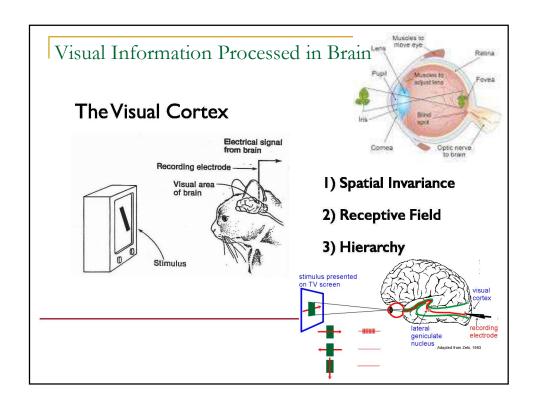
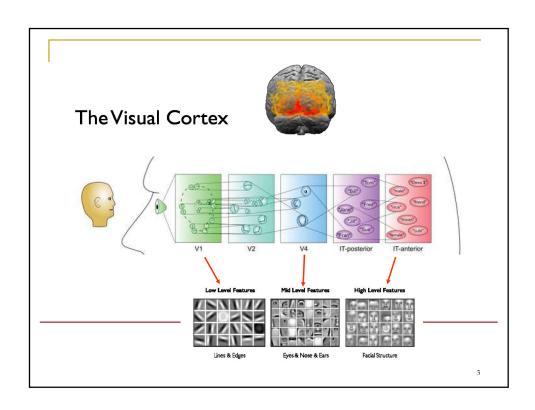
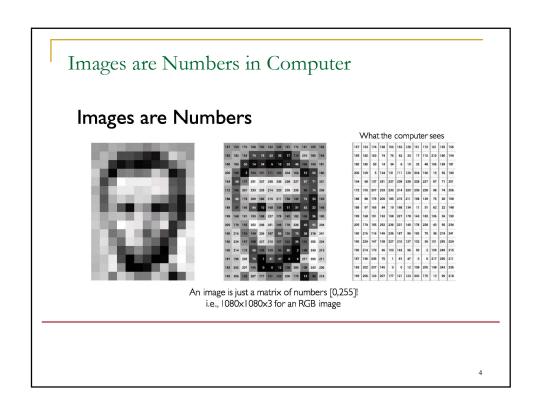
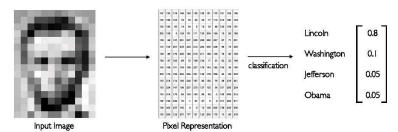
# Convolutional Neural Networks







#### Tasks in Computer Vision



- Regression: output variable takes continuous value
- Classification: output variable takes class label. Can produce probability of belonging to a particular class

#### High Level Feature Detection

#### High Level Feature Detection

Let's identify key features in each image category



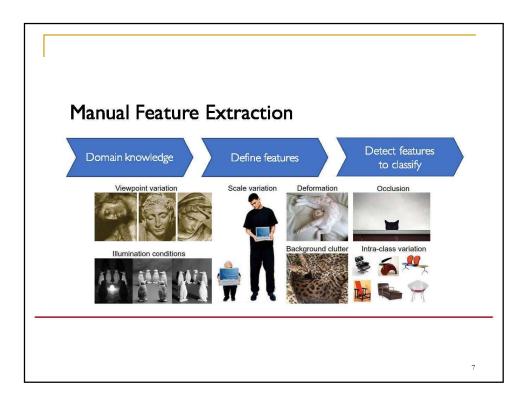
Nose, Eyes, Mouth



Wheels, License Plate, Headlights

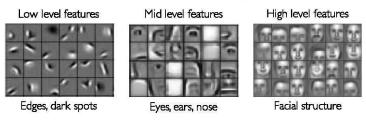


Door, Windows, Steps



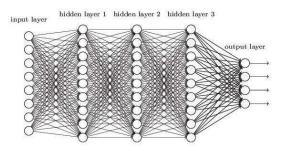
#### Learning Feature Representations

Can we **learn hierarchy of features** directly from the data instead of hand engineering?



