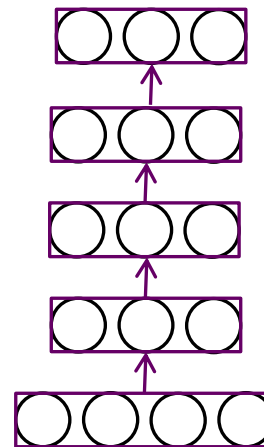# Pretrained Models and Transfer learning

Shailesh S

# Why Stop At One Hidden Layer?

E.g., vision hierarchy for recognizing handprinted text

- Word                         output layer
- Character               hidden layer 3
- Stroke                   hidden layer 2
- Edge                      hidden layer 1
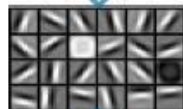- Pixel               input layer

3rd layer
"Objects"

2nd layer
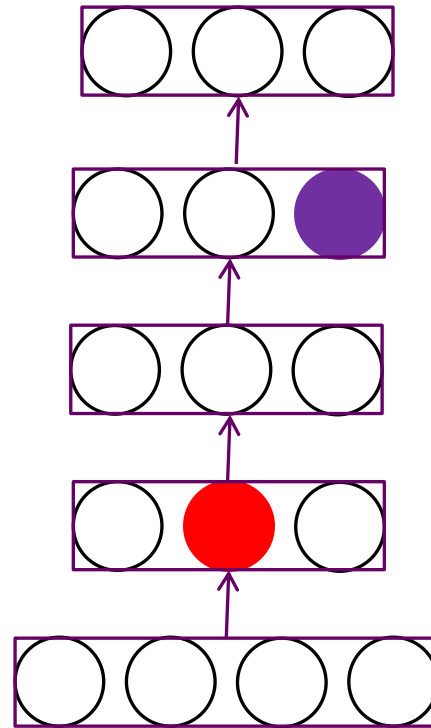"Object parts"

1st layer
"Edges"

Pixels

# Challenges in Deep Learning

- Requires Large dataset for Training
- Overfitting or lack of generalization
- Vanishing or Exploding Gradient Problem
- Appropriate Learning rate
- Covariate Shift
- Effective Training

# Why Deeply Layered Networks Fail

Credit assignment problem
- ◦ How is a neuron in layer 2 supposed to know what it should output until all the neurons above it do something sensible?
- ◦ How is a neuron in layer 4 supposed to know what it should output until all the neurons below it do something sensible?
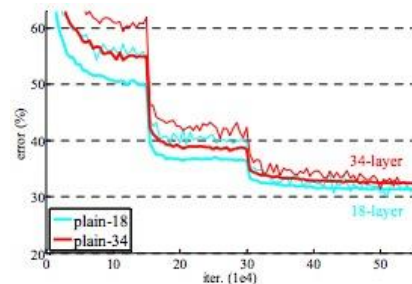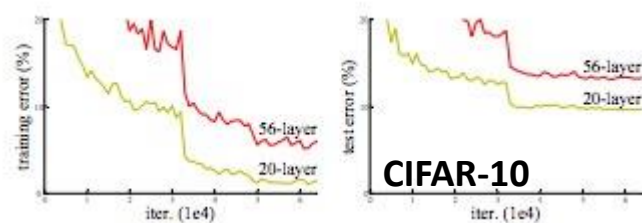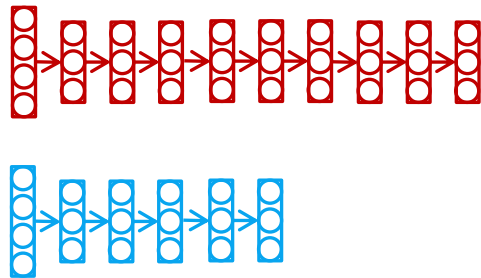
# Deeper Vs. Shallower Nets



Deeper net can represent any mapping that shallower net can
◦ Use identity mappings for the additional layers

Deeper net in principle is more likely to overfit

But in practice it often underfits *on the training set*
◦ Degradation due to harder credit-assignment problem
◦ Deeper isn't always better!



