

Object Recognition – CNN Models

Danna Gurari

University of Colorado Boulder

Fall 2024



Review

- Last lecture: object recognition basics
 - Problem
 - Applications
 - Datasets
 - Evaluation metric
 - A popular solution: convolutional neural networks
- Assignments (Canvas)
 - Reading assignment was due earlier today
 - Next reading assignments due Wednesday and next Monday
 - Project proposal due in 2.5 weeks (review of requirements)
- Questions?

Object Recognition: Today's Topics

- ImageNet Challenge Top Performers
- Baseline Model: AlexNet
- VGG
- ResNet
- Summary of CNN Era

Object Recognition: Today's Topics

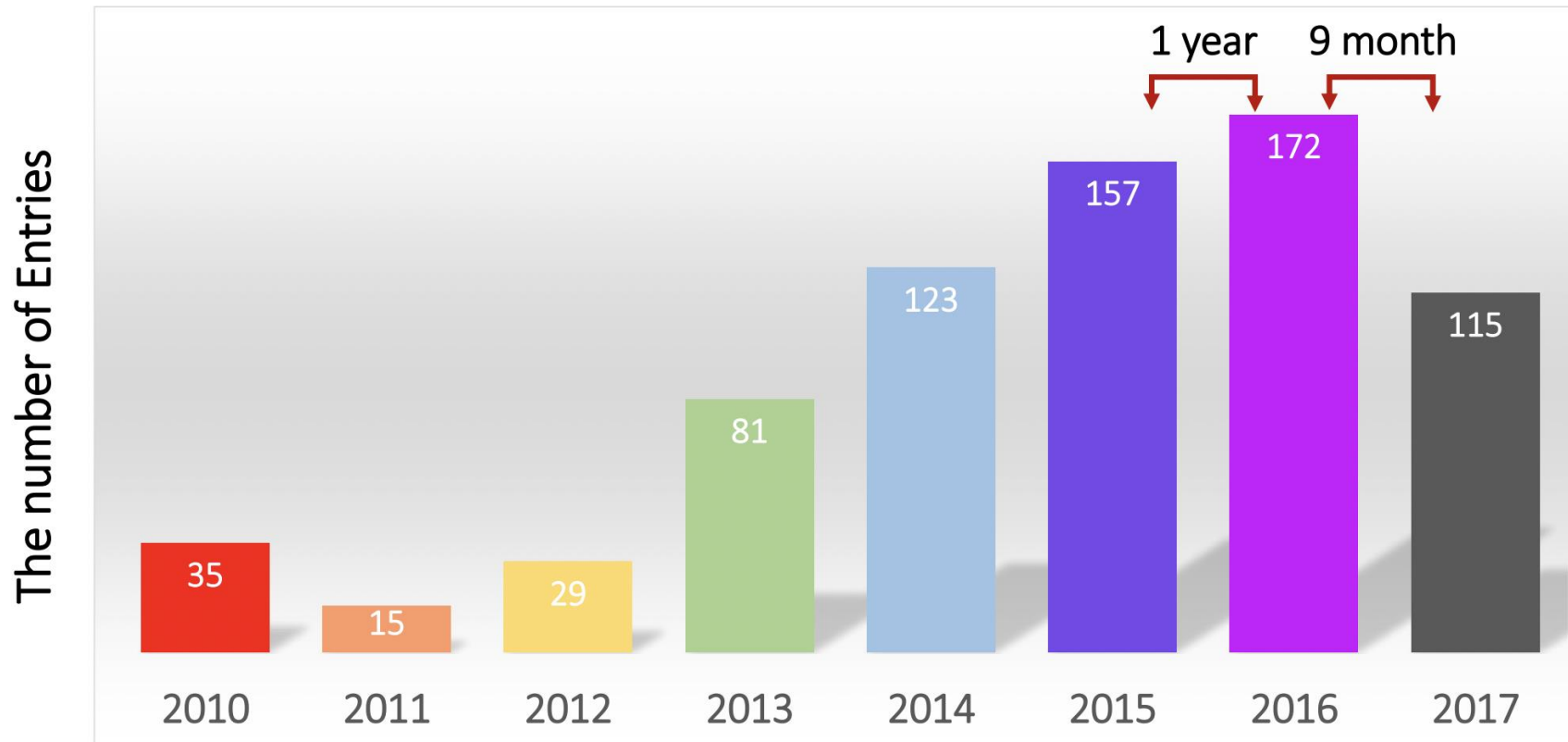
- ImageNet Challenge Top Performers
- Baseline Model: AlexNet
- VGG
- ResNet
- Summary of CNN Era

Recall Catalyst for Computer Vision Revolution: ImageNet Challenge (ILSVRC-2012 version)

- **Goal:** predict a category per image from 1000 options
- **Evaluation metric:** % correct (top-1 and top-5 predictions)
- **Dataset:** ~1.3 million images split into training, validation, and test sets
- **Source:** images scraped from search engines, such as Flickr, and labeled by crowdworkers



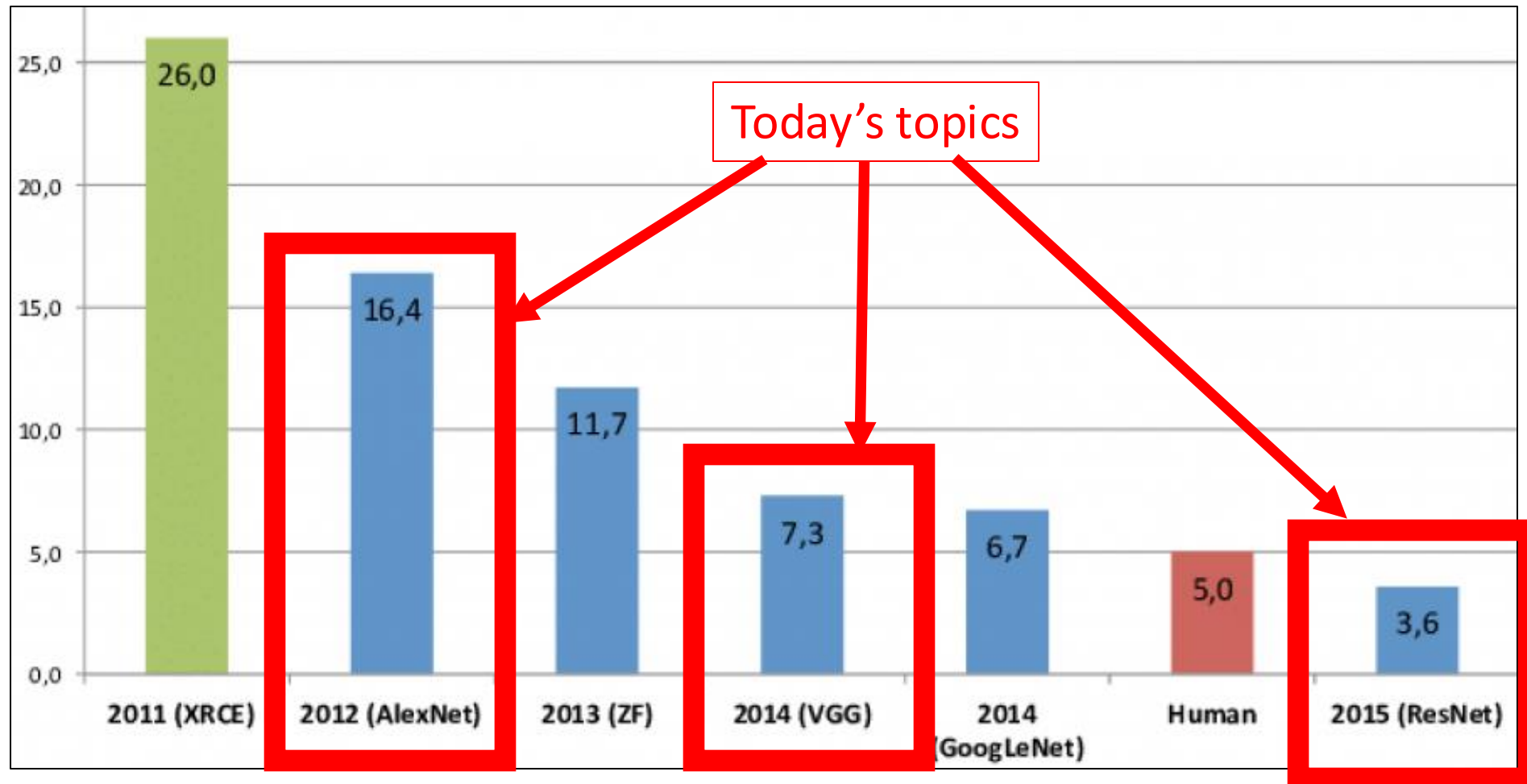
Recall: ImageNet Challenge Submissions



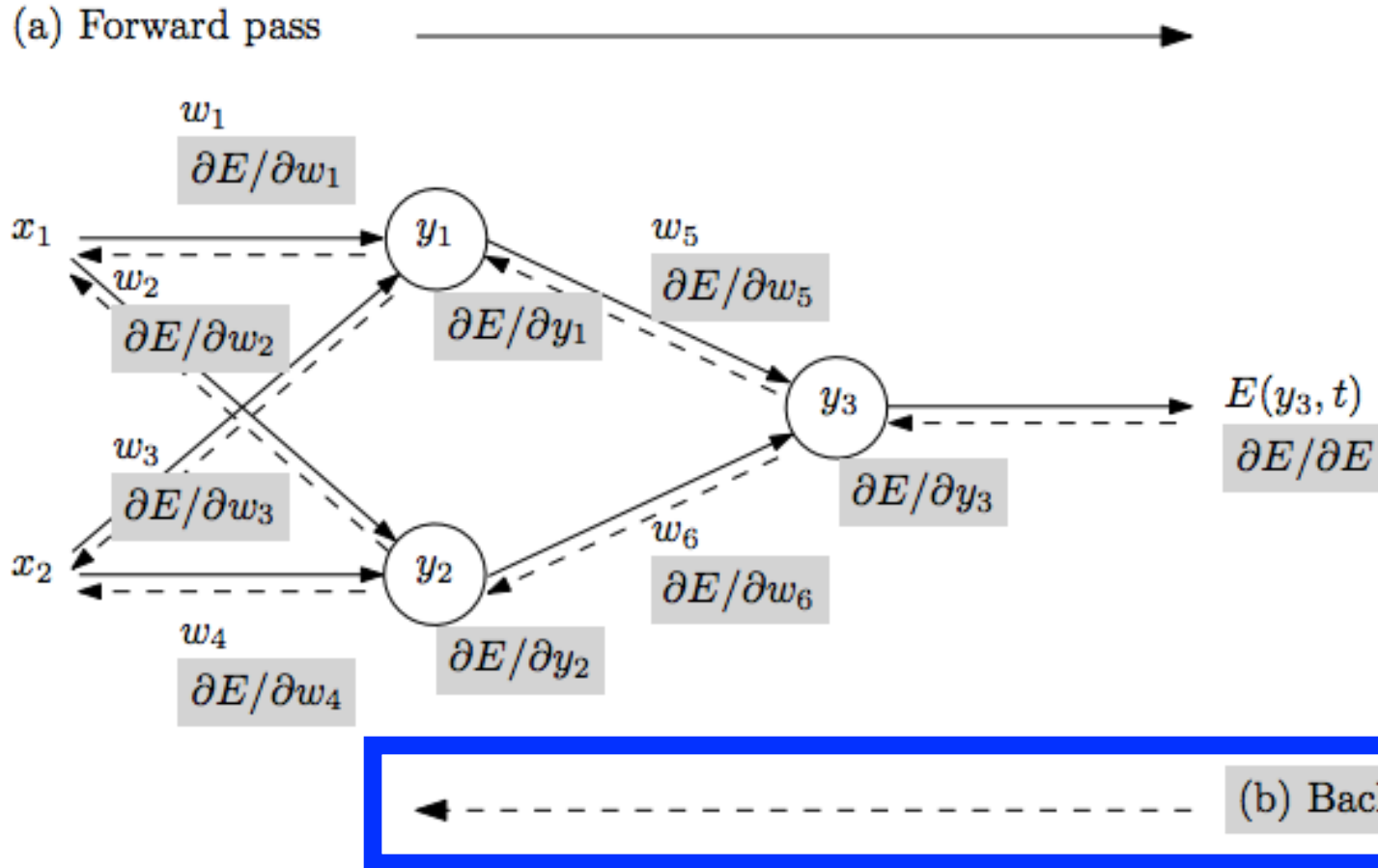
Inspired by AlexNet, many more researchers in the computer vision community proposed neural networks and showed how to make further progress over the years!

