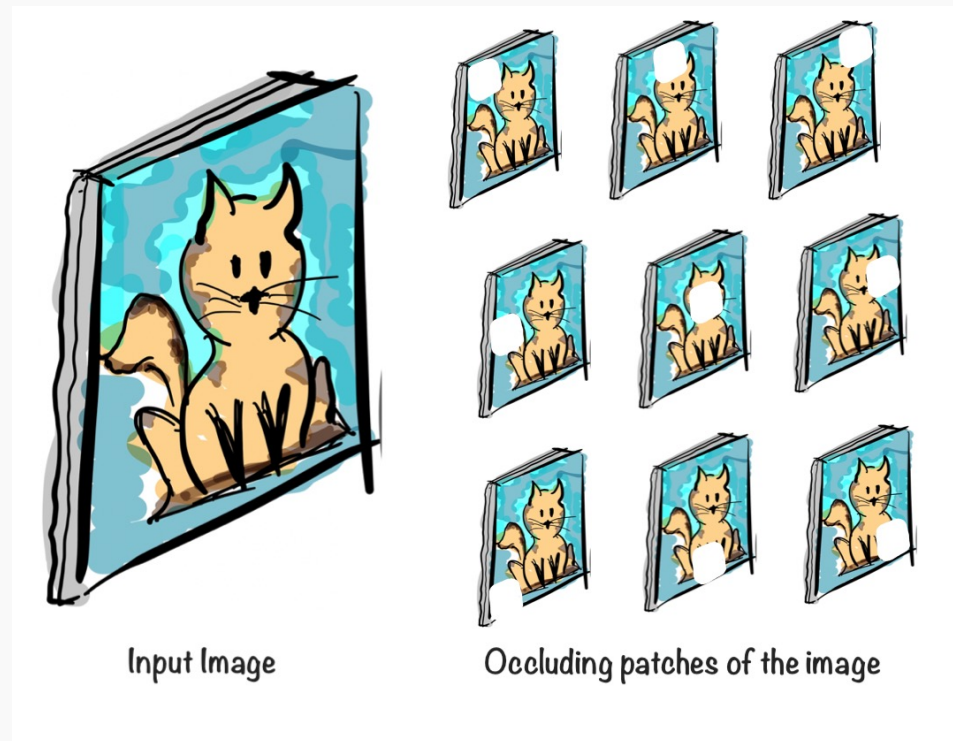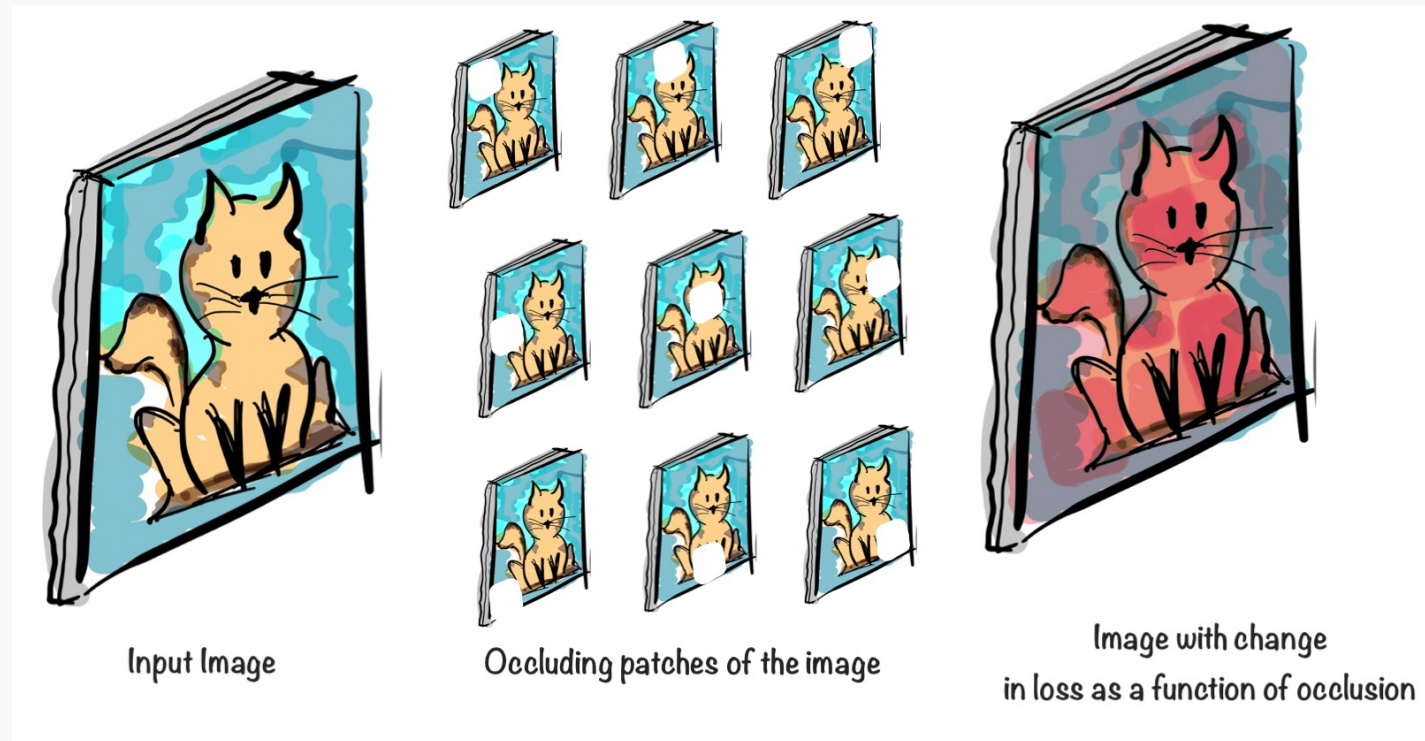