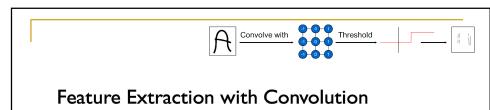
# Learning Visual Features

#### Fully Connected Neural Network



- Connecting neuron in hidden layer to all neurons in input layer
- No spatial information
- Too many parameters

9



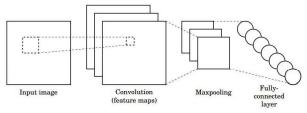
- Filter of size 4x4 : 16 different weights
- Apply this same filter to 4x4 patches in input
- Shift by 2 pixels for next patch

This "patchy" operation is convolution

- 1) Apply a set of weights a filter to extract local features
  - 2) Use multiple filters to extract different features
    - 3) Spatially share parameters of each filter

#### Convolutional Neural Network

#### CNNs for Classification



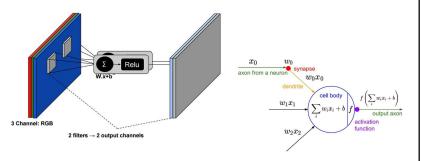
- Convolution: Apply filters with learned weights to generate feature maps.
   Non-linearity: Often ReLU.
- 3. Pooling: Downsampling operation on each feature map.

Train model with image data. Learn weights of filters in convolutional layers.

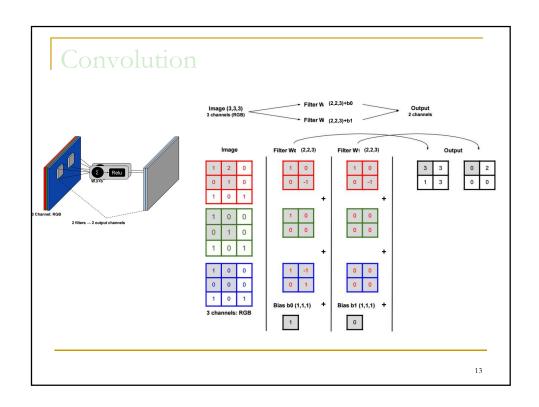
11

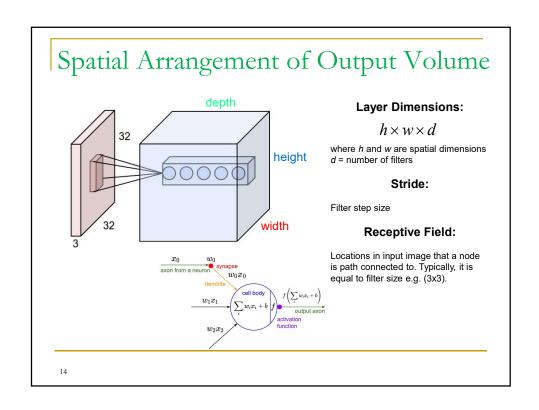
#### Convolution

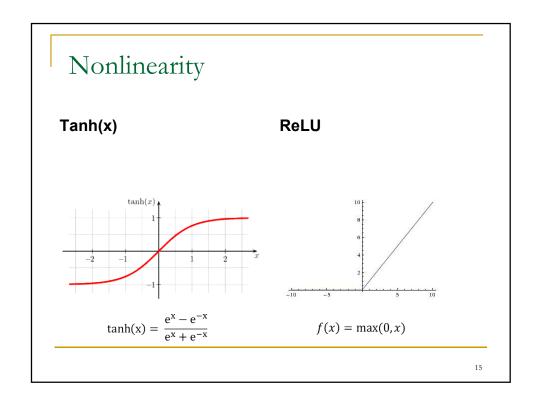
A convolution is a neighborhood operation in which each output pixel is the weighted sum of neighboring input pixels. The matrix of weights is called the convolution kernel, also known as a filter.

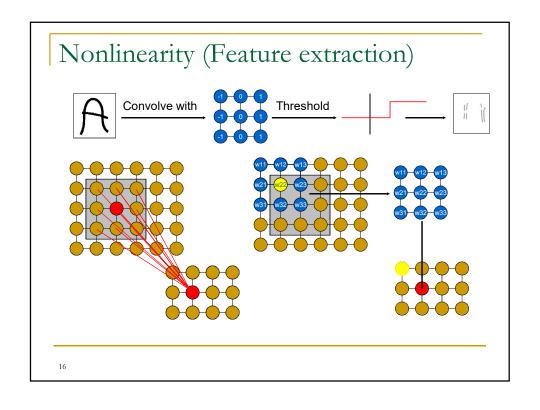


Actually, each filter corresponds to one perceptron. The perceptron is moved around the whole image to extract local features.



