

15-494/694: Cognitive Robotics

Dave Touretzky

Lecture xx:
ImageNet and Transfer
Learning

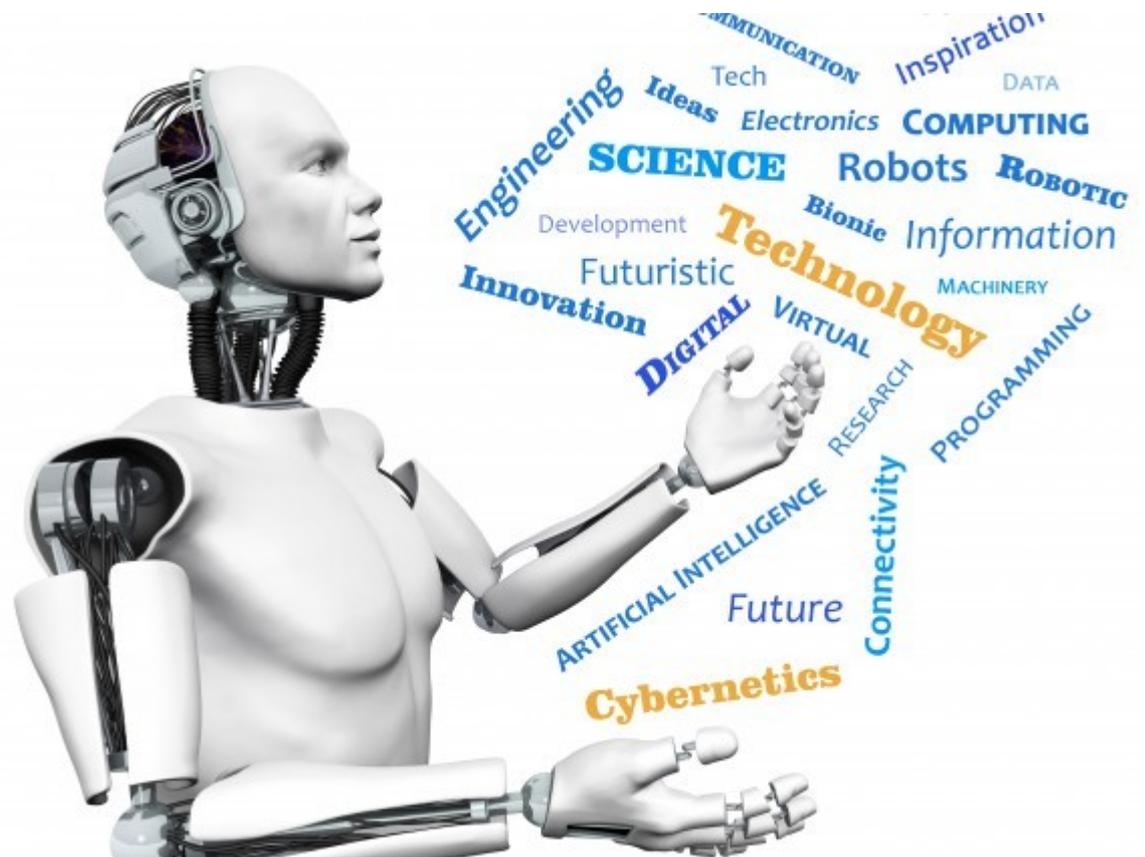


Image from <http://www.futuristgerd.com/2015/09/10>

Object Recognition Challenge

- Computer vision researchers use challenge events to measure progress in the state of the art.
- PASCAL VOC (Visual Object Classes) Challenge:
 - Ran from 2005 to 2012
 - 2005 version had 4 categories (bicycles, motorcycles, people, cars) and 1578 training images
 - 2012 version had 20 categories and 5717 training images