Secret Sauce for State-of-Art: Deeper CNNs

Progress of models on ImageNet (Top 5 Error)



Why It Is Difficult to Achieve Better Performance with CNNs That Are Deeper: Vanishing Gradients



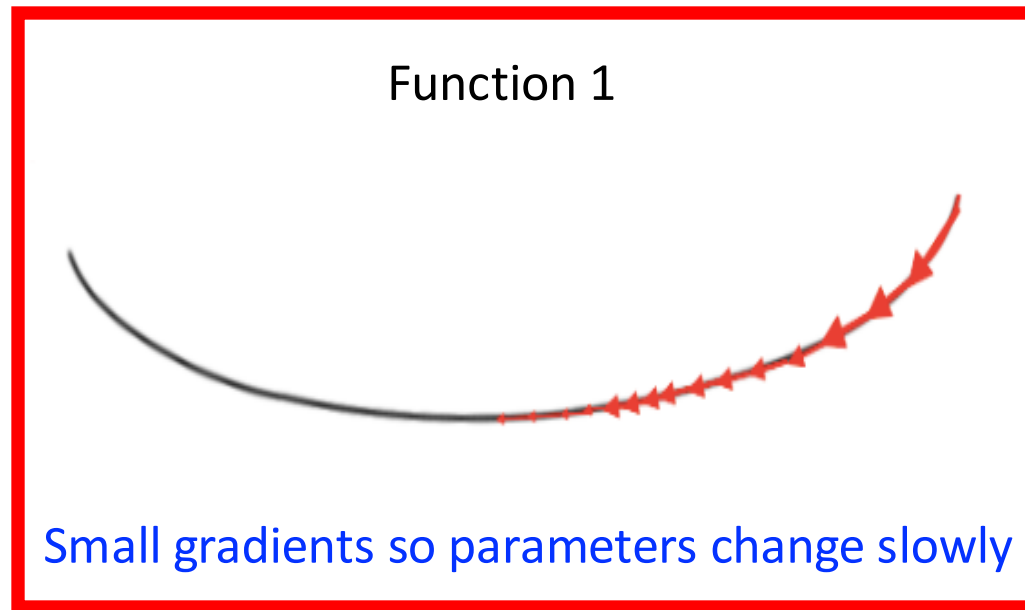
- Repeat until stopping criterion met:
 - 1. Forward pass:** propagate training data through model to make predictions
 - 2. Error quantification:** measure dissatisfaction with a model's predictions on training data
 - 3. Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
 - 4. Update each parameter using calculated gradients**

$$W_x = W_x - \alpha \left(\frac{\partial \text{Error}}{\partial W_x} \right)$$

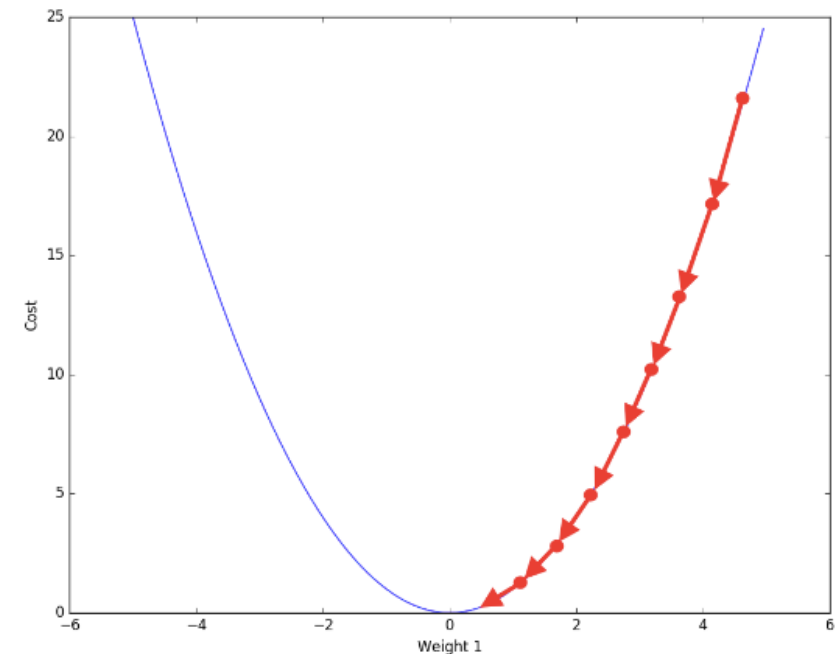
Why It Is Difficult to Achieve Better Performance with CNNs That Are Deeper: Vanishing Gradients

- **Want:** objective function with a gradient large enough to support (efficient) learning; e.g.,

Problem



Function 2

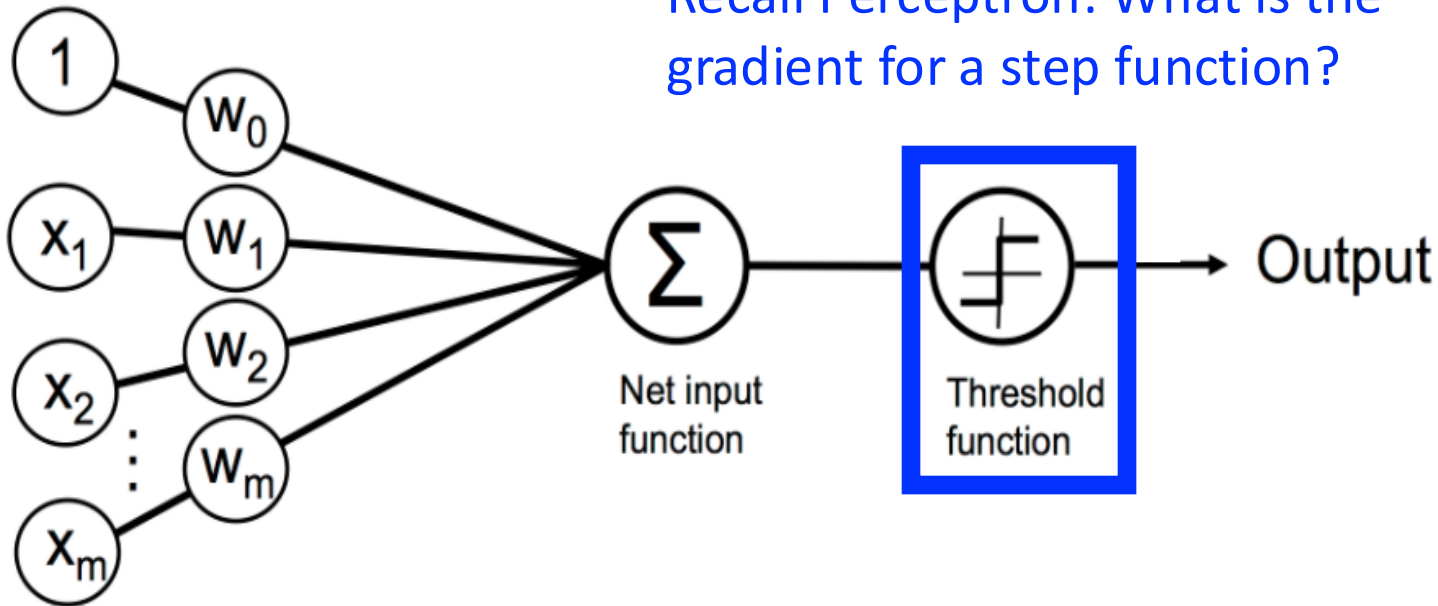


Large gradients so parameters change quickly

$$W_x = W_x - \alpha \left(\frac{\partial \text{Error}}{\partial W_x} \right)$$

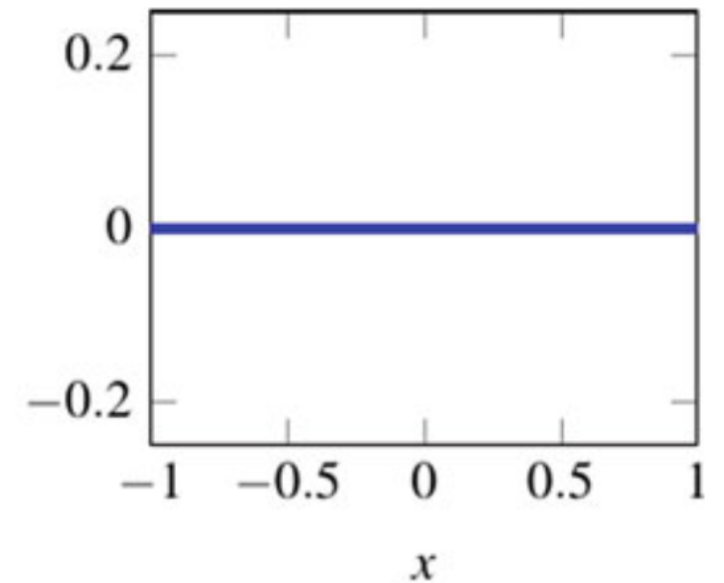
Why It Is Difficult to Achieve Better Performance with CNNs That Are Deeper: Vanishing Gradients

Recall Perceptron: What is the gradient for a step function?



Python Machine Learning; Raschka & Mirjalili

0 everywhere except 0 where it is non-differentiable



Deep Learning for NLP and Speech Recognition; Kamath, Liu, & Whitaker

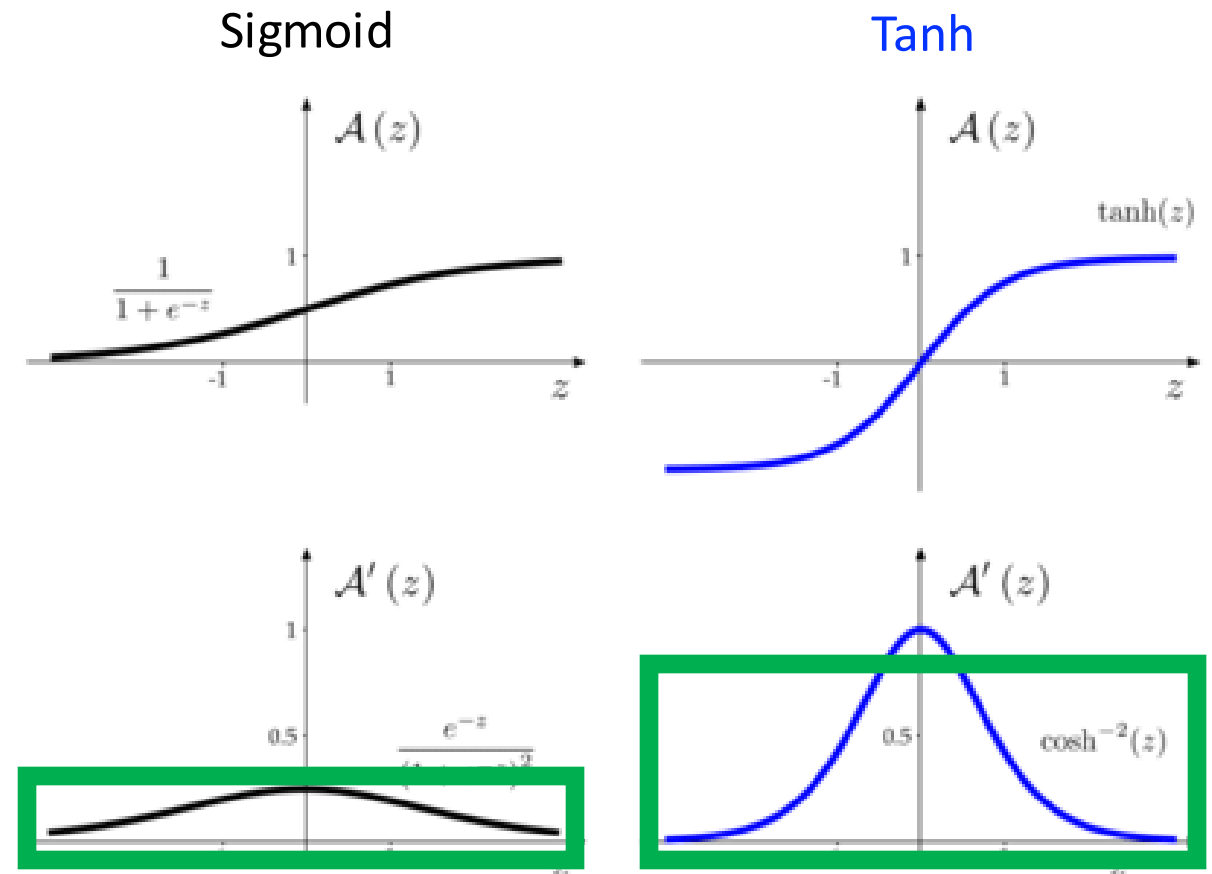
No gradient means model parameters wouldn't change with gradient descent!