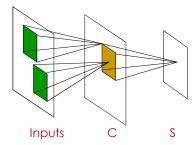






### Nonlinearity (feature extraction)

- Shared weights: all neurons in a feature share the same weights (but not the biases).
- In this way all neurons detect the same feature at different positions in the input image.
- Reduce the number of free parameters.



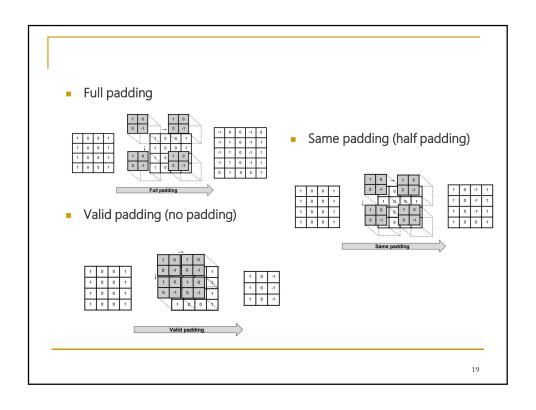
17

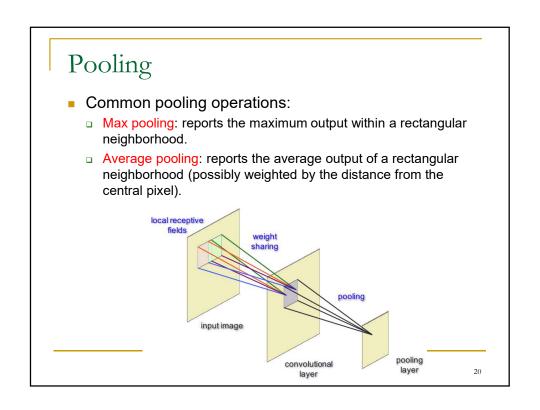
# **Padding**

Padding is basically adding rows or columns of zeros to the borders of an image input. It helps control the output size of the convolution layer. The formula to calculate the output size is: (N – F) / stride + 1.

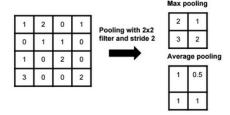
For a 32x32x3 image and using  $10.5 \times 5$  filters with stride 1 and pad 2, we get an output with size 32x32x10.

| F . |   |   |   |   |   |
|-----|---|---|---|---|---|
| 0   | 0 | 0 | 0 | 0 | 0 |
| 0   | 0 | 1 | 1 | 3 | 0 |
| 0   | 1 | 2 | 3 | 5 | 0 |
| 0   | 2 | 3 | 5 | 1 | 0 |
| 0   | 0 | 1 | 1 | 1 | 0 |
| 0   | 0 | 0 | 0 | 0 | 0 |





#### An Example

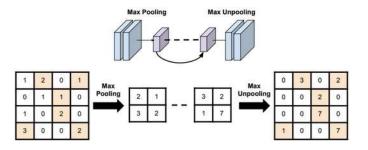


With pooling we reduce the size of the data without changing the depth. Max pooling preserves edge, resulting in spatial invariance.

The output size of a polling operation on a 8x8x10 representation using a  $2 \times 2$  filter and with stride 2 is 4x4x10 (We can use the same formula: (N - F) / stride + 1).