**CIFAR-10**



**ImageNet**
thin=train,
thick=valid.

**He, Zhang, Ren, and Sun (2015)**

# Why Deeply Layered Networks Fail

Vanishing gradient problem
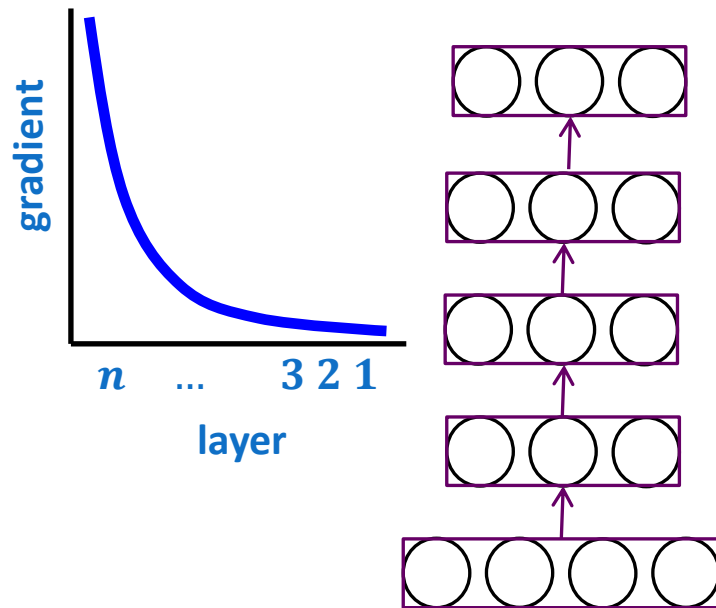- With logistic or tanh units