Why It Is Difficult to Achieve Better Performance with CNNs That Are Deeper: Vanishing Gradients

Recall: Activation Functions

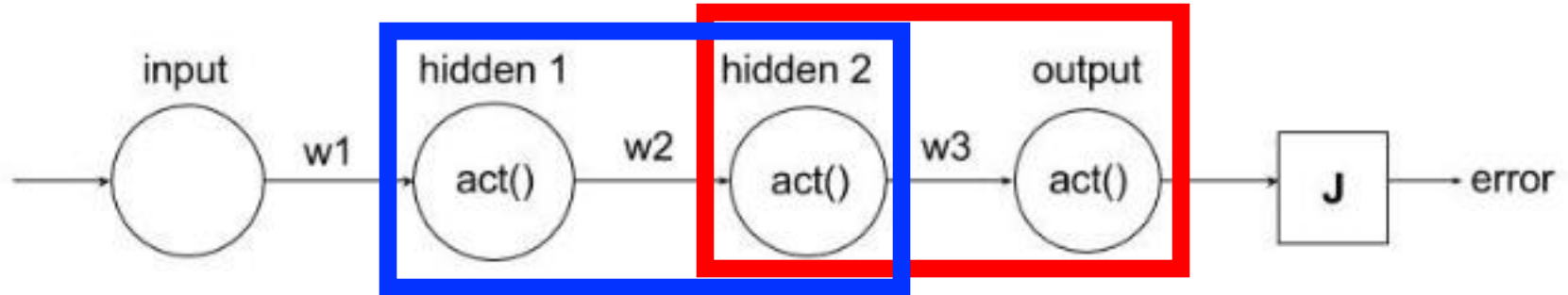
Small gradients limit amount model parameters change with gradient descent

$$W_x = W_x - \alpha \left(\frac{\partial \text{Error}}{\partial W_x} \right)$$



Why It Is Difficult to Achieve Better Performance with CNNs That Are Deeper: Vanishing Gradients

- Toy example:



- Error Derivative with respect to weight w1:

$$\frac{\partial error}{\partial w1} = \frac{\partial error}{\partial output} \cdot \frac{\partial output}{\partial hidden2} \cdot \frac{\partial hidden2}{\partial hidden1} \cdot \frac{\partial hidden1}{\partial w1}$$

e.g., derivative of sigmoid

activation function: (0 to 1/4]

e.g., derivative of sigmoid

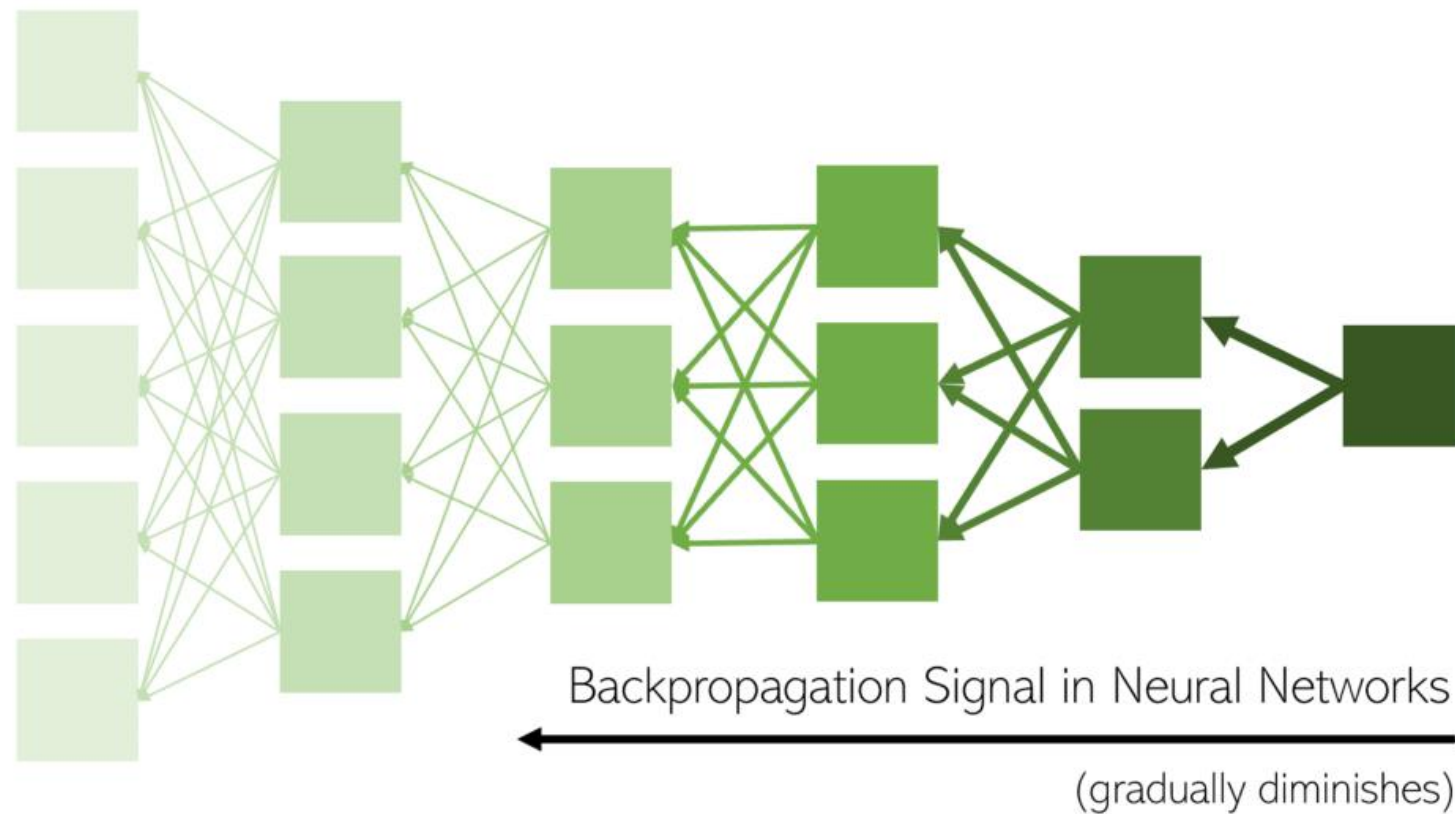
activation function: (0 to 1/4]

Problem: When multiplying more numbers smaller than 1, gradient decreases leading to reduced weight changes!

$$W_x = W_x - \alpha \left(\frac{\partial Error}{\partial W_x} \right)$$

Why It Is Difficult to Achieve Better Performance with CNNs That Are Deeper: Vanishing Gradients

Smallest gradients at **earliest layers make them slowest to train**, yet later layers depend on those earlier layers to do something useful; consequently, NNs struggle with garbage in means garbage out

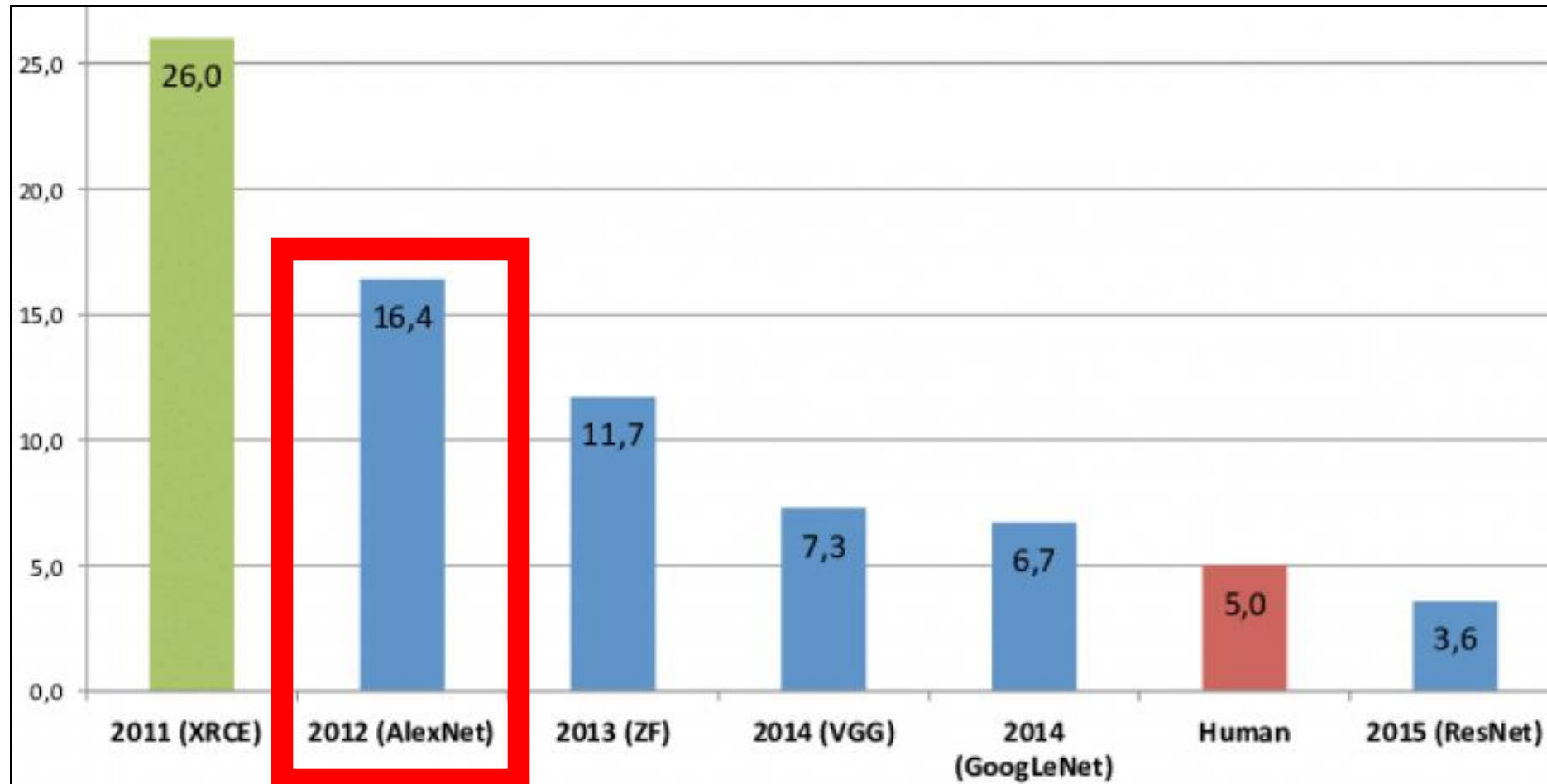


Object Recognition: Today's Topics

- ImageNet Challenge Top Performers
- **Baseline Model: AlexNet**
- VGG
- ResNet
- Summary of CNN Era

AlexNet: A Deeper CNN

Progress of models on ImageNet (Top 5 Error)



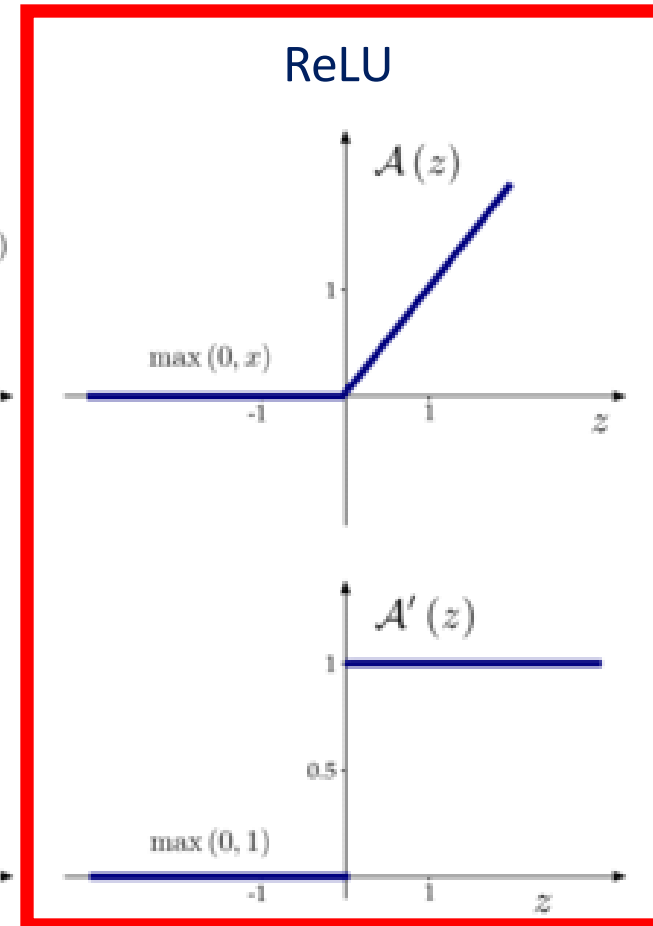
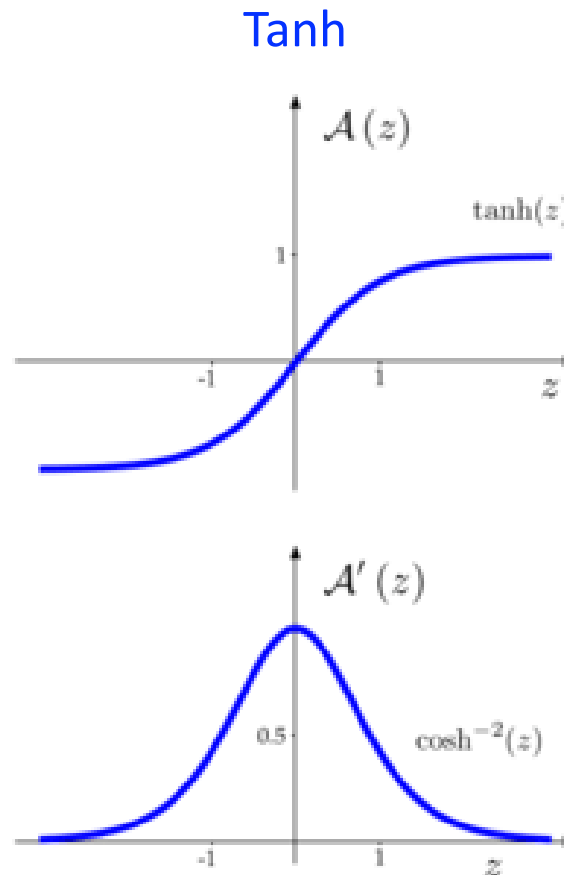
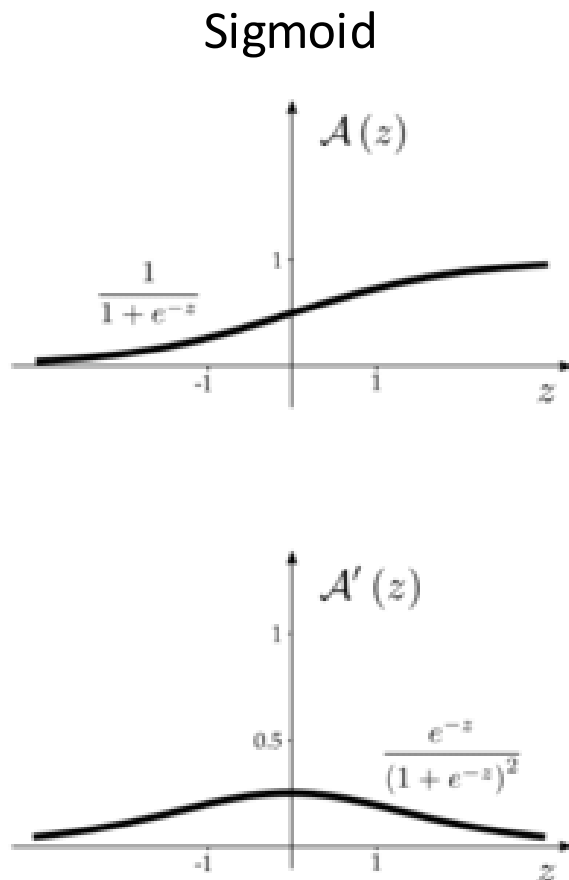
Why AlexNet?