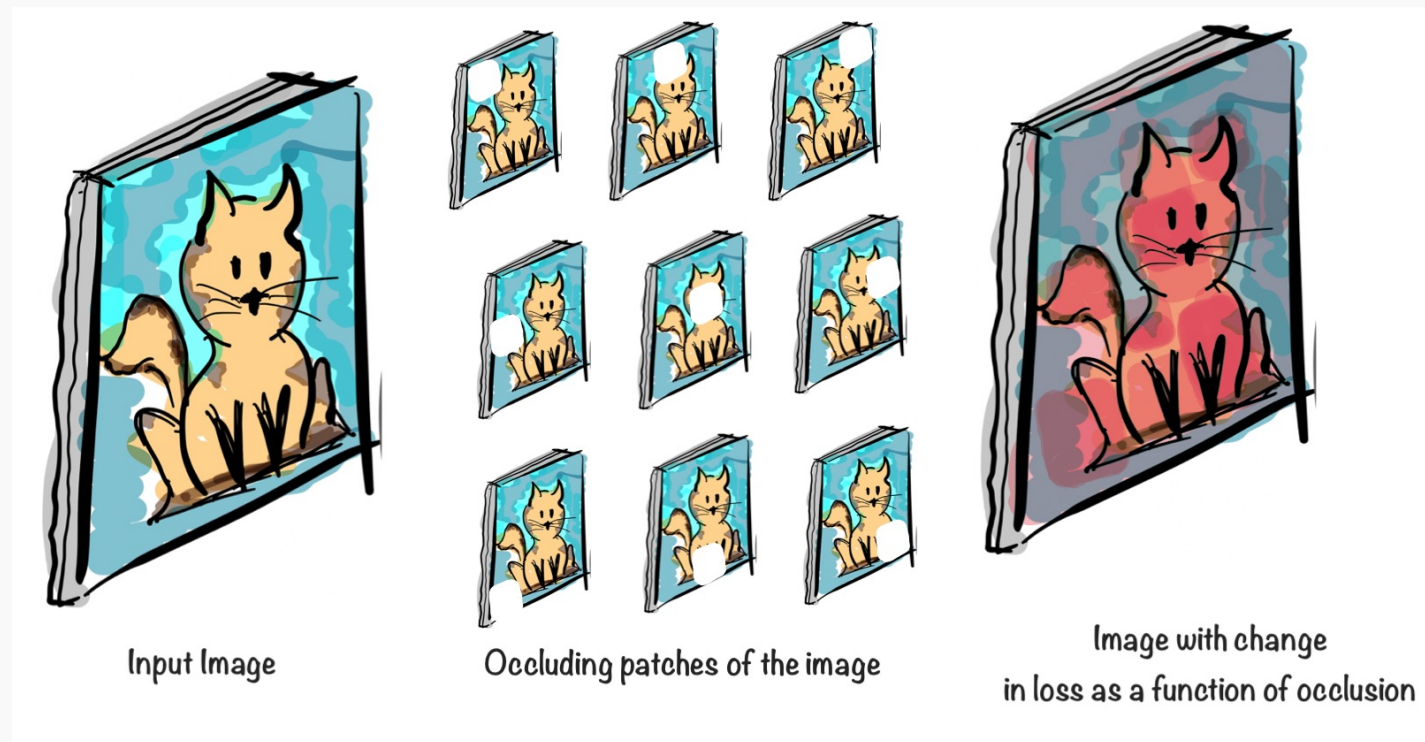
# Occlusion methods