21

# Unpooling

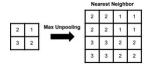




Bed of nails unpooling



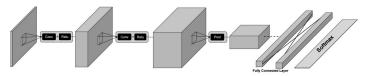
Nearest neighbor unpooling



23

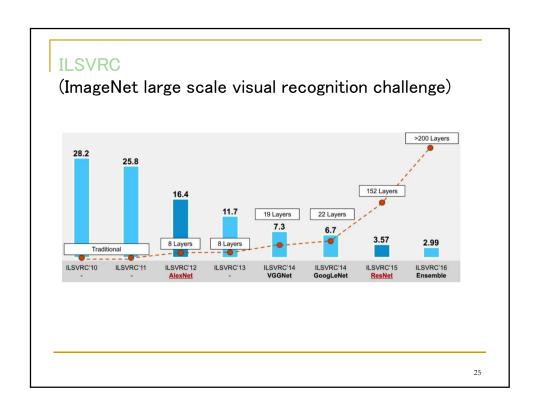
## Architecture

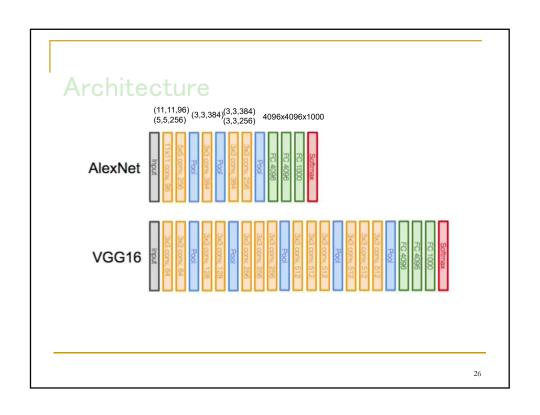
In general the architecture of a convolutional neural network is as below:

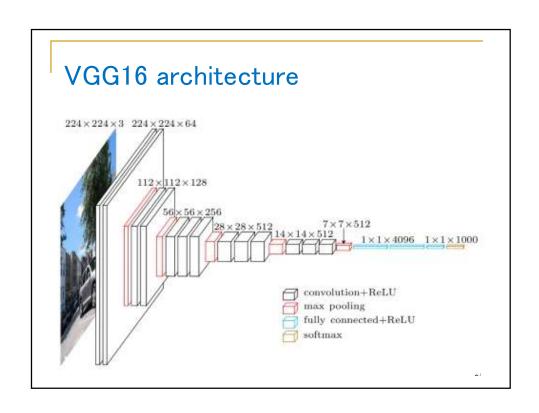


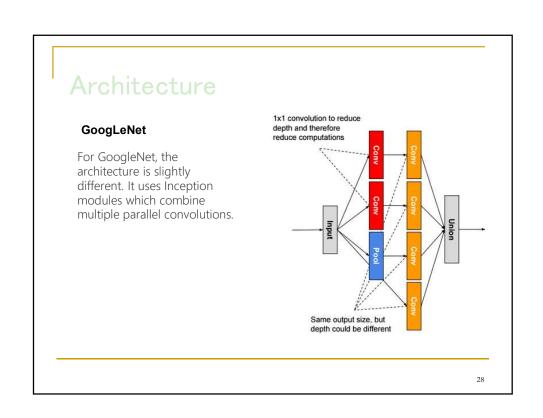
Conv $\rightarrow$  Relu $\rightarrow$  Conv  $\rightarrow$  Relu $\rightarrow$  Pool  $\rightarrow$  ...  $\rightarrow$  Conv  $\rightarrow$  Fully Connected Layer  $\rightarrow$  Softmax

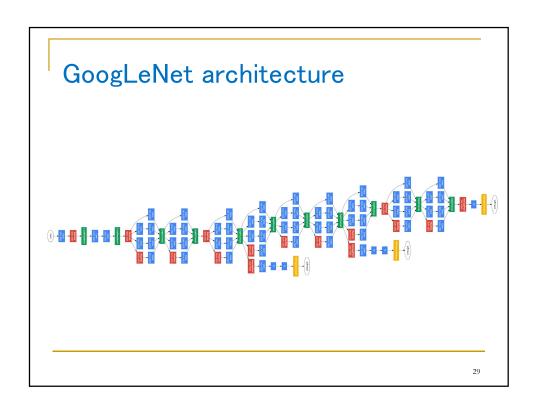
Some well known CNN architectures are: AlexNet (8 layers), VGG (16-19 layers), GoogLeNet (22 layers) and ResNet (152 layers).