Alex is the name of the paper's author 😊

[Alex Krizhevsky](#), Ilya Sutskever, and Geoffrey E. Hinton. "Imagenet classification with deep convolutional neural networks." *Advances in neural information processing systems* (2012).

Key Idea: Non-Saturating Activation Functions

Use activation functions with derivative value equal to 1 (i.e., $1 \times 1 \times 1 \dots$ doesn't vanish)



Benefits:

- Fast to compute
- Can preserve gradient and so support efficient learning

Architecture

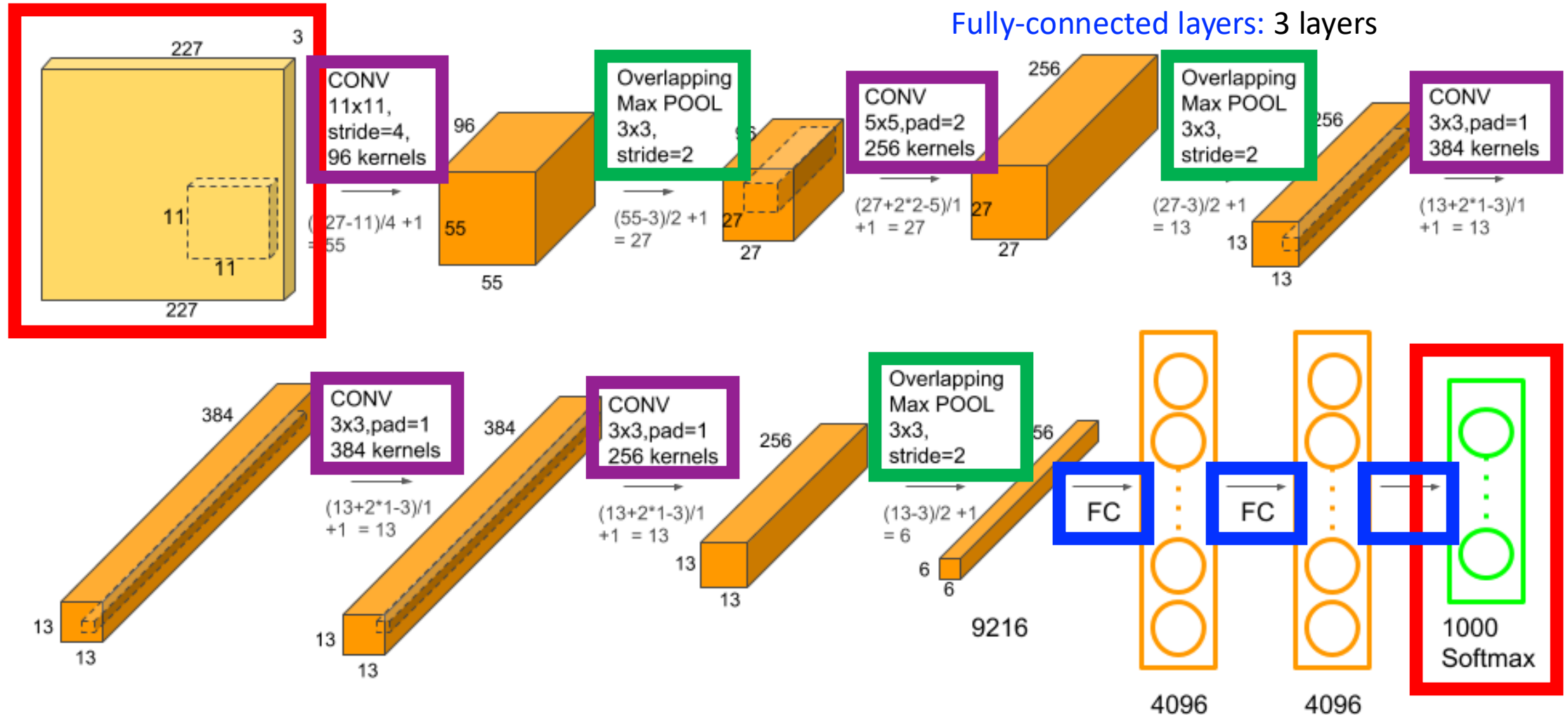
Input: RGB image resized to fixed input size

Output: 1000 class probabilities (sums to 1)

Convolutional layers: 5 layers

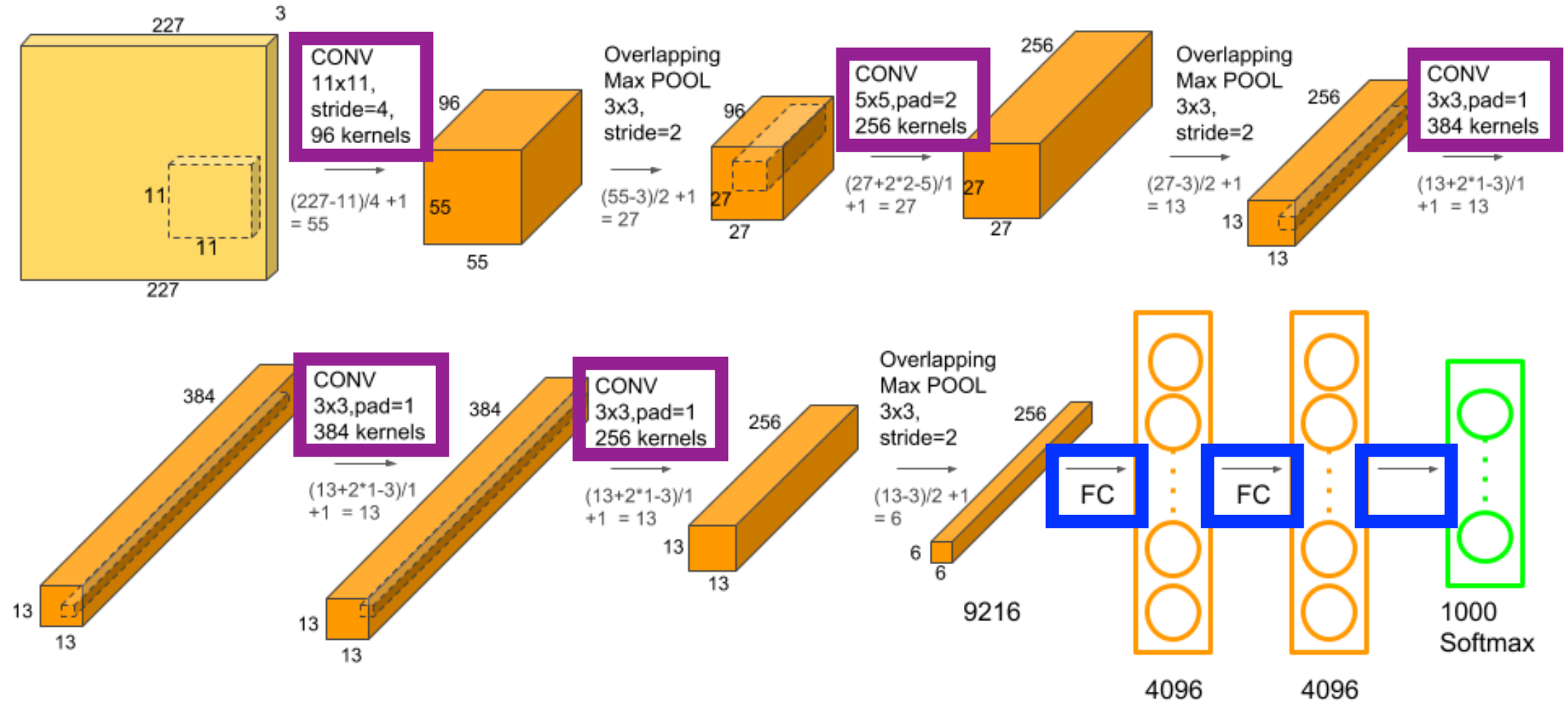
Pooling Layers: 3 layers

Fully-connected layers: 3 layers



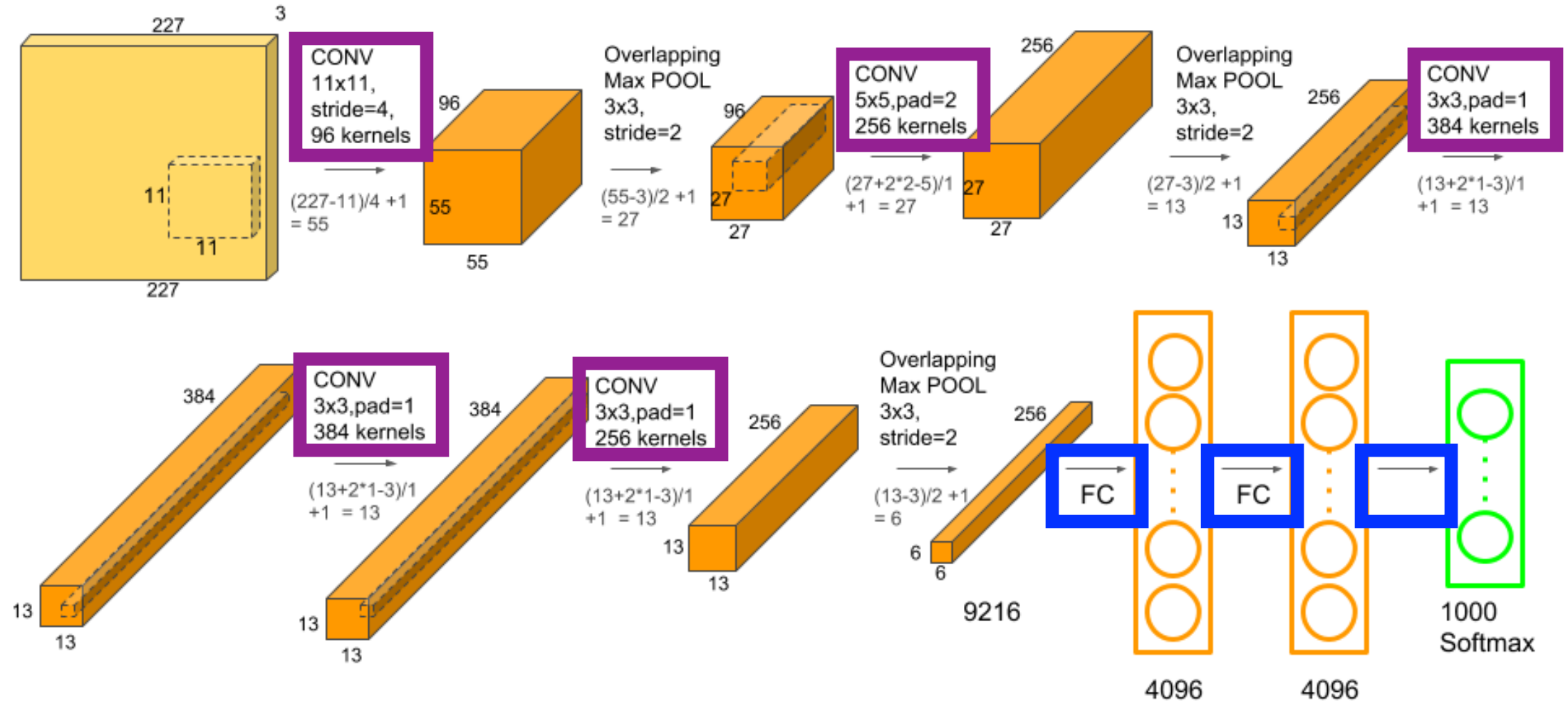
Architecture

How many layers have model parameters that need to be learned?