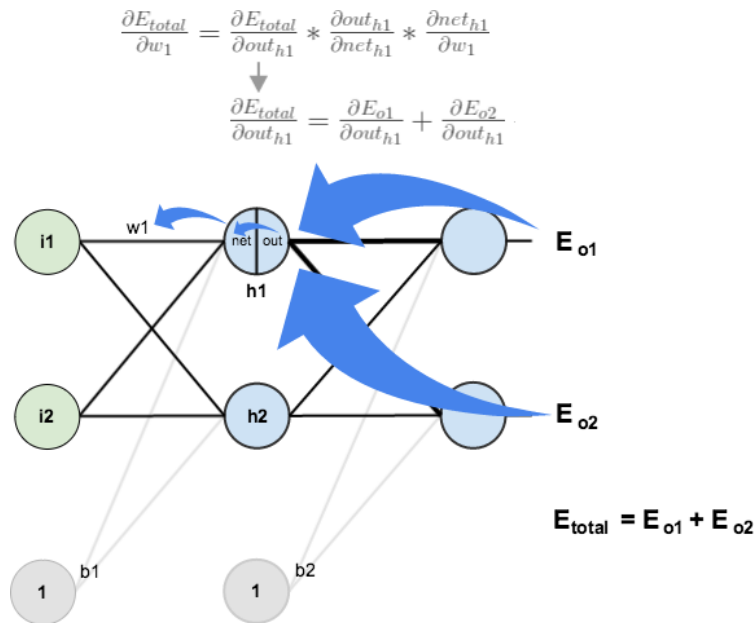$$y_j = \frac{1}{1 + \exp(-z_j)}$$

$$y_j = \tanh(z_j)$$

$$\frac{\partial y_j}{\partial z_j} = y_j(1 - y_j)$$

$$\frac{\partial y_j}{\partial z_j} = (1 + y_j)(1 - y_j)$$

- Error gradients get squashed as they are passed back through a deep network

# Why Deeply Layered Networks Fail

# Why Deeply Layered Networks Fail



$$y = \max(\mathbf{0}, \mathbf{z})$$

$$\frac{\partial y}{\partial z} = \begin{cases} 0 & z \leq 0 \\ 1 & \text{otherwise} \end{cases}$$
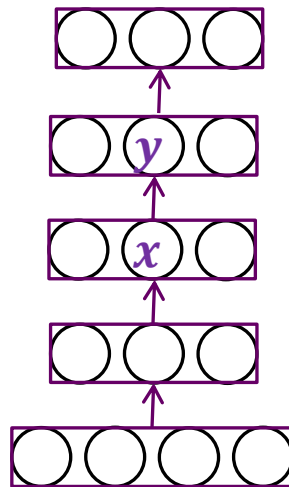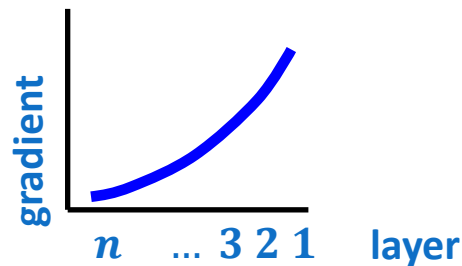
Exploding gradient problem
◦ with linear or ReLU units

$$\frac{\partial y}{\partial x} = \begin{cases} 0 & y = 0 \\ w_{yx} & \text{otherwise} \end{cases}$$

◦ can be a problem when $\left| \sum_y \frac{\partial y}{\partial x} \right| > 1$

# Trick From Long Ago To Avoid Local Optima

Add direct connections from input to output layer

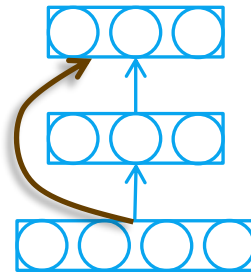Easy bits of the mapping are learned by the direct connections
◦ Easy bit = linear or saturating-linear functions of the input
◦ Often captures 80-90% of the variance in the outputs

Hidden units are reserved for learning the hard bits of the mapping
◦ They don't need to learn to copy activations forward for the linear portion of the mapping

Problem
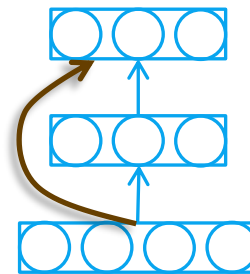◦ Adds a lot of free parameters

# Latest Tricks

Novel architectures that
◦ skip layers
◦ have linear connectivity between layers

Advantage over direct-connection architectures
◦ no/few additional free parameters
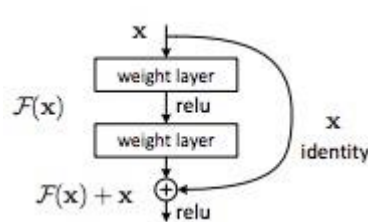
# Deep Residual Networks (ResNet)

Add linear short-cut connections to architecture

Basic building block
- $\mathcal{F}(x) = W_2 \sigma(W_1 x)$
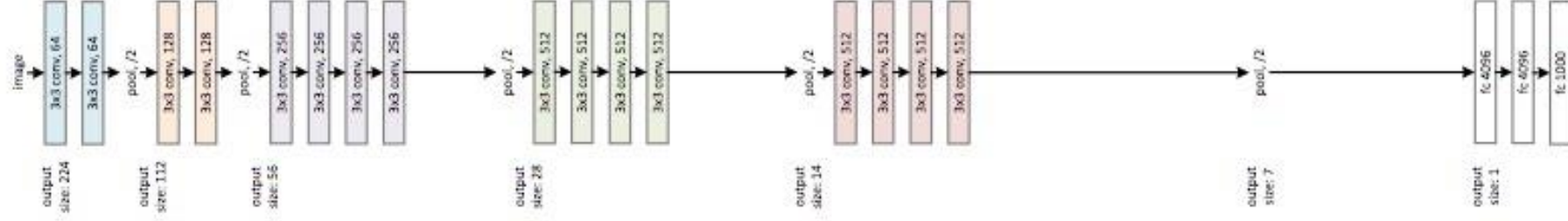- $y = \sigma(\mathcal{F}(x) + x)$