Occlusion methods attributes importance for the classification of the image. Occlusion involves running a patch over part of the image to see which pixels affect the classification the most.



Input Image

# Occlusion methods

Occlusion methods attributes importance for the classification of the image. Occlusion involves running a patch over part of the image to see which pixels affect the classification the most.



Input Image

Occluding patches of the image

# Occlusion methods

Occlusion methods attributes importance for the classification of the image. Occlusion involves running a patch over part of the image to see which pixels affect the classification the most.



Input Image

Occluding patches of the image
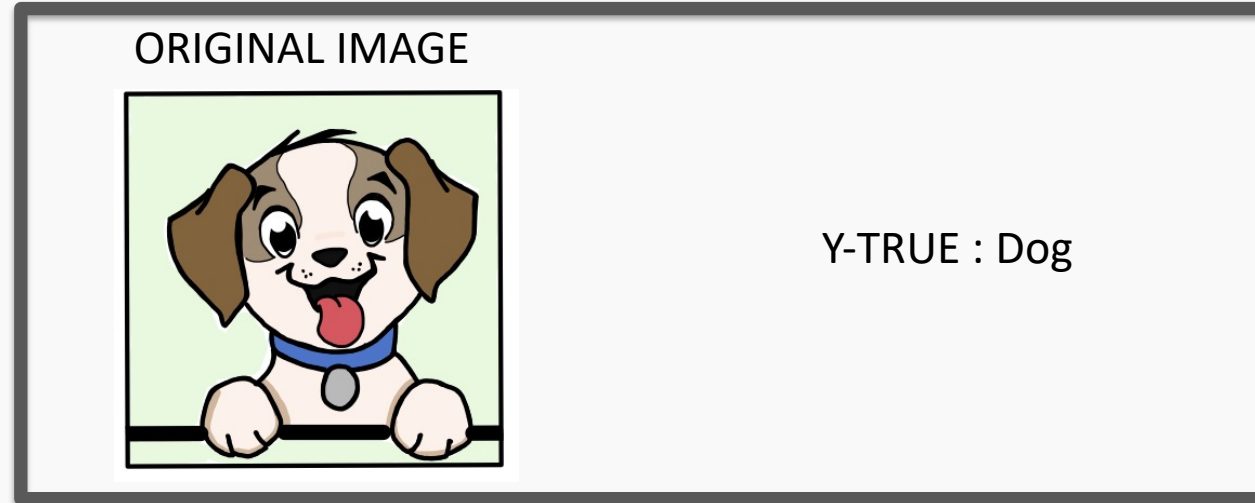
Image with change
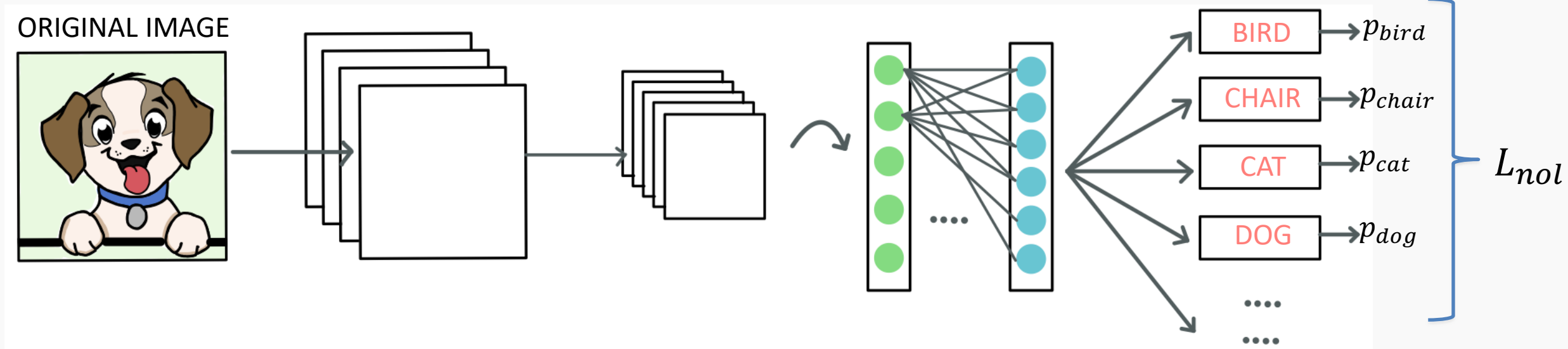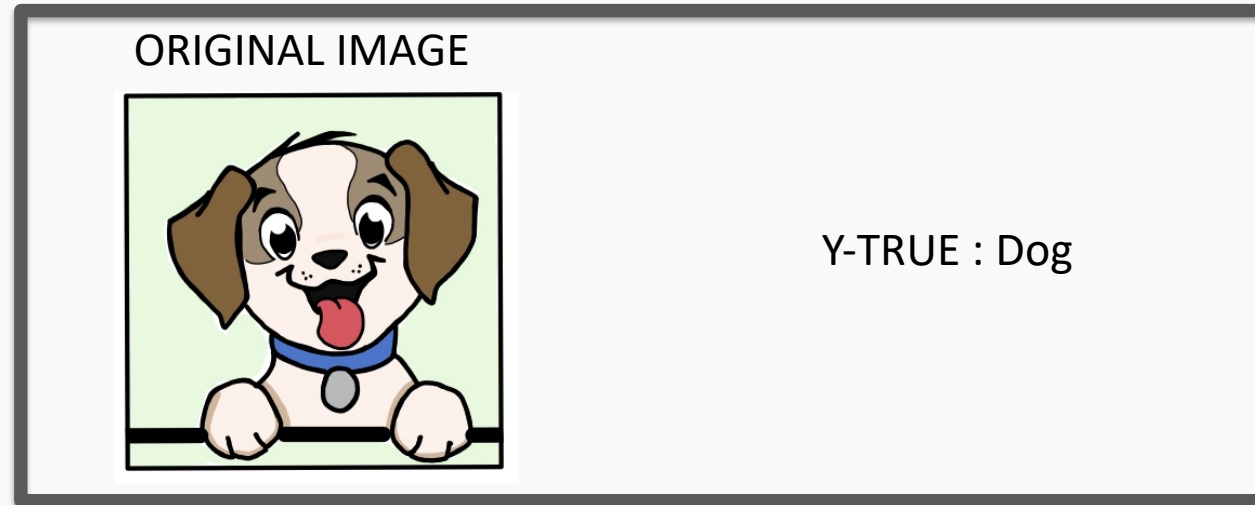in loss as a function of occlusion