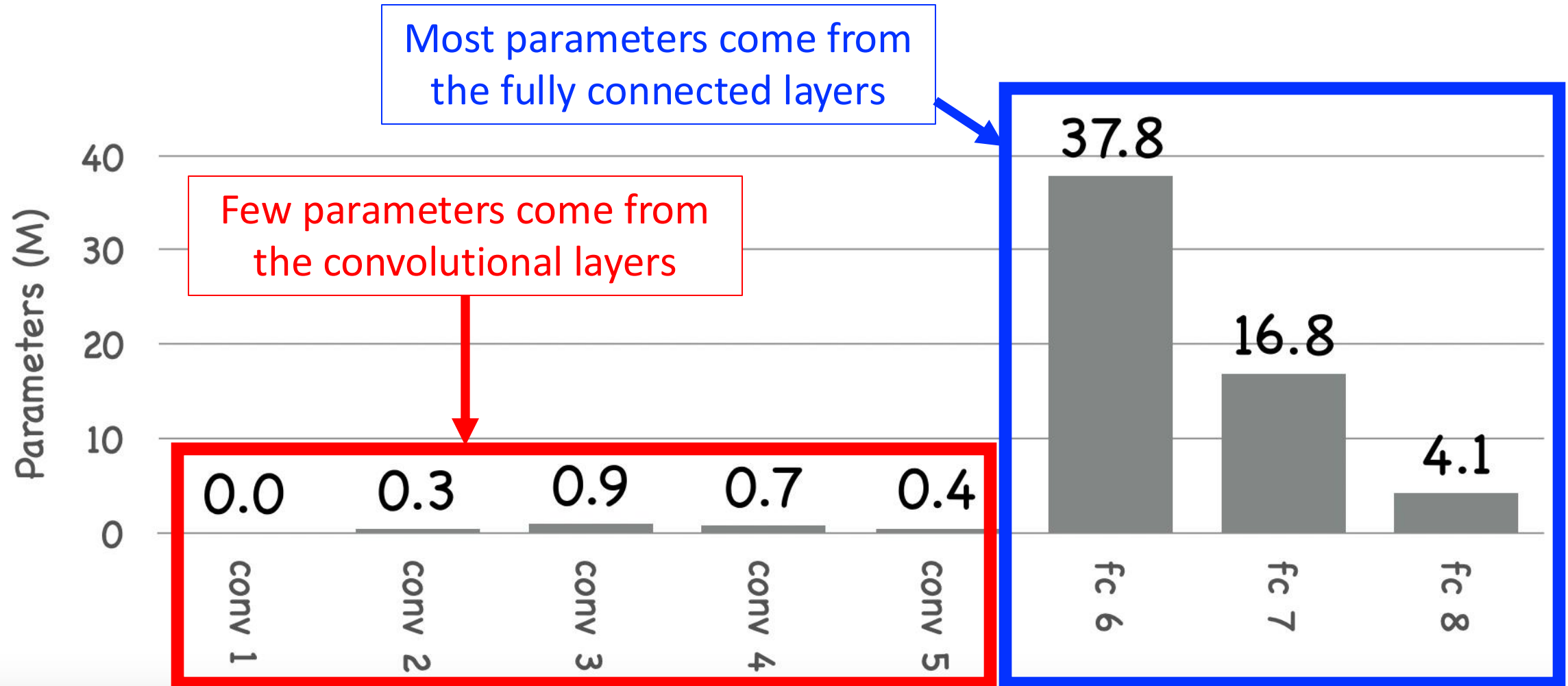
Architecture

Altogether, 60 million model parameters must be learned!



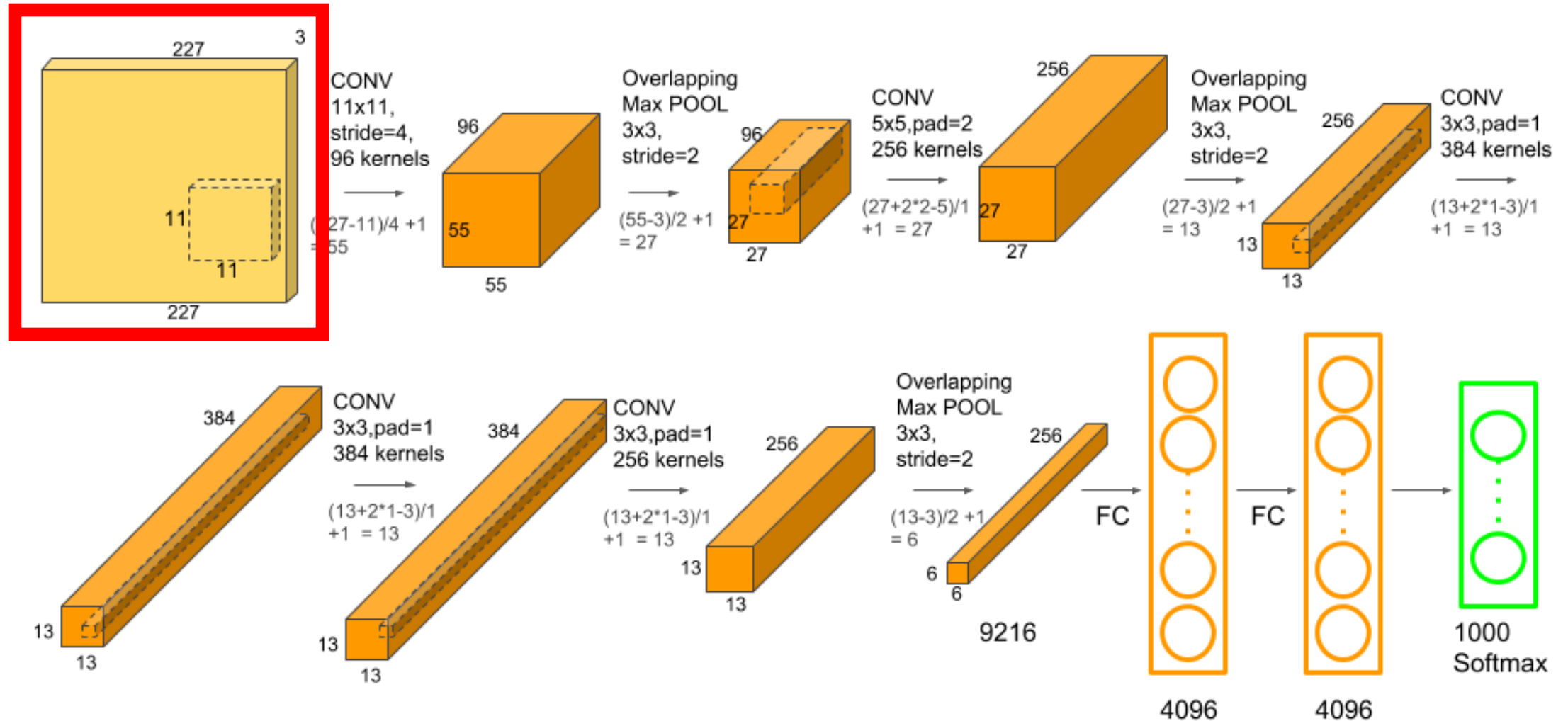
Architecture

Altogether, 60 million model parameters must be learned!