ImageNet

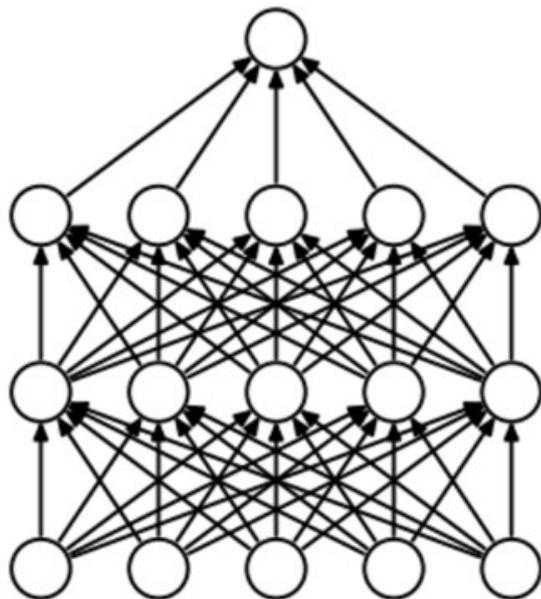
- Created by Fei-Fei Li at Stanford.
- See www.image-net.org
- 15 million labeled images, 22,000 categories
- ILSVRC: ImageNet Large Scale Visual Recognition Challenge: 2009-2017
 - 1000 categories, including 118 dog breeds
 - 1.2 million training images

AlexNet

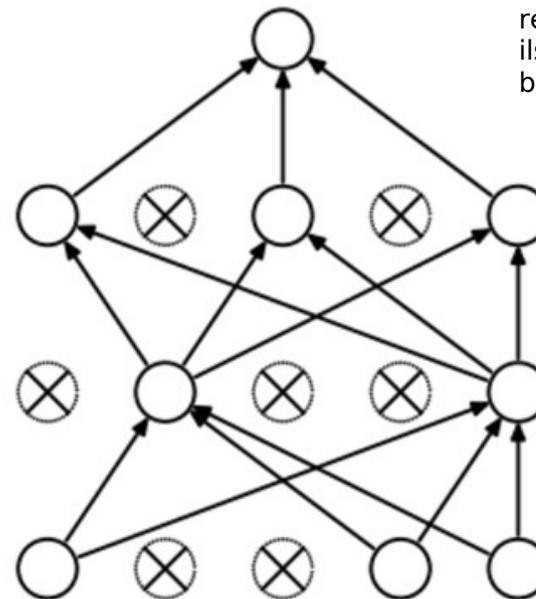
- The winners of the 2012 ILSVRC:
 - Alex Krizhevsky, Ilya Sutsker, and Geoffren Hinton
 - Deep convolutional neural net (DCNN) called AlexNet
 - Trained using two GPU boards
 - Introduced ReLU in place of tanh
 - Used “dropout” to reduce overfitting
 - Error rate of 15.3% was 10% better than the runner-up
 - Put deep neural nets on the map

Dropout in AlexNet

- For each training step, disable 50% of the neurons for both the forward and backward pass.
- Reduces overfitting.