# Occlusion methods

However, to obtain fine details we need to use a small occlusion area, increasing the number of model evaluations. This can become impractical for a fine resolution and many test images.
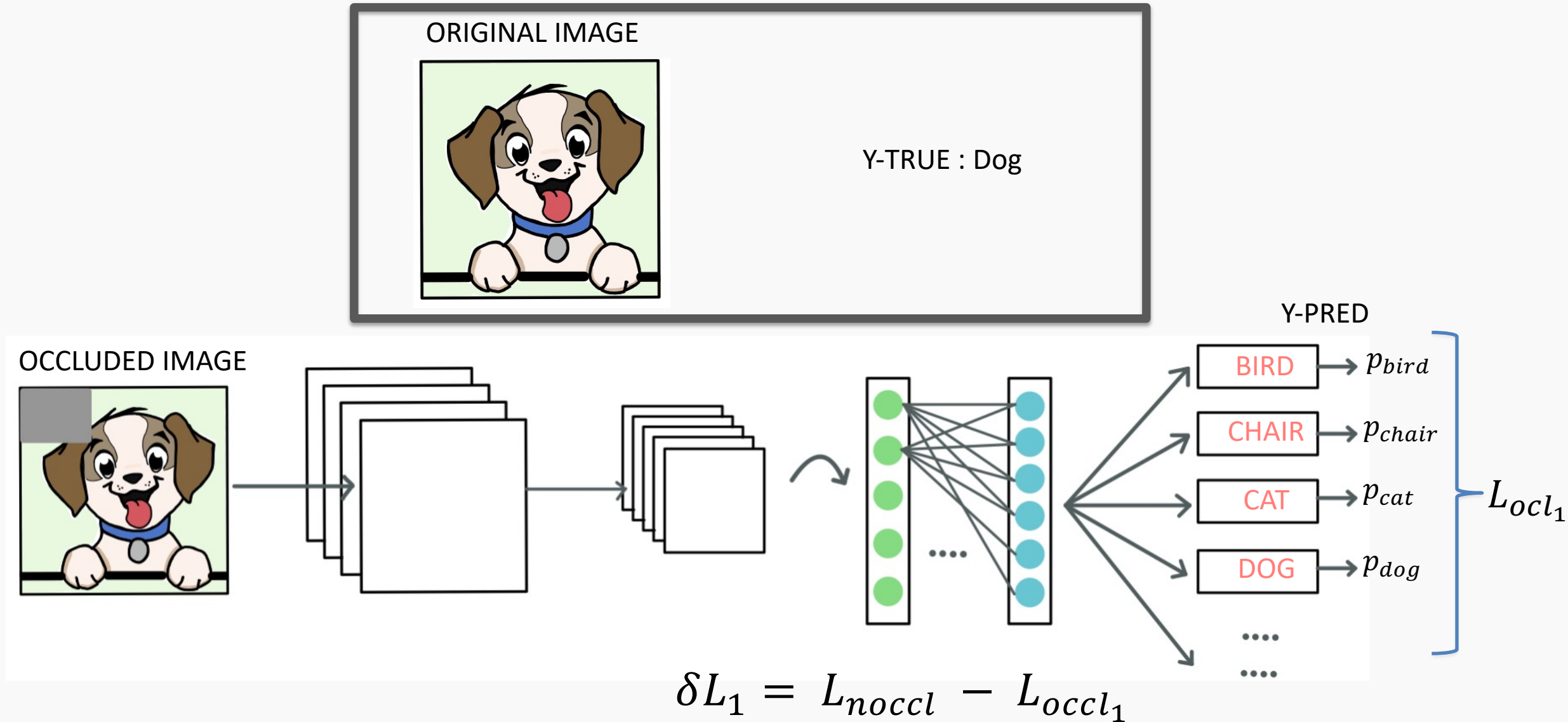