**activation flow**

Variations
- allow different input and output dimensionalities: $y = \sigma(\mathcal{F}(x) + W_3 x)$
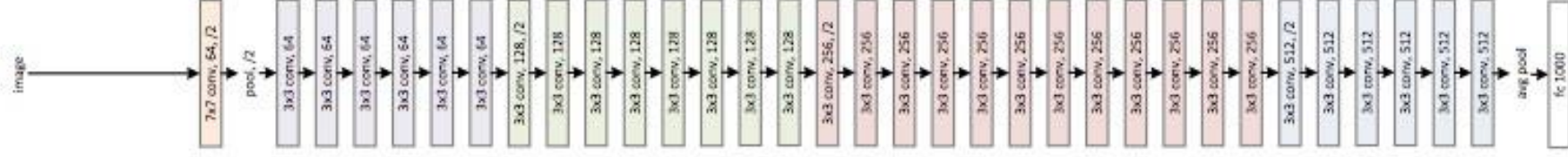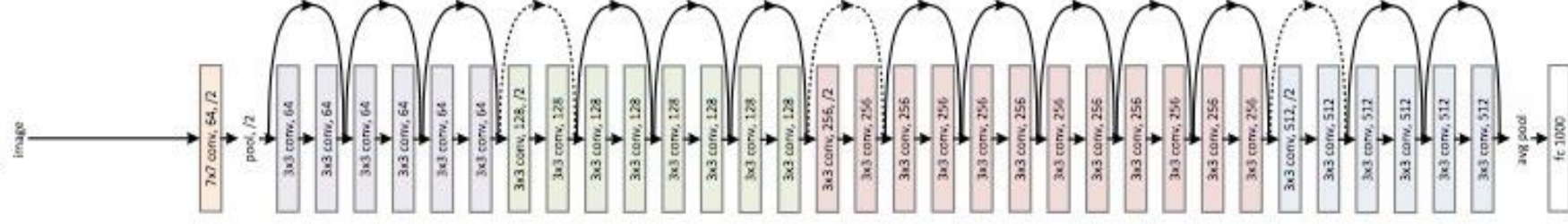- $\mathcal{F}(x)$ can have an arbitrary number of layers

He, Zhang, Ren, and Sun (20...
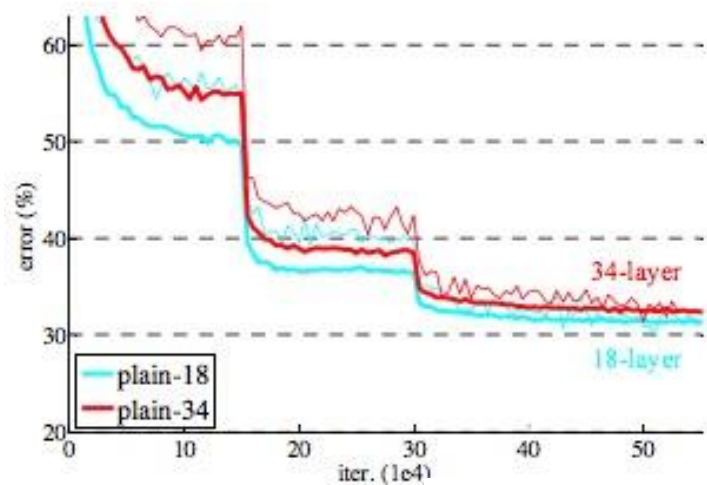
# ResNet



| | plain | ResNet |
|---|---|---|
| 18 layers | 27.94 | 27.88 |
| 34 layers | 28.54 | **25.03** |

**Top-1 Error %**

He, Zhang, Ren, and Sun (2015)

# Inception Networks

# Inception

- traditional deep neural networks, which is well known for stacked convolutions
- Inception networks are well known for parallel stacked convolutions
- Inception block contains different size of convolutions and one max pooling layers which is arranged in a parallel manner.It is aggregated for the final output for the inception block
- I.e, Computing 1*1,3*3, 5*5  convolutions with in the same module of the network which is arranged in a parallel manner