(a) Standard Neural Net



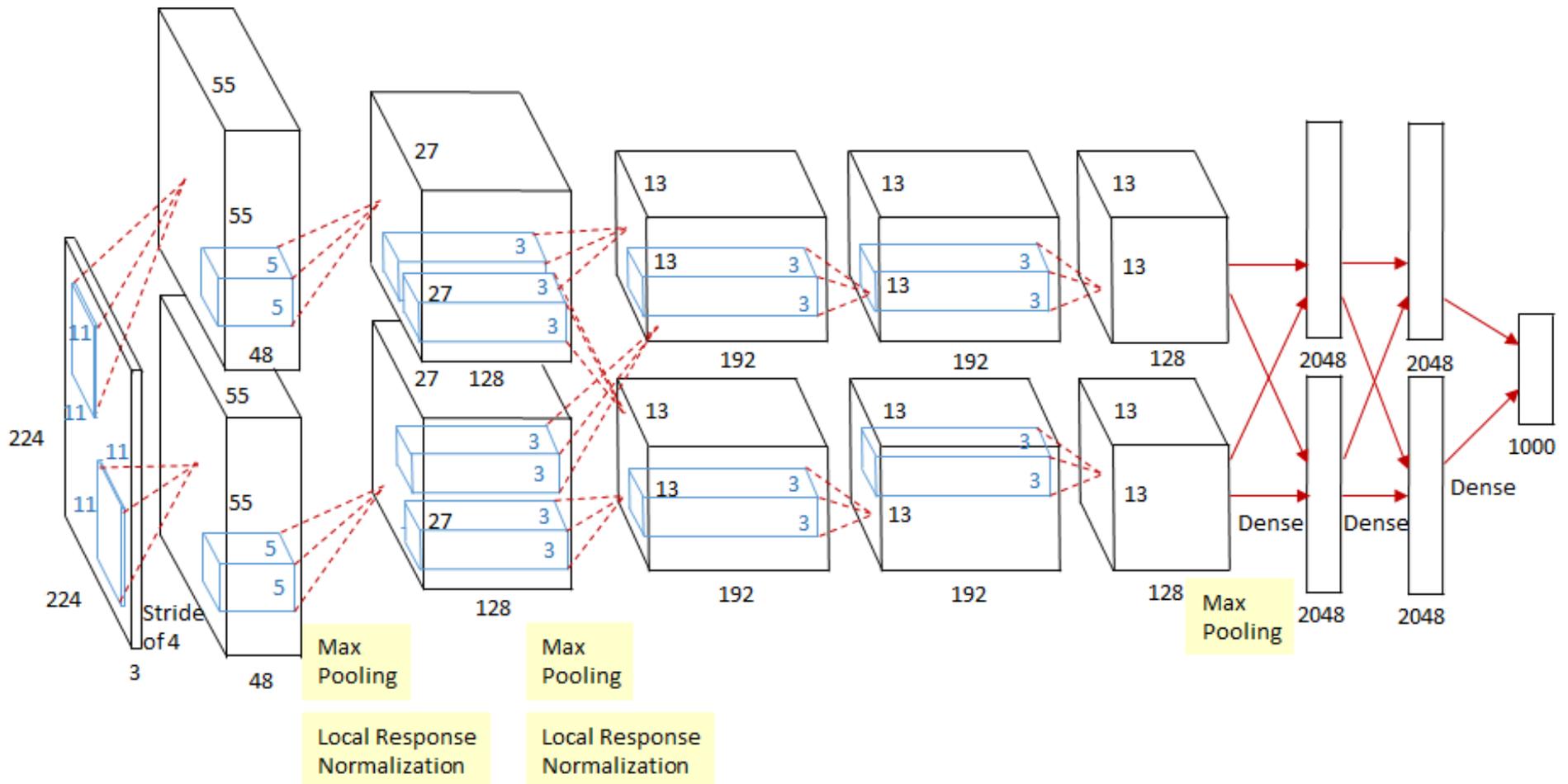
(b) After applying dropout.

Figure from
<https://medium.com/coinmonks/paper-review-of-alexnet-caffenet-winner-in-ilsvrc-2012-image-classification-b93598314160>

Data Augmentation in AlexNet

- Take random 224×224 crops of a 256×256 image, plus their horizontal reflections. Increases training set size by a factor of $32^2 \times 2 = 2048$.
- Add random factors to RGB values to simulate variations in lighting.
- These steps help the network generalize better.

AlexNet Architecture



All hidden layers were split in two and trained on different GPU boards due to GPU memory limitations.