Input Image

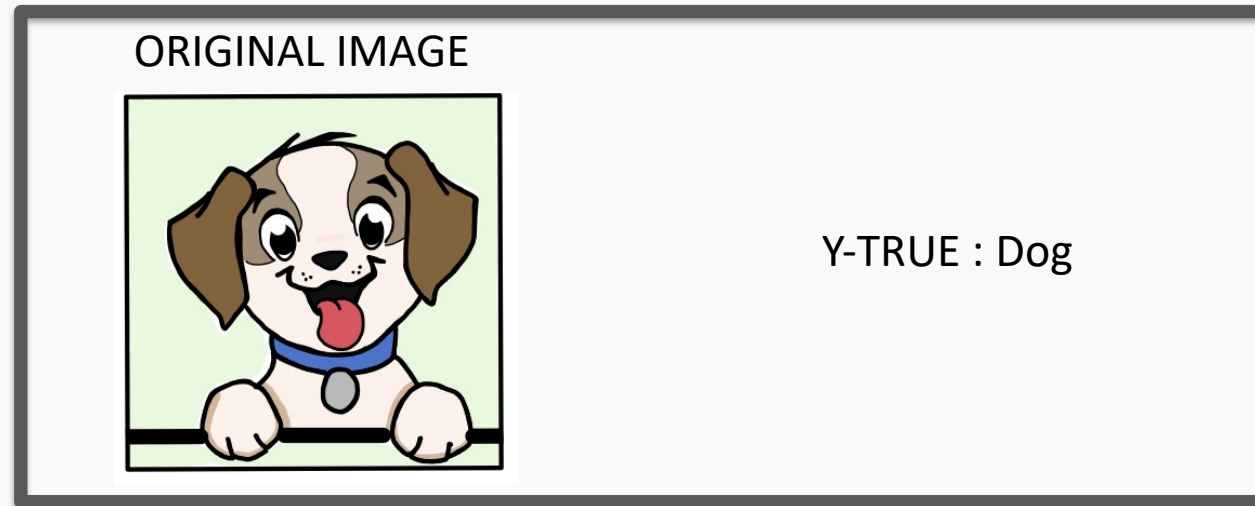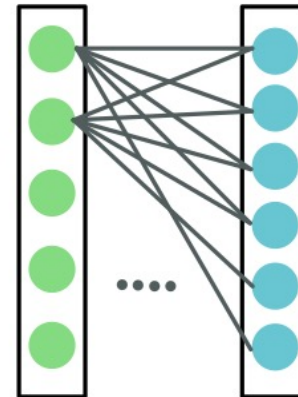Occluding patches of the image

Image with change
in loss as a function of occlusion

# Occlusion process

ORIGINAL IMAGE



Y-TRUE : Dog

# Occlusion process



ORIGINAL IMAGE

Y-TRUE : Dog

Y-PRED

ORIGINAL IMAGE

BIRD $\rightarrow p_{bird}$

CHAIR $\rightarrow p_{chair}$

CAT $\rightarrow p_{cat}$

DOG $\rightarrow p_{dog}$

$L_{nol}$

# Occlusion process



ORIGINAL IMAGE

Y-TRUE : Dog

OCCLUDED IMAGE

Y-PRED

BIRD $\rightarrow p_{bird}$

CHAIR $\rightarrow p_{chair}$

CAT $\rightarrow p_{cat}$

DOG $\rightarrow p_{dog}$

$L_{ocl_1}$

$$\delta L_1 = L_{noccl} - L_{occl_1}$$

# Occlusion process



ORIGINAL IMAGE

Y-TRUE : Dog

OCCLUDED IMAGE

Y-PRED

BIRD $\rightarrow$ $p_{bird}$

CHAIR $\rightarrow$ $p_{chair}$

CAT $\rightarrow$ $p_{cat}$

DOG $\rightarrow$ $p_{dog}$

....

....