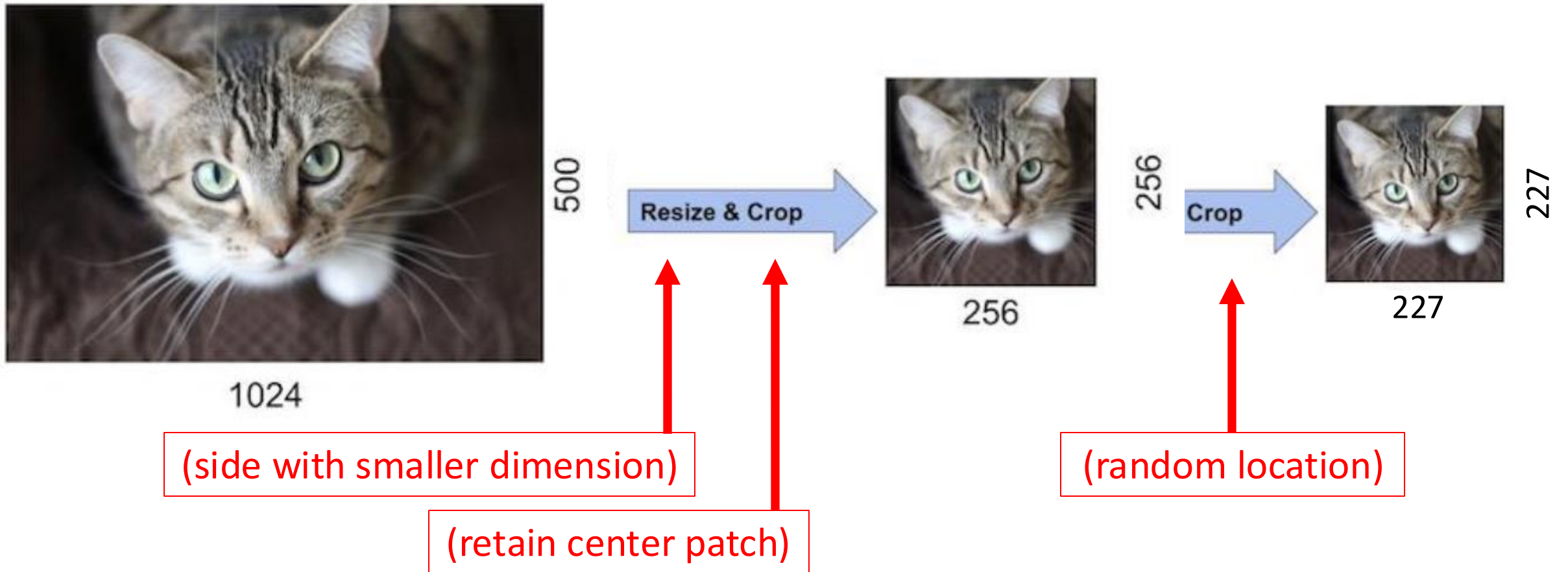


Input: RGB image resized to fixed input size

Architecture



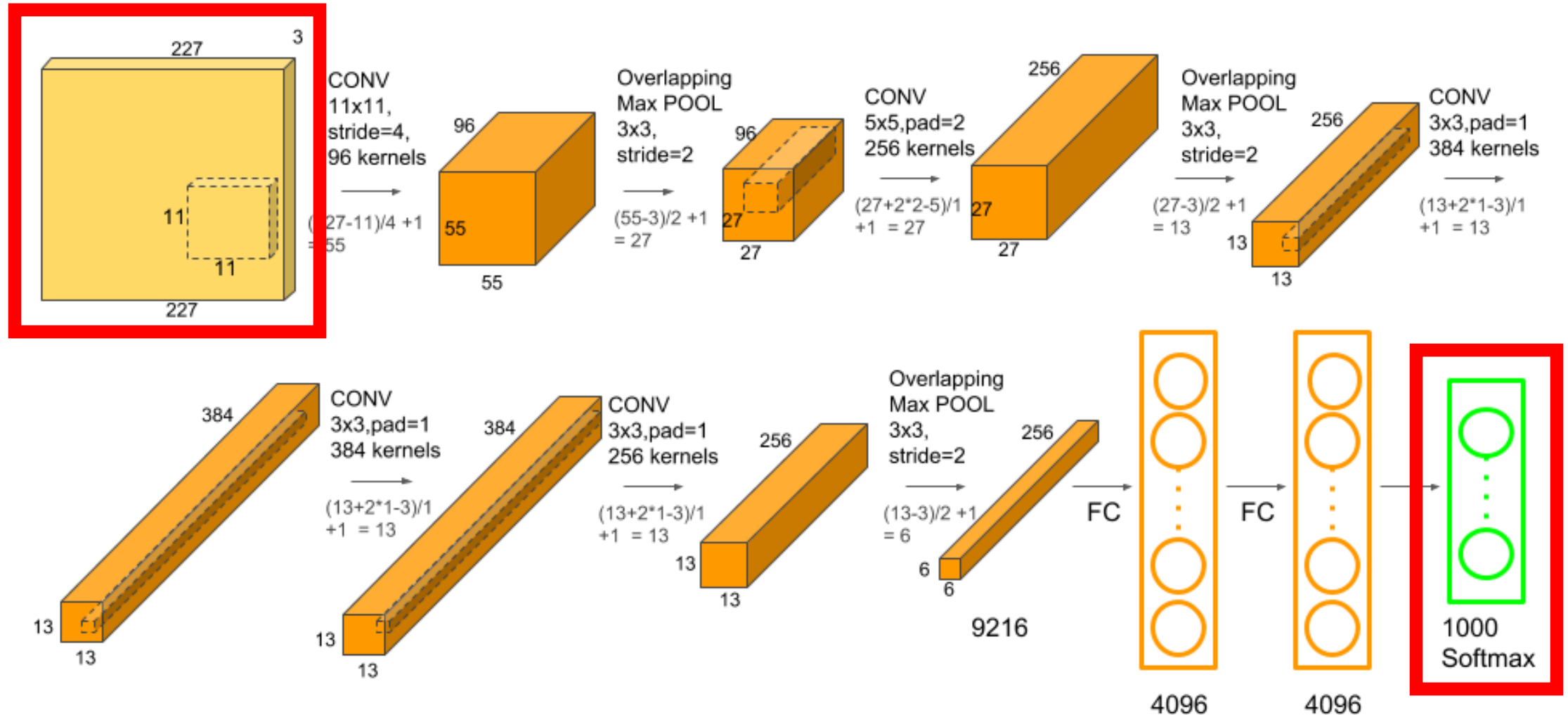
Input Preprocessing



Architecture

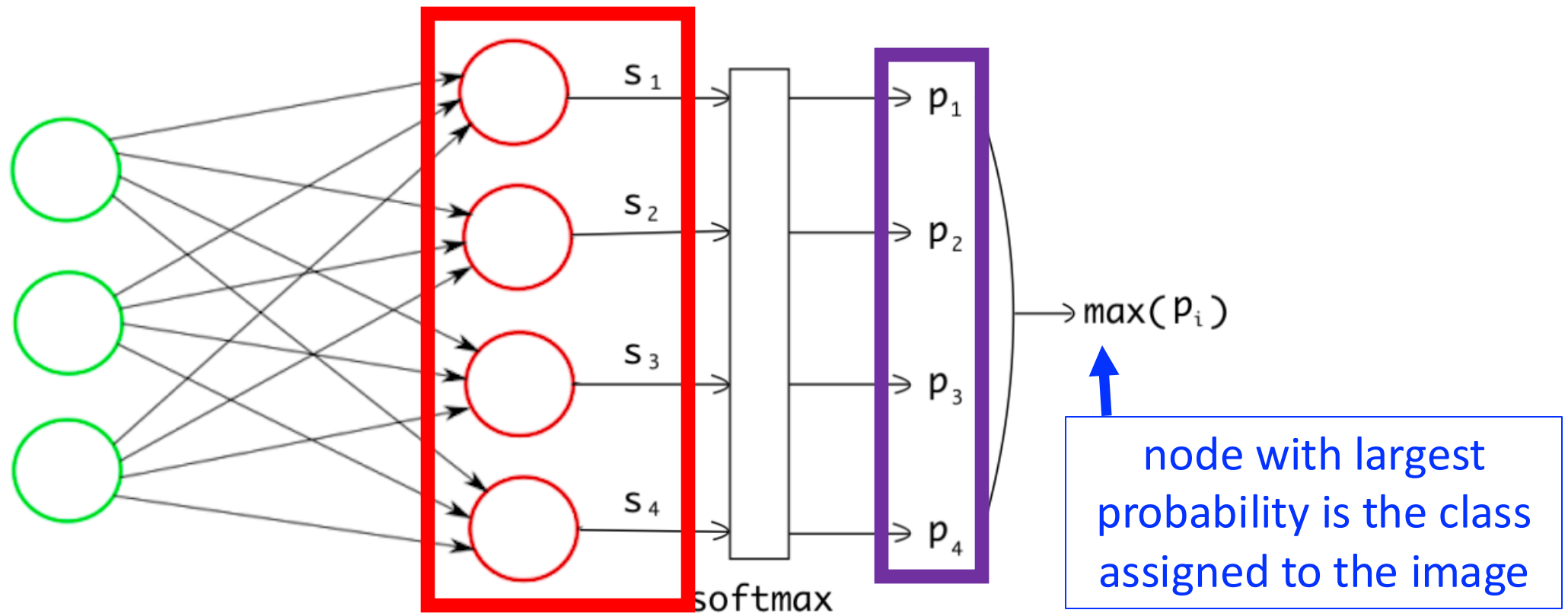
Input: RGB image resized to fixed input size

Output: 1000 class probabilities (sums to 1)



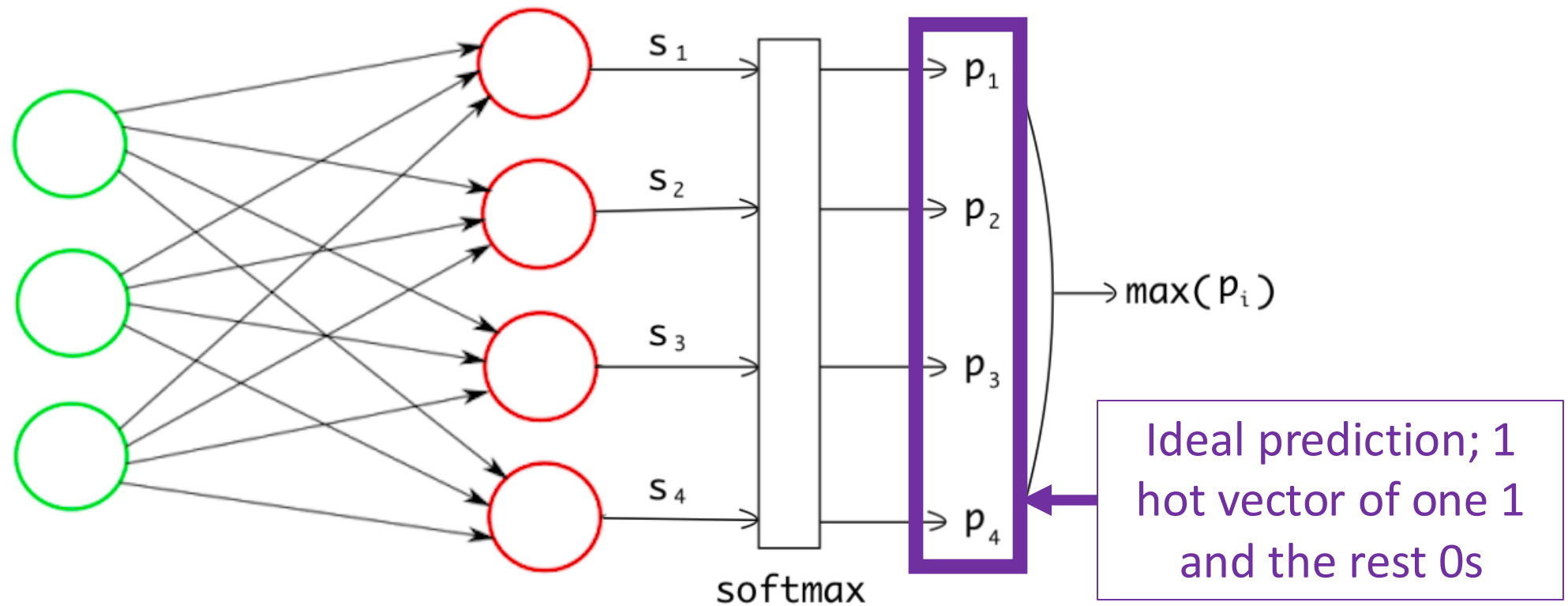
Architecture: Output Softmax Layer (Multiclass Classification)

Softmax: converts vector of **scores** into a **probability distribution** that sums to 1; e.g.,



Architecture: Output Softmax Layer (Multiclass Classification)

Softmax: converts vector of **scores** into a **probability distribution** that sums to 1; e.g.,



Architecture: Output Softmax Layer (Multiclass Classification)

Softmax: converts vector of **scores** into a probability distribution that sums to 1

Remove negative values while preserving original order of scores using e (negative scores become slightly larger than 0 and positive values grow exponentially; choosing e rather than another exponent base simplifies subsequent math)

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

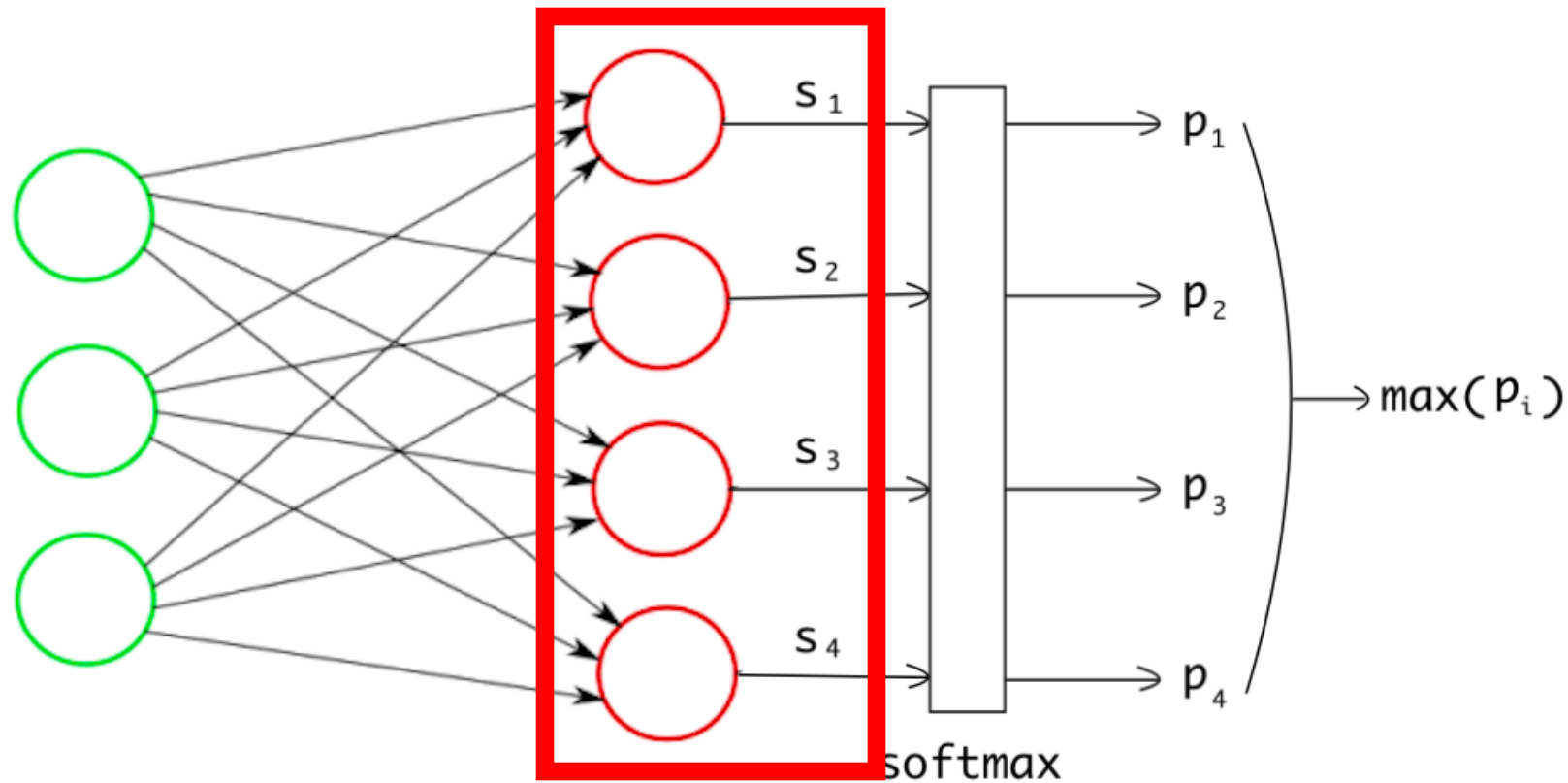
$i = 1, \dots, K$

Number of classes (i.e., 1000)

Want to divide each node's score by sum of all entries to make them sum to 1 (normalization)

Architecture: Output Softmax Layer (Multiclass Classification)

Softmax: converts vector of **scores** into a probability distribution that sums to 1; e.g.,



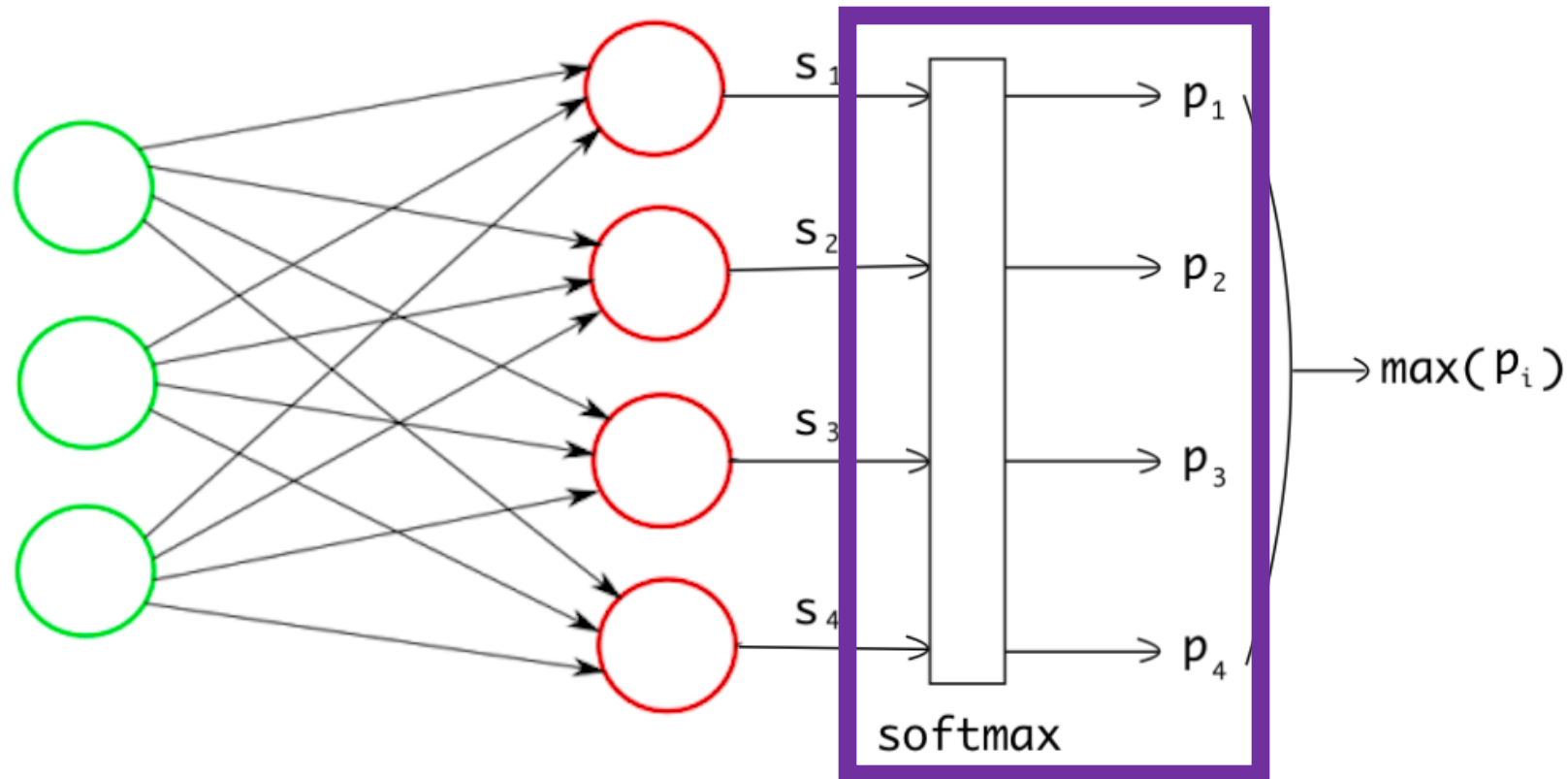
Architecture: Output Softmax Layer (Multiclass Classification)