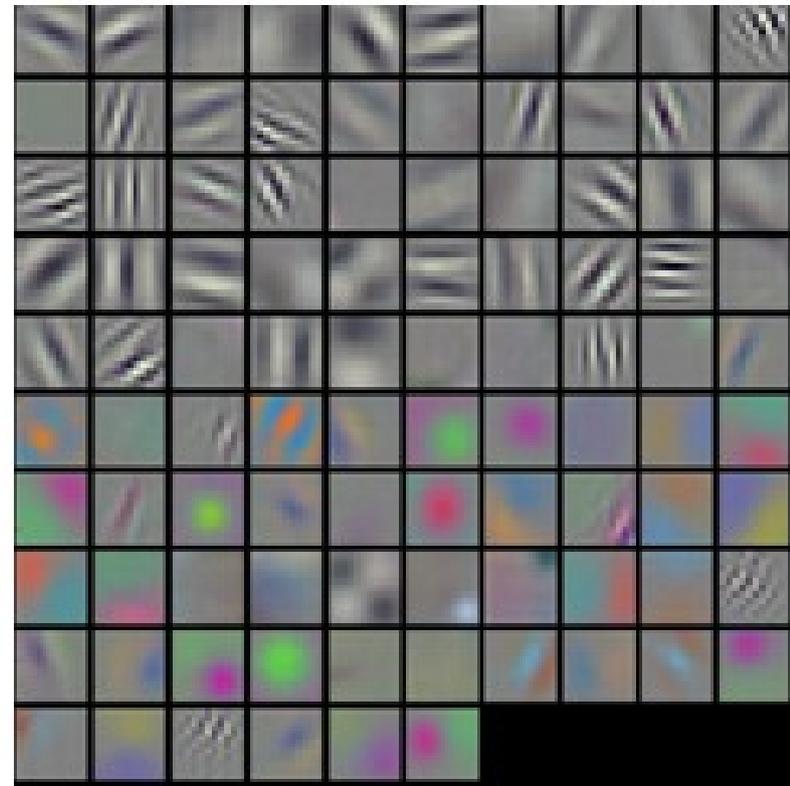
AlexNet Layer 1 Kernels

AlexNet's 96 11x11 layer 1 kernels.

First 48 trained on GPU 1 look for edges.

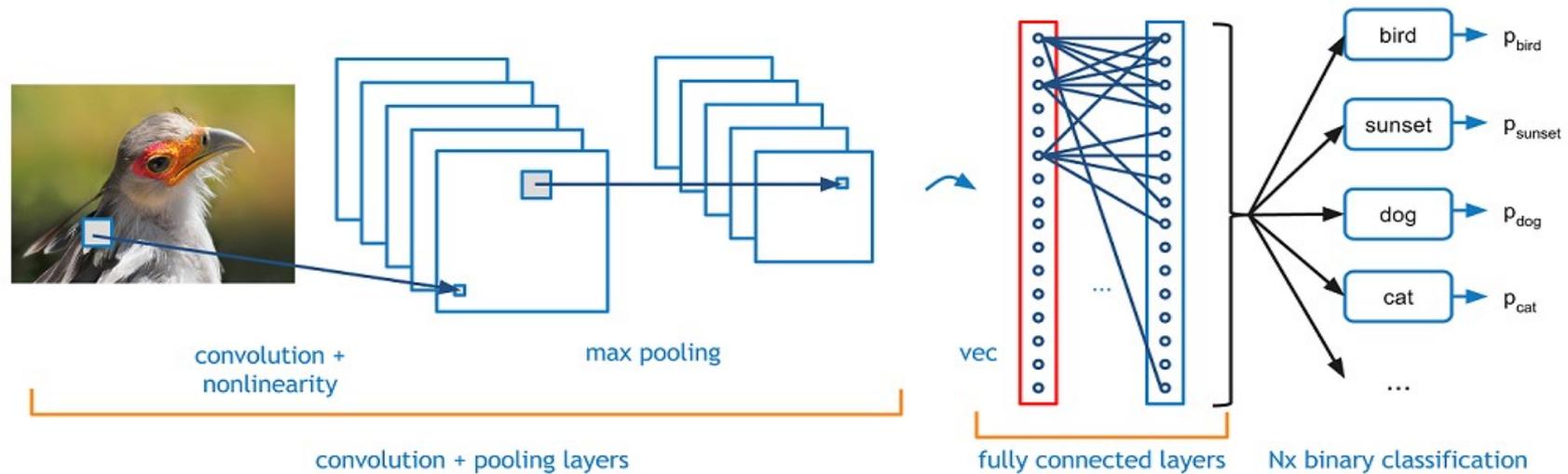
Second 48 trained on GPU 2 look for color.

This separation is a natural consequence of the normalization terms in the early layers.