$L_{ocl_2}$

$$\delta L_2 = L_{noccl} - L_{occl_2}$$

# Occlusion process



ORIGINAL IMAGE

Y-TRUE : Dog

Y-PRED

OCCLUDED IMAGE

BIRD $\rightarrow p_{bird}$

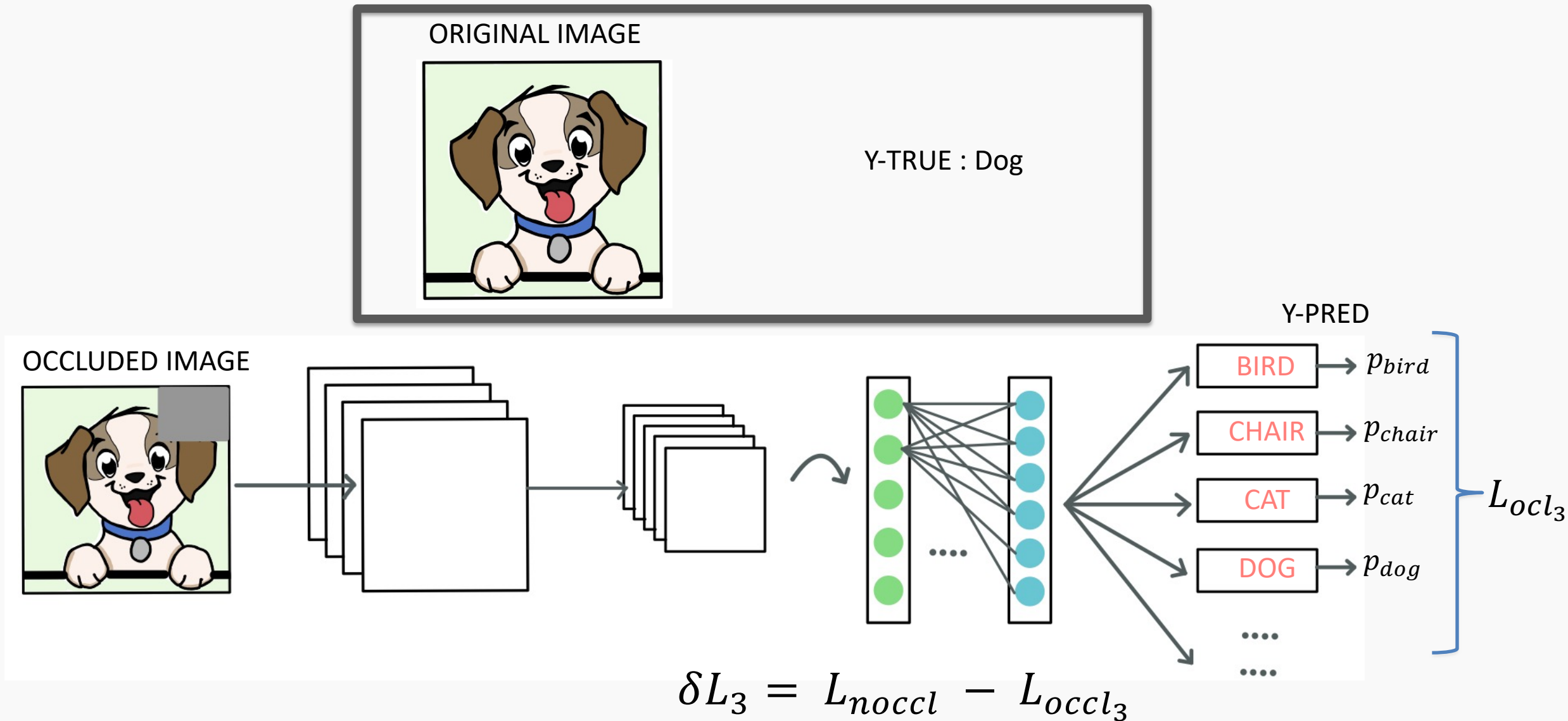CHAIR $\rightarrow p_{chair}$

CAT $\rightarrow p_{cat}$

DOG $\rightarrow p_{dog}$

$L_{ocl_3}$

$$\delta L_3 = L_{noccl} - L_{occl_3}$$
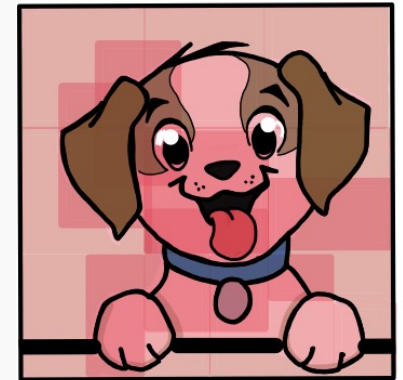
# Occlusion process