Softmax: converts vector of **scores** into a probability distribution that sums to 1; e.g.,

	Scoring Function
Dog	-3.44
Cat	1.16
Boat	-0.81
Airplane	3.91

Architecture: Output Softmax Layer (Multiclass Classification)

Softmax: converts vector of scores into a **probability distribution** that sums to 1; e.g.,



Architecture: Output Softmax Layer (Multiclass Classification)

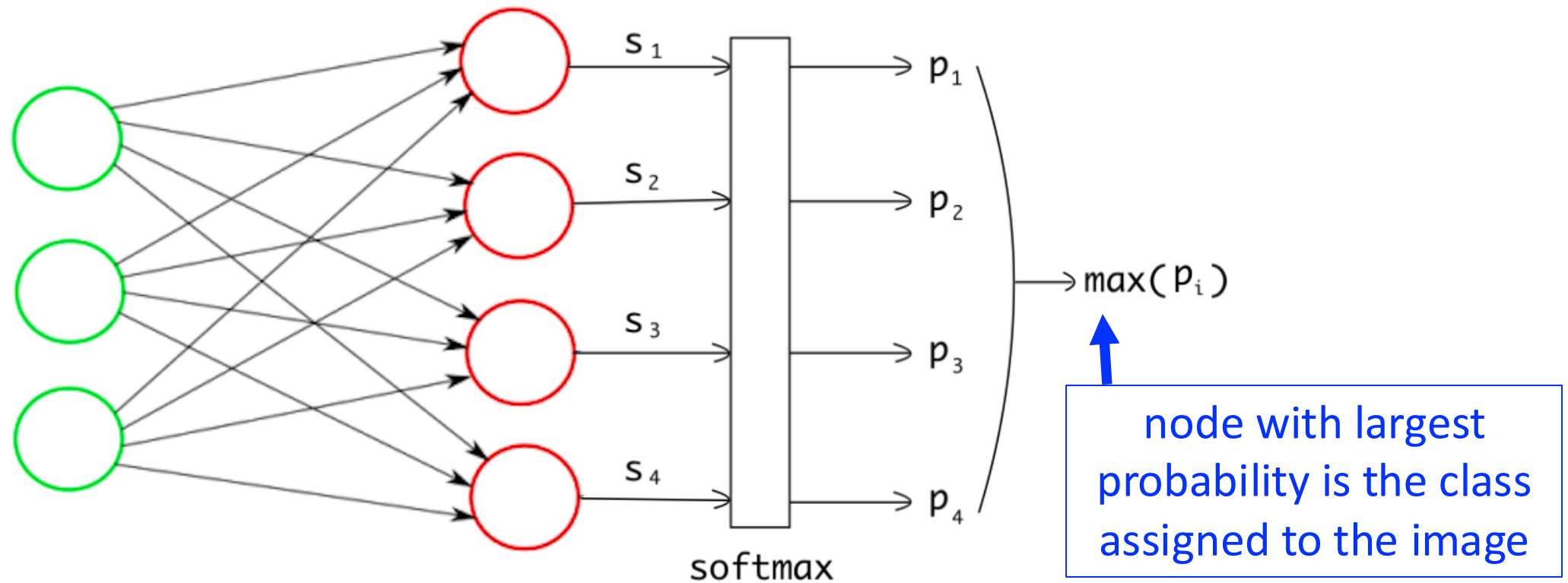
Softmax: converts vector of scores into a **probability distribution** that sums to 1; e.g.,

$$e^{z_i} \quad \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

	Scoring Function	Unnormalized Probabilities	Normalized Probabilities
Dog	-3.44	0.0321	0.0006
Cat	1.16	3.1899	0.0596
Boat	-0.81	0.4449	0.0083
Airplane	3.91	49.8990	0.9315

Architecture: Output Softmax Layer (Multiclass Classification)

Softmax: converts vector of **scores** into a **probability distribution** that sums to 1; e.g.,

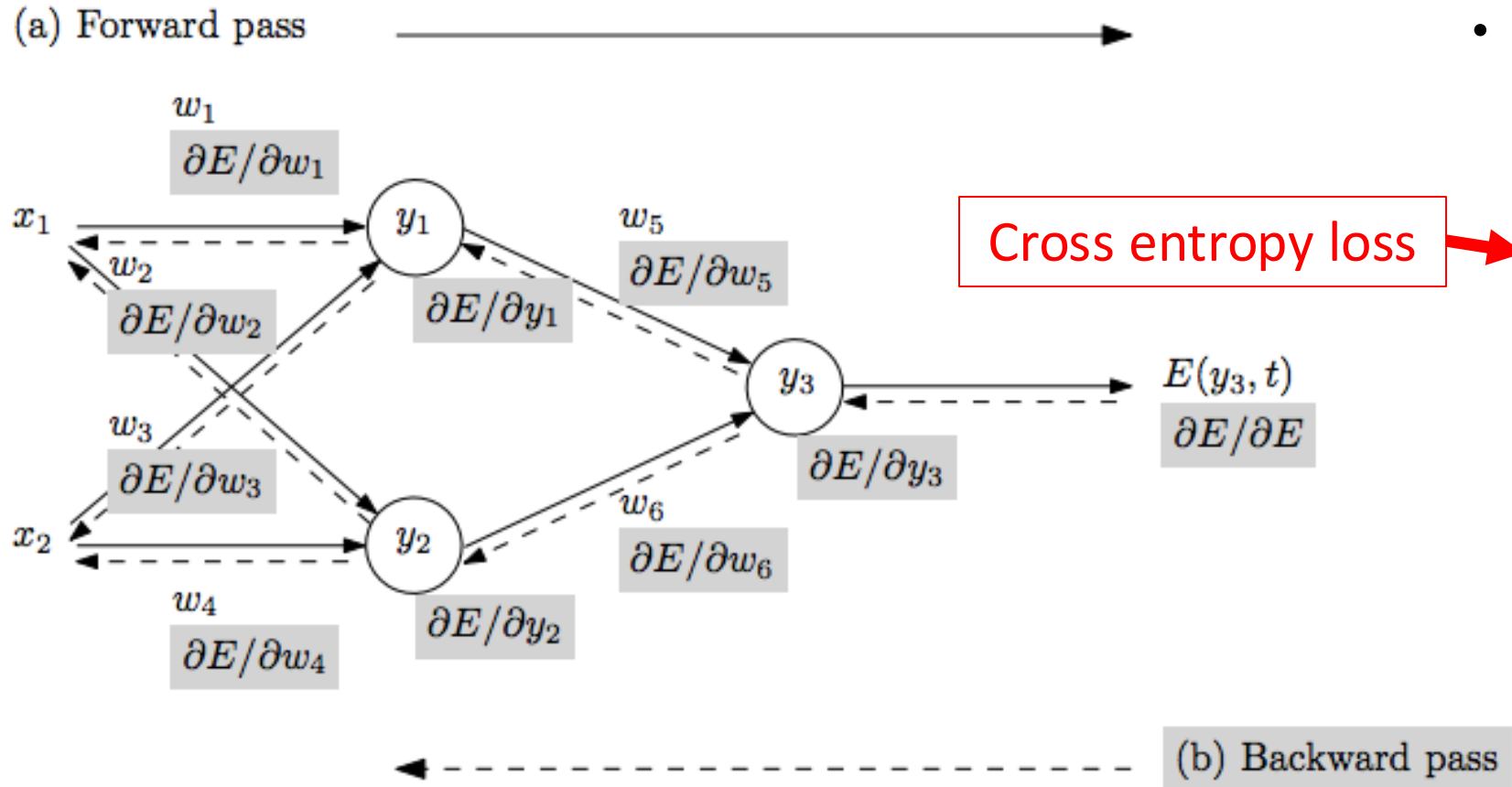


Architecture: Output Softmax Layer (Multiclass Classification)

Softmax: converts vector of scores into a probability distribution that sums to 1; e.g.,

	Scoring Function	Unnormalized Probabilities	Normalized Probabilities
Dog	-3.44	0.0321	0.0006
Cat	1.16	3.1899	0.0596
Boat	-0.81	0.4449	0.0083
Airplane	3.91	49.8990	0.9315

Algorithm Training: Recall How NNs Learn

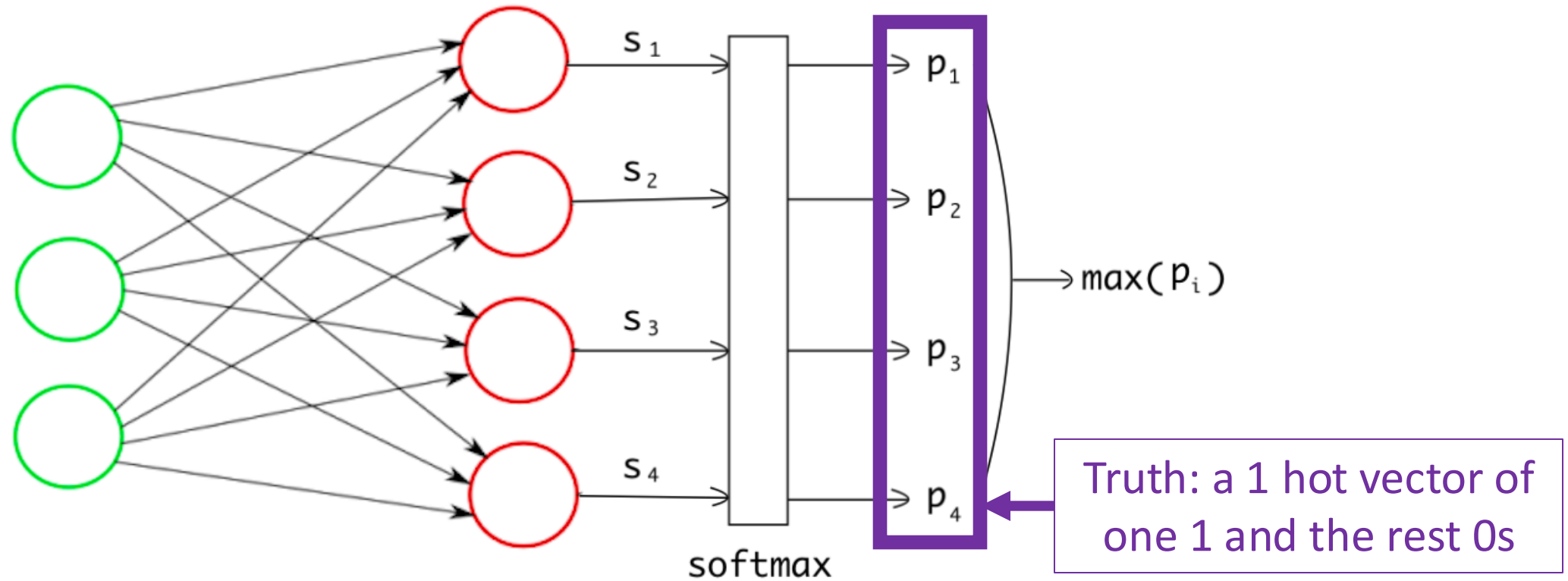


- Repeat until stopping criterion met:

1. **Forward pass:** propagate training data through model to make predictions
2. **Error quantification:** measure dissatisfaction with a model's predictions on training data
3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
4. Update each parameter using calculated gradients

Algorithm Training: Measure Cross Entropy Loss

Measure distance between predicted and true class distribution for each example



Algorithm Training: Measure Cross Entropy Loss

Probability distribution of predicted class

Probability distribution of true class

Number of classes?

Recall, truth is set to 1 for one class and 0 otherwise

Observed features

Simplifies to the log of the predicted probability for the correct class (i.e., negative log likelihood loss)

$$\begin{aligned} L_{\text{CE}}(\hat{y}, y) &= - \sum_{k=1}^K y_k \log \hat{y}_k \\ &= - \sum_{k=1}^K y_k \log \hat{p}(y = k | x) \\ &= - \log \hat{y}_k, \quad (\text{where } k \text{ is the correct class}) \\ &= - \log \frac{\exp(w_k \cdot x + b_k)}{\sum_{j=1}^K \exp(w_j \cdot x + b_j)} \end{aligned}$$

Algorithm Training: Measure Cross Entropy Loss

Probability distribution of predicted class

Probability distribution of true class

$$\begin{aligned} L_{\text{CE}}(\hat{y}, y) &= - \sum_{k=1}^K y_k \log \hat{y}_k \\ &= - \sum_{k=1}^K y_k \log \hat{p}(y = k | x) \\ &= - \log \hat{y}_k, \quad (\text{where } k \text{ is the correct class}) \\ &= - \log \frac{\exp(w_k \cdot x + b_k)}{\sum_{j=1}^K \exp(w_j \cdot x + b_j)} \end{aligned}$$

Number of classes?