Visualizations of filters

Generic Object Recognition CNN



<https://adeshpande3.github.io/A-Beginner%27s-Guide-To-Understanding-Convolutional-Neural-Networks/>

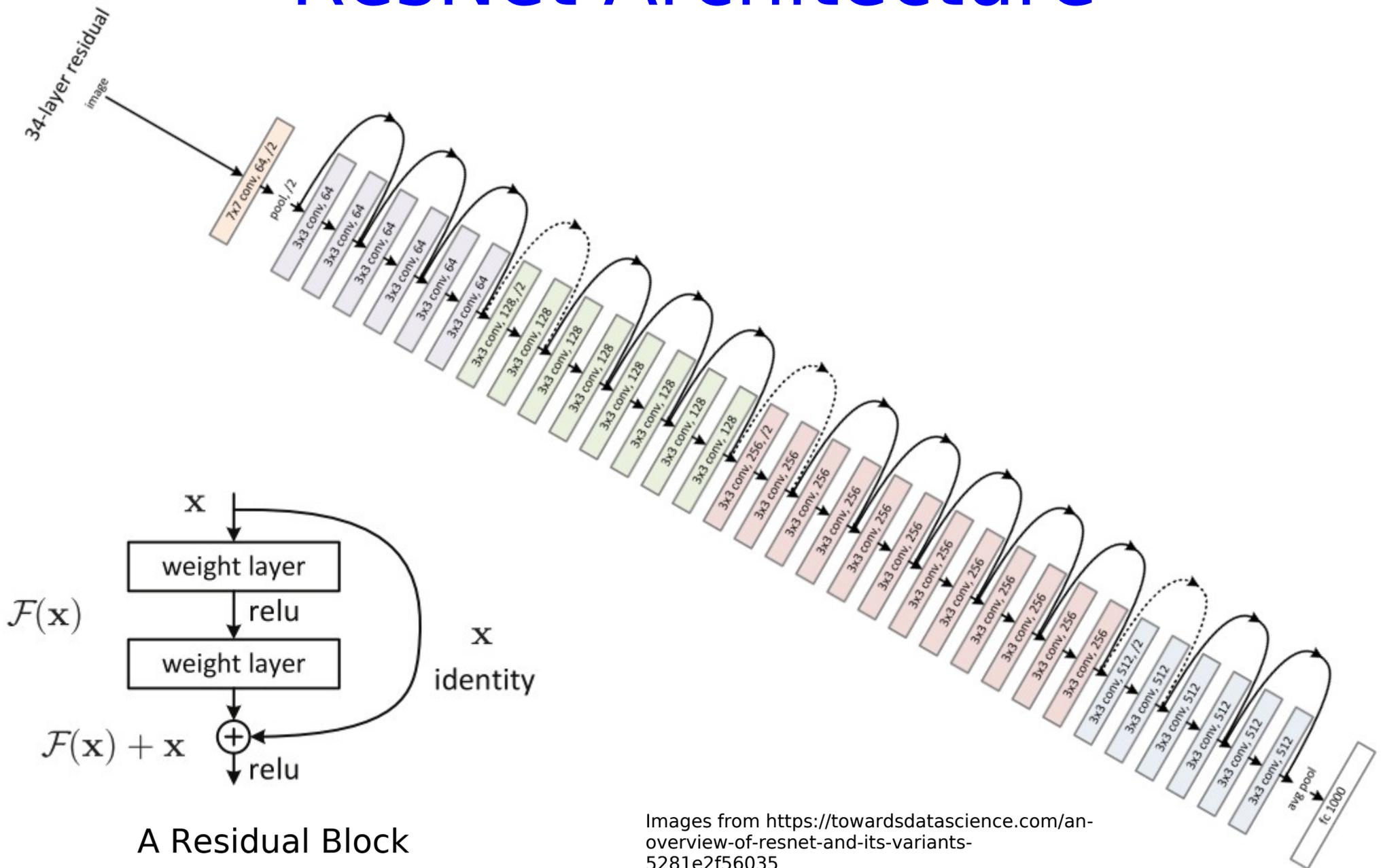
After AlexNet

- AlexNet had 8 layers: 5 convolutional and 3 fully connected.
- In 2015 Microsoft won the ILSVRC using a deep neural network with 100 layers.
- By the end of the ILSVRC in 2017, the best entrants were seeing accuracies of over 95% (error rate $< 5\%$).

Residual Blocks

- Residual blocks were introduced in ResNet:
 - For really deep networks, it's hard for the error signal to propagate backwards through many layers.
 - Solution: add shortcut connections, e.g., from layer i to layer $i+2$, so that error can back-propagate more quickly.
 - A residual block contains hidden layers with a shortcut connection.