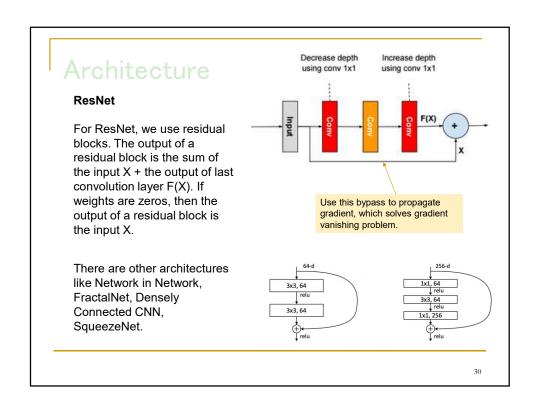


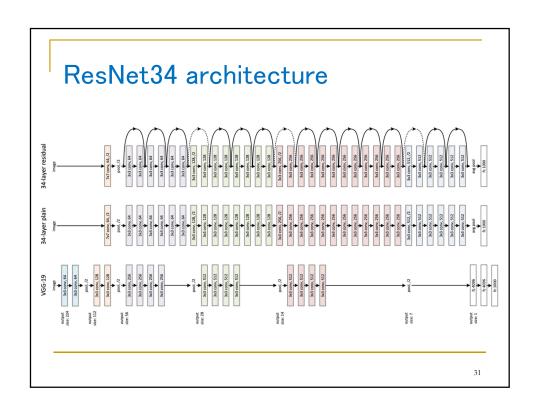


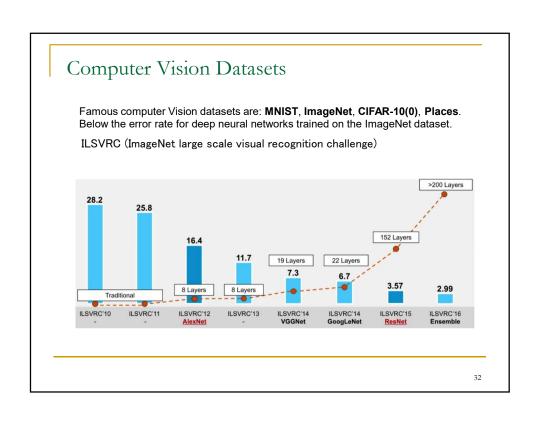


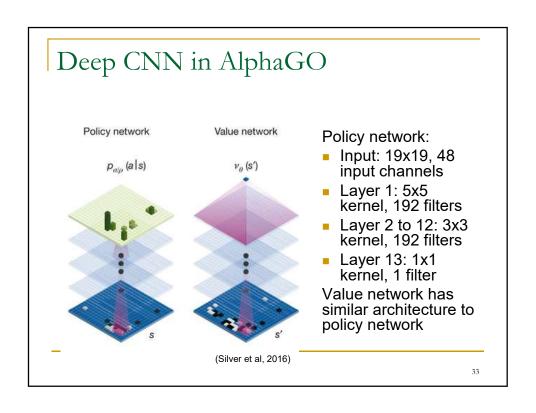


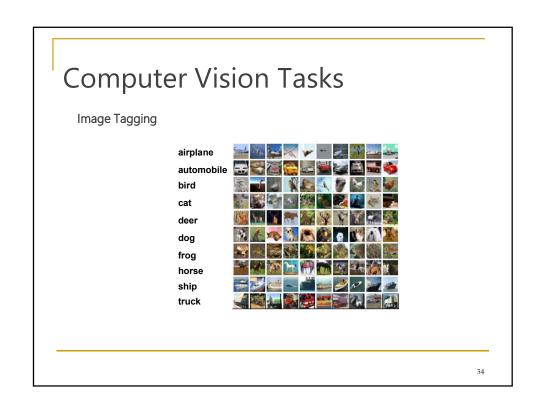




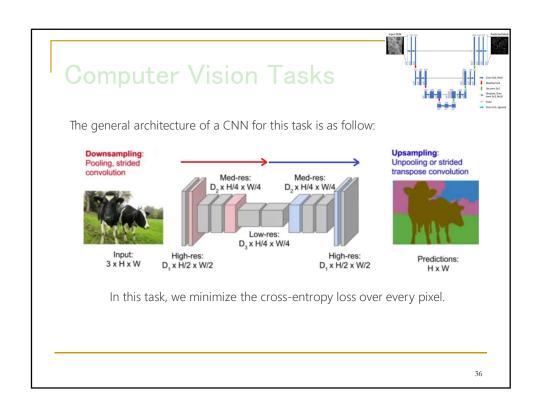


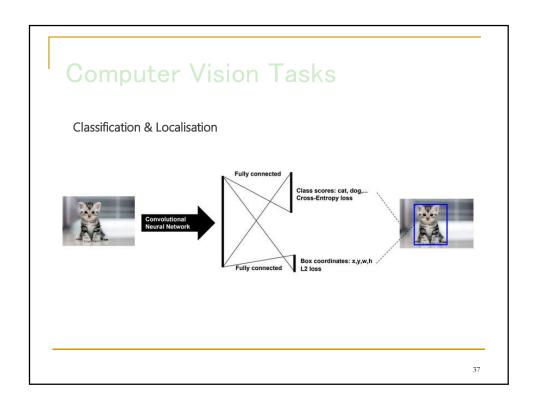


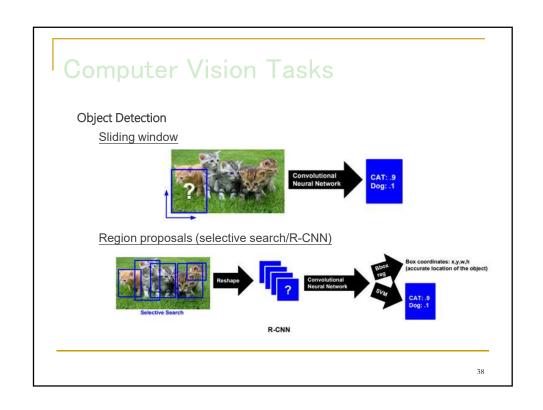


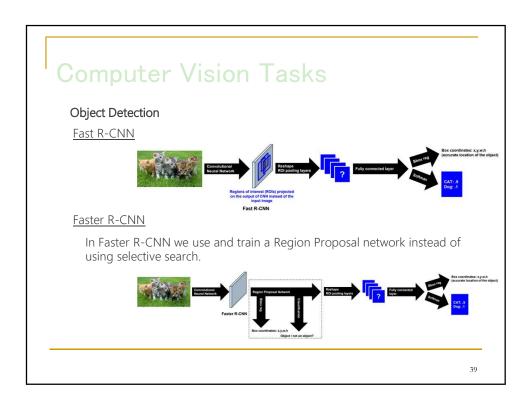


# Computer Vision Tasks Semantic Segmentation Semantic segmentation is the task of assigning a class-label to each pixel in an image. Sky Cat Cow Sky Grass





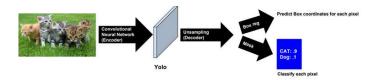