ORIGINAL IMAGE



Occlusion →

$$\begin{bmatrix} \delta L_1 & \delta L_2 & \delta L_3 & \delta L_4 \\ \delta L_5 & \delta L_6 & \delta L_7 & \delta L_8 \\ .. & .. & .. & .. \\ \delta L_{i-3} & \delta L_{i-2} & \delta L_{i-1} & \delta L_i \end{bmatrix}$$
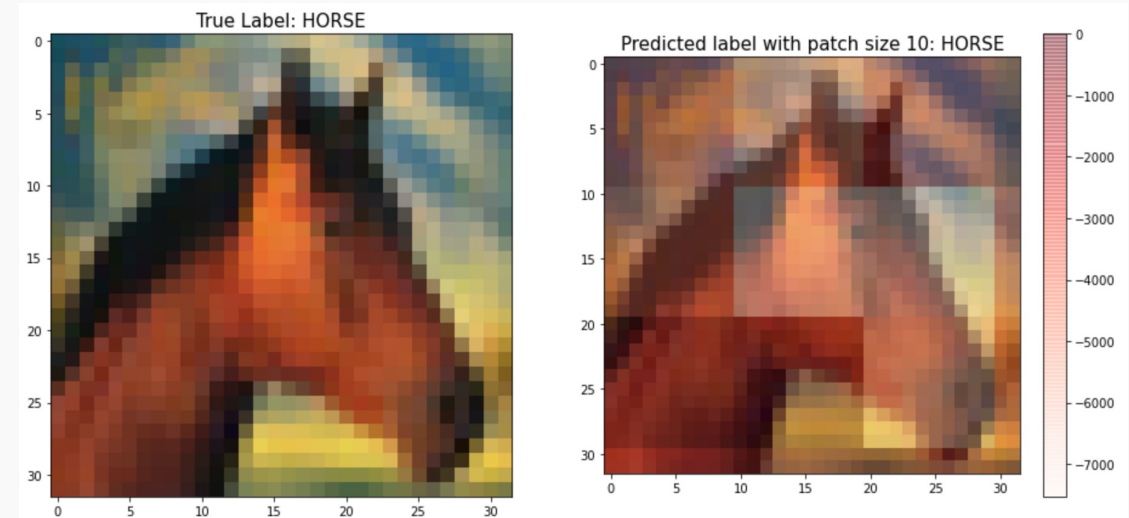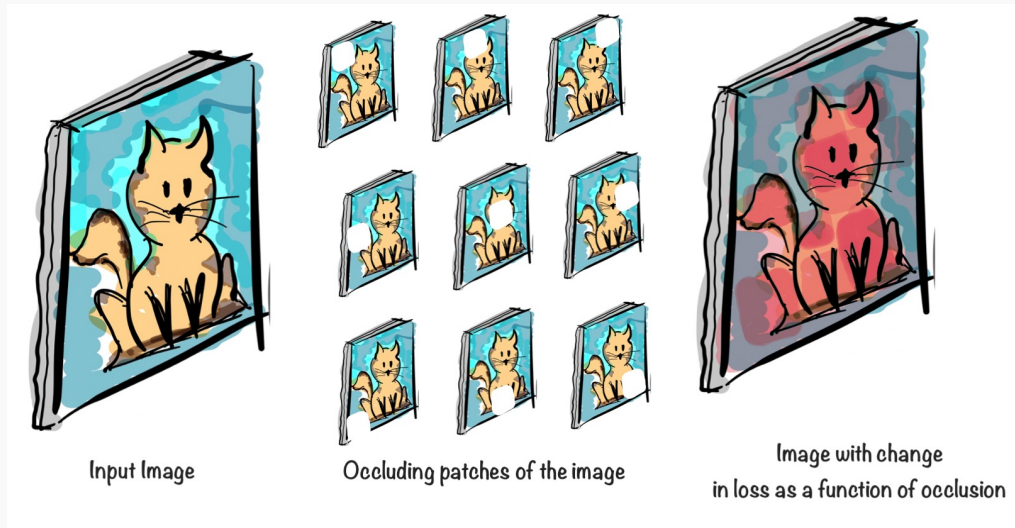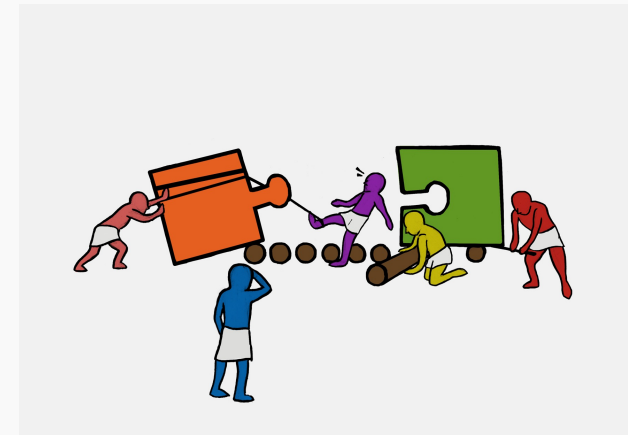
OCCLUSION LOSS MAP