ResNet Architecture

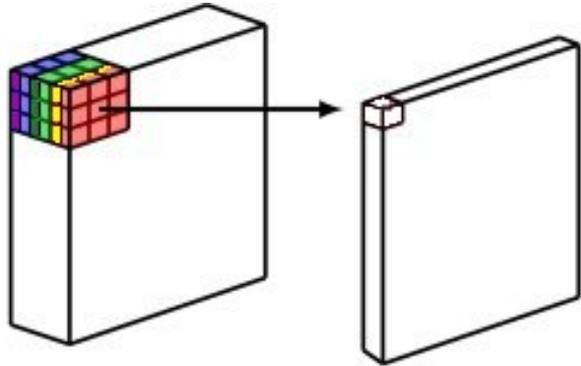


Images from <https://towardsdatascience.com/an-overview-of-resnet-and-its-variants-5281e2f56035>

Mobile Implementations

- People want to implement computer vision on mobile phones. Networks must be reduced in size.
- Various architectures explore ways to reduce the size of the network and the number of multiply-add operations.
 - Separable convolutions
 - Bottlenecks
- Examples: MobileNet, SqueezeNet

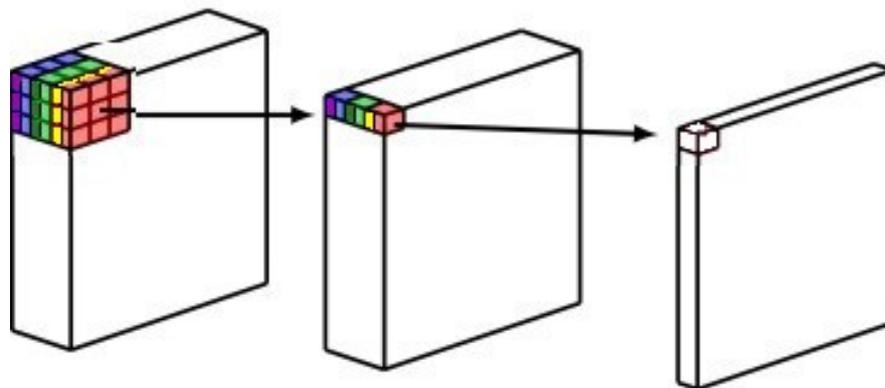
Separable Convolutions



(a) Conventional Convolutional Neural Network

3x3 kernel over 6 channels

$$3 \times 3 \times 6 = 54 \text{ weights}$$



Depthwise Convolution

Pointwise Convolution

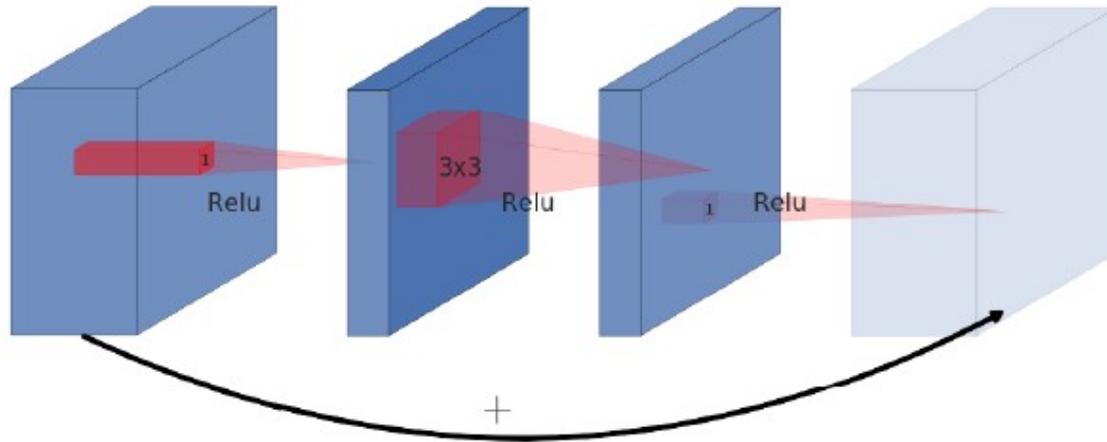
(b) Depthwise Separable Convolutional Neural Network