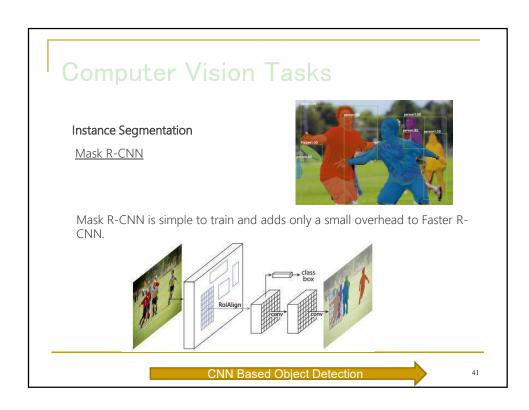
#### Object Detection

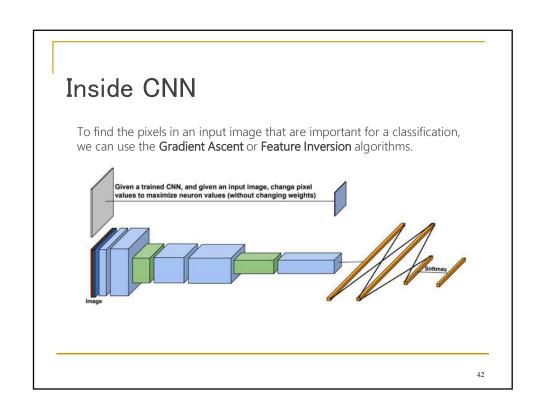
Yolo



It's recommended to use **Focal loss** function when training the model. **Focal loss** is a loss function considered class imbalance defined by

$$FL(p_t) = -\alpha_t \frac{(1 - p_t)^{\gamma}}{\log(p_t)}$$





#### Adversarial Attacks on Neural Networks

#### Adversarial Image:

Modify image to increase error

$$x \leftarrow x + \eta \frac{\partial J(\overrightarrow{w}, x, \overrightarrow{y})}{\partial x}$$
 Fix your weights  $\theta$ , and true label  $y$ 

"How does a small change in the input increase our loss"