# Occlusion process

- Input image to a trained network
- Take note of the true label, K
- Get the prediction of the true image, Q
- Compute the loss, $L_{nol} = -\log P(y = Q)$
- Occlude patches of the image with gray blocks starting at the top left
  - Get the prediction of this occluded version of the image
  - Compute the loss, $L_{ocl_i} = -\log P_{occ}(y = Q)$
  - Compute the difference of the losses, $L_{nol} - L_{ocl_i}$

If K=Q, we answer the what parts of the image have contributed to correctly predict. If K<>Q, we answer what parts of the image contributed to predict the incorrect class.

# Exercise: Image Occlusion

The aim of this exercise is to understand occlusion. Occlusion involves running a patch over the entire image to see which pixels affect the classification the most.



Input Image

Occluding patches of the image

Image with change in loss as a function of occlusion



True Label: HORSE

Predicted label with patch size 10: HORSE

# Taylor series expansion

Any differentiable function $f(x)$ can be approximated as a series around $x_0$ as:

$$f(x) = f(x_0) + \frac{(x - x_0)^1}{1!} \left.\frac{\partial f}{\partial x}\right|_{x_o} + \frac{(x - x_0)^2}{2!} \left.\frac{\partial^2 f}{\partial x^2}\right|_{x_o} + \cdots$$

This function can be the logistic regression or even a complex neural network.

**Note**: Including more terms will improve the approximation.

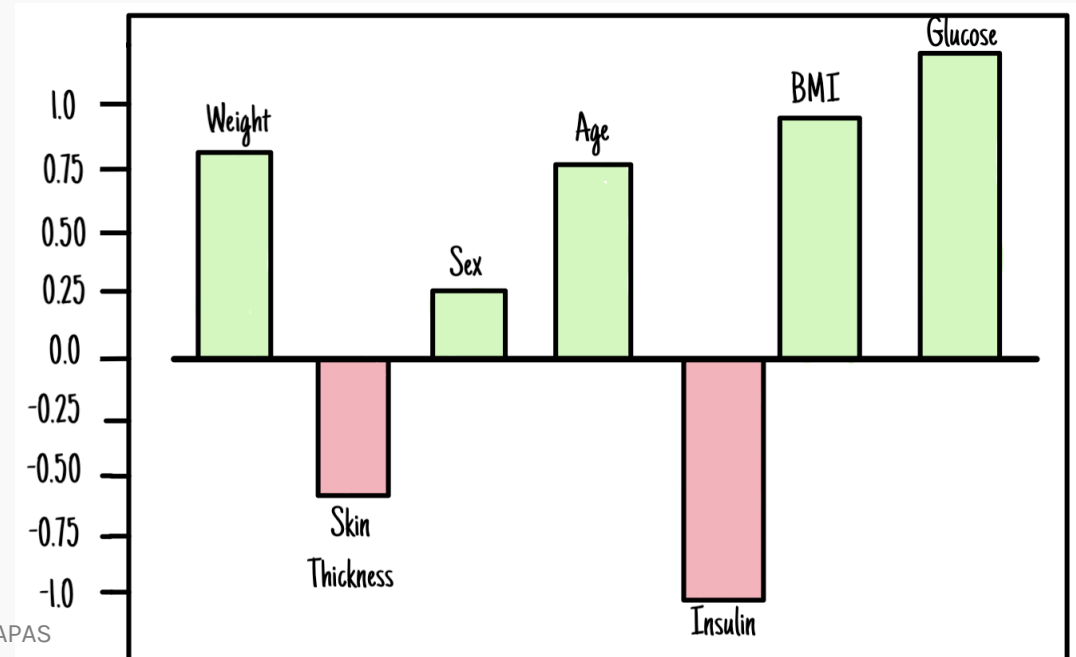# Visualizing Top Predictors by Input Gradient

Since the input gradient of an objective function for a trained model indicates which input dimensions has the greatest effect on the model decision at an input **x**, we can visualize the "top predictors" of outcome for a particular input **x**.

We can think of this as approximating our neural network model with a linear model locally at an input **x** and then interpreting the weights of this linear approximation.

$$NN(x) \approx NN(x_0) + w^T(x - x_0)$$
$$\approx NN(0) + w^T(x)$$
$$\approx w^T x + b$$

$$w = \frac{\partial NN}{\partial x}\bigg|_{x_o}$$

58

# Thank you