One 3x3 kernel applied to all 6 channels (depthwise convolution)

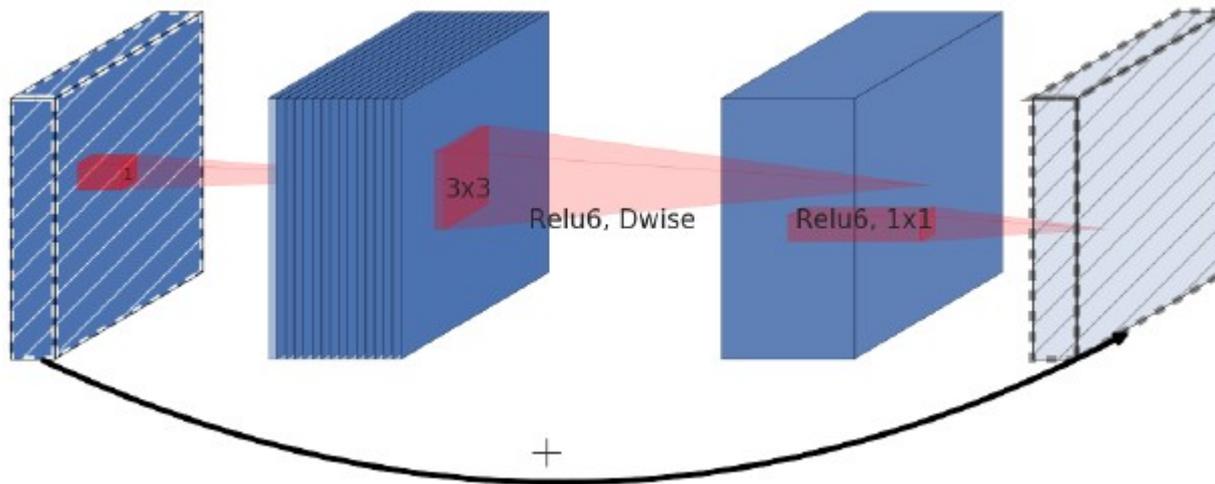
Linear weighted combination of the results (pointwise convolution)

$$3 \times 3 + 6 = 15 \text{ weights}$$

Bottlenecks with Residuals



MobileNet:
residual
bottleneck



MobileNetV2:
inverted residual
bottleneck

PyTorch Vision Models

- PyTorch contains several pre-trained object recognition models, including AlexNet, ResNet, Inception, VGG, and MobileNetV2.
- Look in `torchvision.models` for a list.
- Models are trained on the ImageNet dataset.

MobileNetV2 on Cozmo

- See the course's demos folder.
- Uses pre-trained MobileNetV2 model from `torchvision.models`.
- Feeds a 224x224 Cozmo camera image into the network and reports the top 5 labels.

Transfer Learning

- How can we quickly train a visual classifier for a new object class?
- Use the last hidden layer of a pre-trained ImageNet classifier as a feature vector.
- Train a classifier on the new categories using just 1-2 layers of trainable weights, or just use k-nearest neighbor.
- This is how Teachable Machine works.

Teachable Machine

<https://teachablemachine.withgoogle.com>

The screenshot displays the Teachable Machine 2.0 interface. At the top left, the Google logo is followed by the text "Teachable Machine 2.0: Making AI easier for everyone" and the name "Jordan". In the top right corner, there is a "Watch later" button. The main interface is divided into two sections. On the left, a large image of a woman with long brown hair is shown against a red background. Below this image is a blue button labeled "Hold to Record" and a gear icon. A "MORE VIDEOS" button is also visible. On the right, a panel titled "30 Image Samples" displays a grid of 30 small images of the same woman, arranged in 5 rows and 6 columns. Below the image grid, a "Training" panel is visible, featuring a "Train Model" button and an "Advanced" dropdown menu. At the bottom of the screen, a video player controls bar shows a play button, a volume icon, and the time "0:59 / 2:08". The YouTube logo and a full-screen icon are also present in the bottom right corner.