# Inception

- In every image, there is a huge variation in the location of information
- Filter size should be chosen according to the information content

# Google Net



(a) Inception module, naïve version

# Architecture based on Inception network

- Inception v1(2014)
- Inception v3(2015)
- Inception v3(2015)
- Inception v4(2016)
- Inception ResNet(2016)

Adaptability of the different scales can be accommodated

Overfitting can be prevented

(a) Inception module, naïve version

Deep neural networks are computationally expensive. To make it cheaper, the author limit the number of input channels by adding an extra 1*1 convolution before 5*5 and 3*3 convolutions

# Transfer Learning

# Transfer Learning

- Transfer learning generally refers to a process where a model trained on one problem is used in some way on a second related problem.

- The weights in re-used layers may be used as the starting point for the training process and adapted in response to the new problem.

- This usage treats transfer learning as a type of weight initialization scheme.

- This may be useful when the first related problem has a lot more labeled data than the problem of interest and the similarity in the structure of the problem may be useful in both contexts.

- **Transfer learning** consists of taking features learned on one problem, and leveraging them on a new, similar problem.

- Taking model learned on one problem , and update it for the new ,similar problem
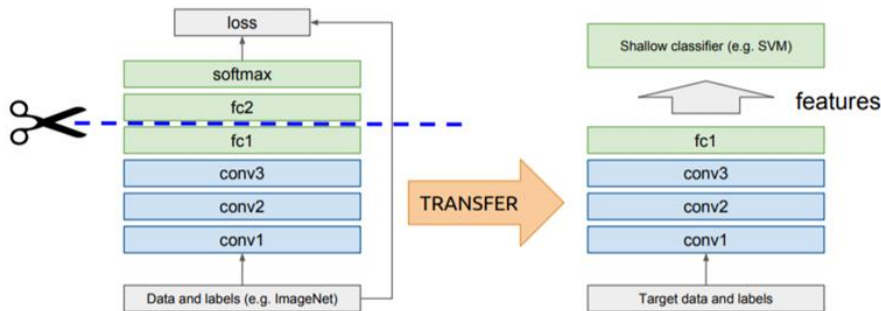
# Feature Extraction

The pre-trained model can be used as a separate feature extraction program, in which case input can be pre-processed by the model or portion of the model to a given an output (e.g. vector of numbers) for each input image, that can then use as input when training a new model.

Idea: use outputs of one or more layers of a network trained on a different task as generic feature detectors. Train a new shallow model on these features.

Assumes that $D_S = D_T$

# Freezing layers:

- Take layers from previously trained model
- Freezing them so as to avoid destroying any of the information during future training rounds
- Add some new trainable layers on top of the frozen layers
- Train the new layers on your dataset
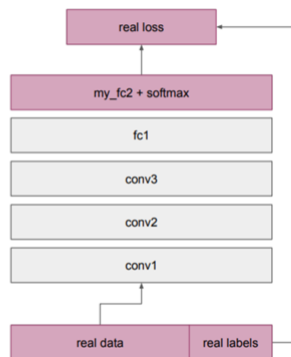
## Fine-tuning: supervised domain adaptation

Train deep net on "nearby" task for which it is easy to get labels using standard backprop

- E.g. ImageNet classification
- Pseudo classes from augmented data
- Slow feature learning, ego-motion

Cut off top layer(s) of network and replace with supervised objective for target domain

**Fine-tune** network using backprop with labels for target domain until validation loss starts to increase

Aligns $D_S$ with $D_T$

## Freeze or fine-tune?

Bottom $n$ layers can be frozen or fine tuned.

- **Frozen**: not updated during backprop
- **Fine-tuned**: updated during backprop

Which to do depends on target task:

- **Freeze**: target task labels are scarce, and we want to avoid overfitting
- **Fine-tune**: target task labels are more plentiful

In general, we can set learning rates to be different for each layer to find a tradeoff between freezing and fine tuning