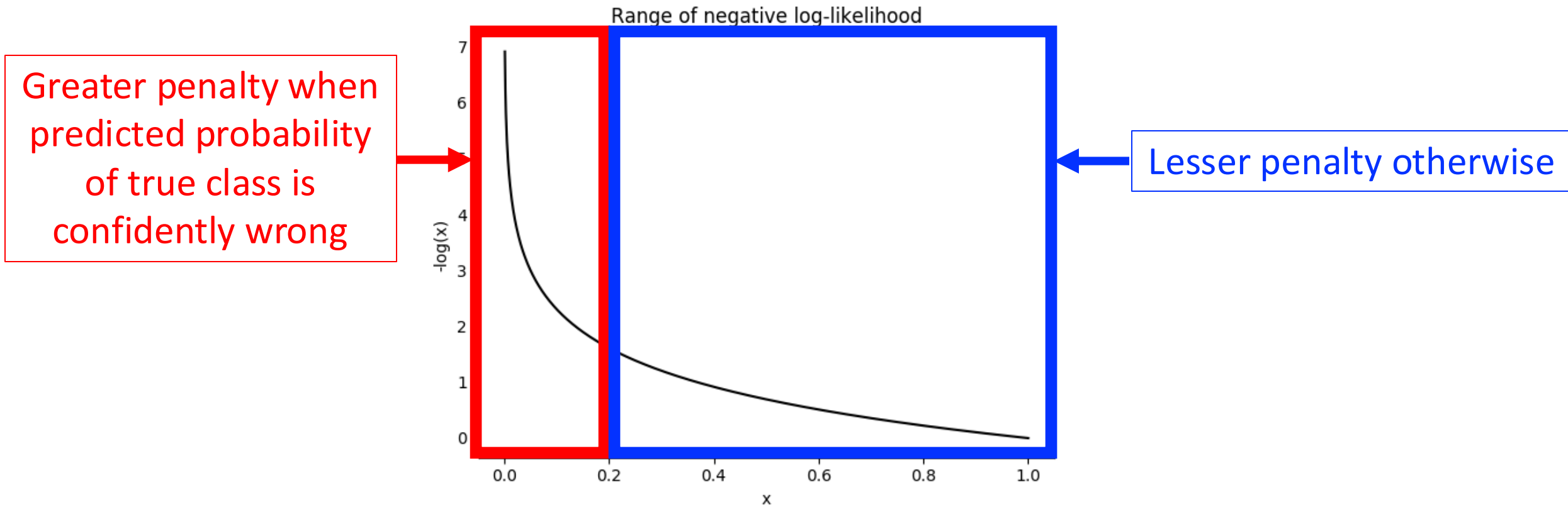
Recall, truth is set to 1 for one class and 0 otherwise

Observed features

Range of possible values:

- Minimum: 0 (negative log of 1)
- Maximum: Infinity (negative log of 0)

Algorithm Training: Measure Cross Entropy Loss



What is the range of possible values?

- Minimum: 0 (negative log of 1)
- Maximum: Infinity (negative log of 0)

$$= -\log \frac{\exp(w_k \cdot x + b_k)}{\sum_{j=1}^K \exp(w_j \cdot x + b_j)}$$

Algorithm Training: Measure Cross Entropy Loss

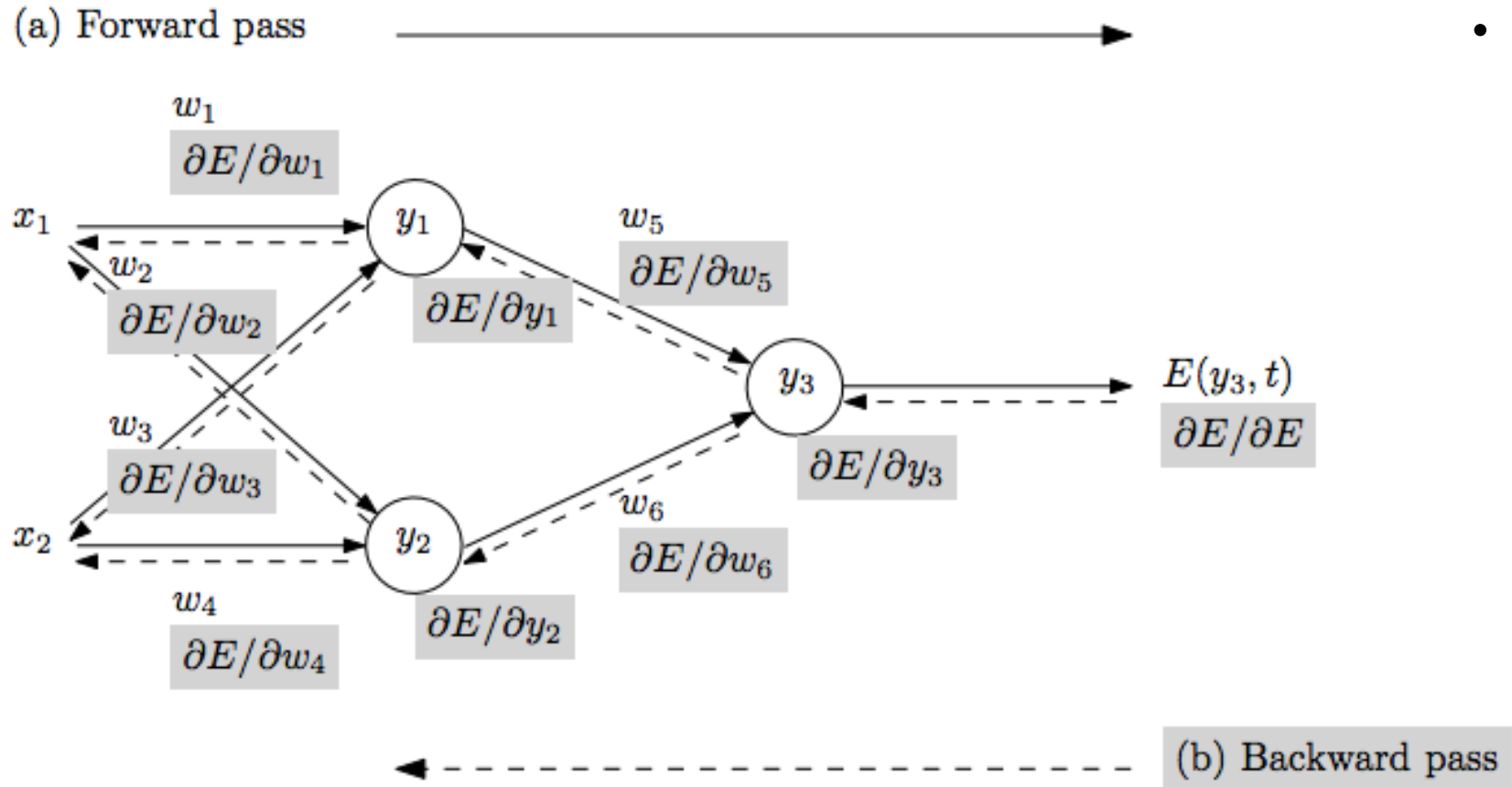
e.g., What would be the loss for this example if the true class label is cat?

$$= -\log \hat{y}_k, \quad (\text{where } k \text{ is the correct class})$$

$$= -\log(0.0596) = 2.82$$

	Scoring Function	Unnormalized Probabilities	Normalized Probabilities
Dog	-3.44	0.0321	0.0006
Cat	1.16	3.1899	0.0596
Boat	-0.81	0.4449	0.0083
Airplane	3.91	49.8990	0.9315

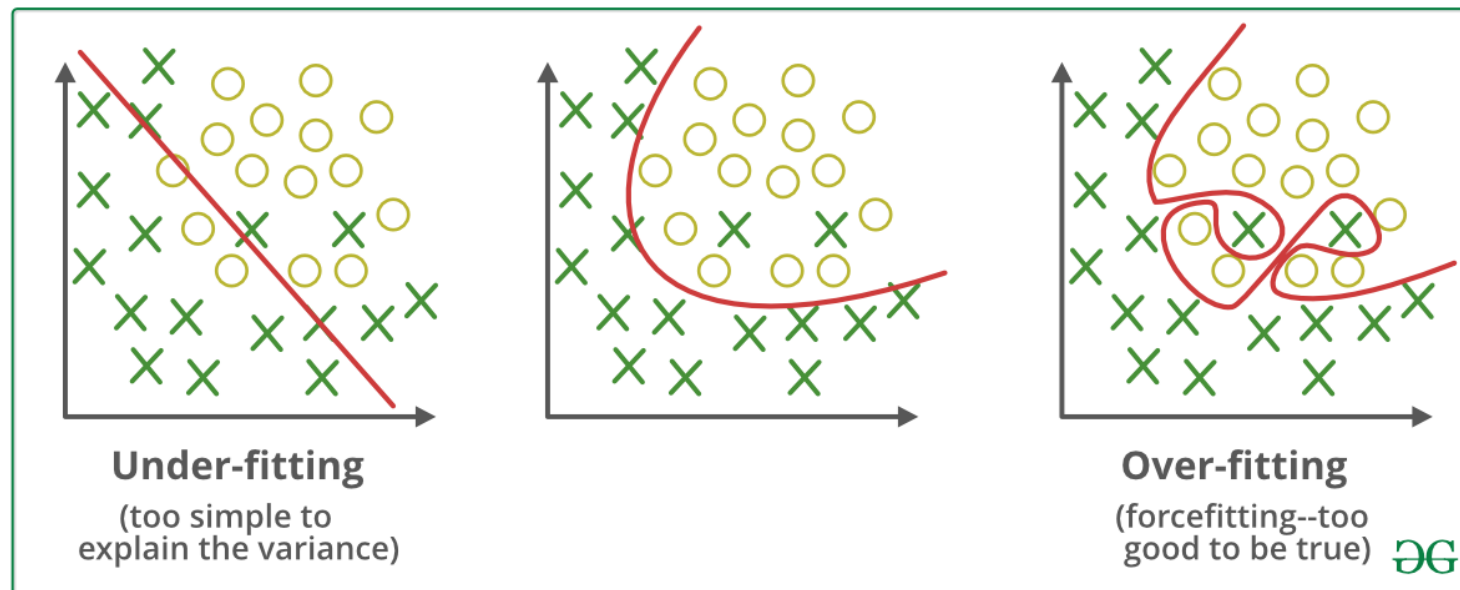
Algorithm Training



- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make predictions
 2. **Error quantification:** measure dissatisfaction with a model's predictions on training data
 3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
 4. Update each parameter using calculated gradients

Algorithm Training: Challenge Is Overfitting

- Overfitting is risk for models with larger representational capacity (i.e., # of parameters); AlexNet has 60 million parameters!

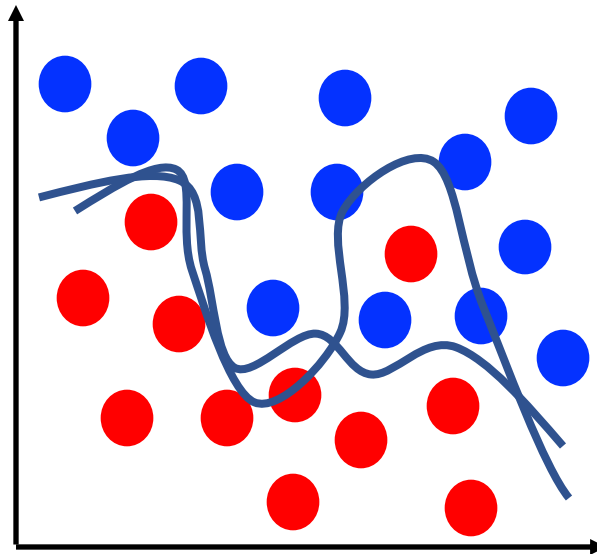


- Model learns to model **noise!**

AlexNet Remedies for Overfitting

- Overfitting is risk for models with larger representational capacity (i.e., # of parameters); AlexNet has 60 million parameters!
 1. Data augmentation: add more training data; e.g., intuitively,

Adding training data



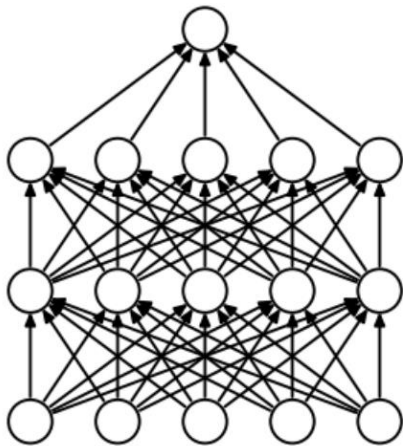
AlexNet Remedies for Overfitting

- Overfitting is risk for models with larger representational capacity (i.e., # of parameters); AlexNet has 60 million parameters!
 1. Data augmentation
 1. Random patches and their mirror images (2048x more data)
 2. Adjust RGB channels (using PCA to add multiples of principal components)

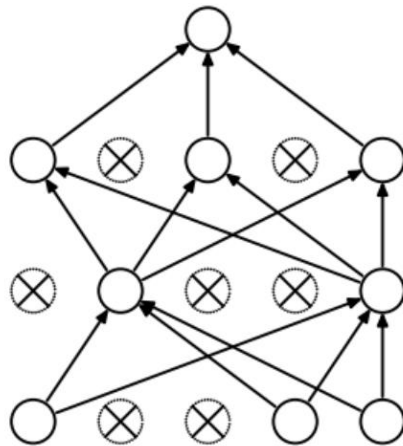


AlexNet Remedies for Overfitting

- Overfitting is risk for models with larger representational capacity (i.e., # of parameters); AlexNet has 60 million parameters!
 1. Data augmentation
 1. Random patches and their mirror images (2048x more data)
 2. Adjust RGB channels (using PCA to add multiples of principal components)
 2. Dropout (50% of nodes for first two fully connected layers); mimics ensembles by learning to solve same problem with different subnetworks



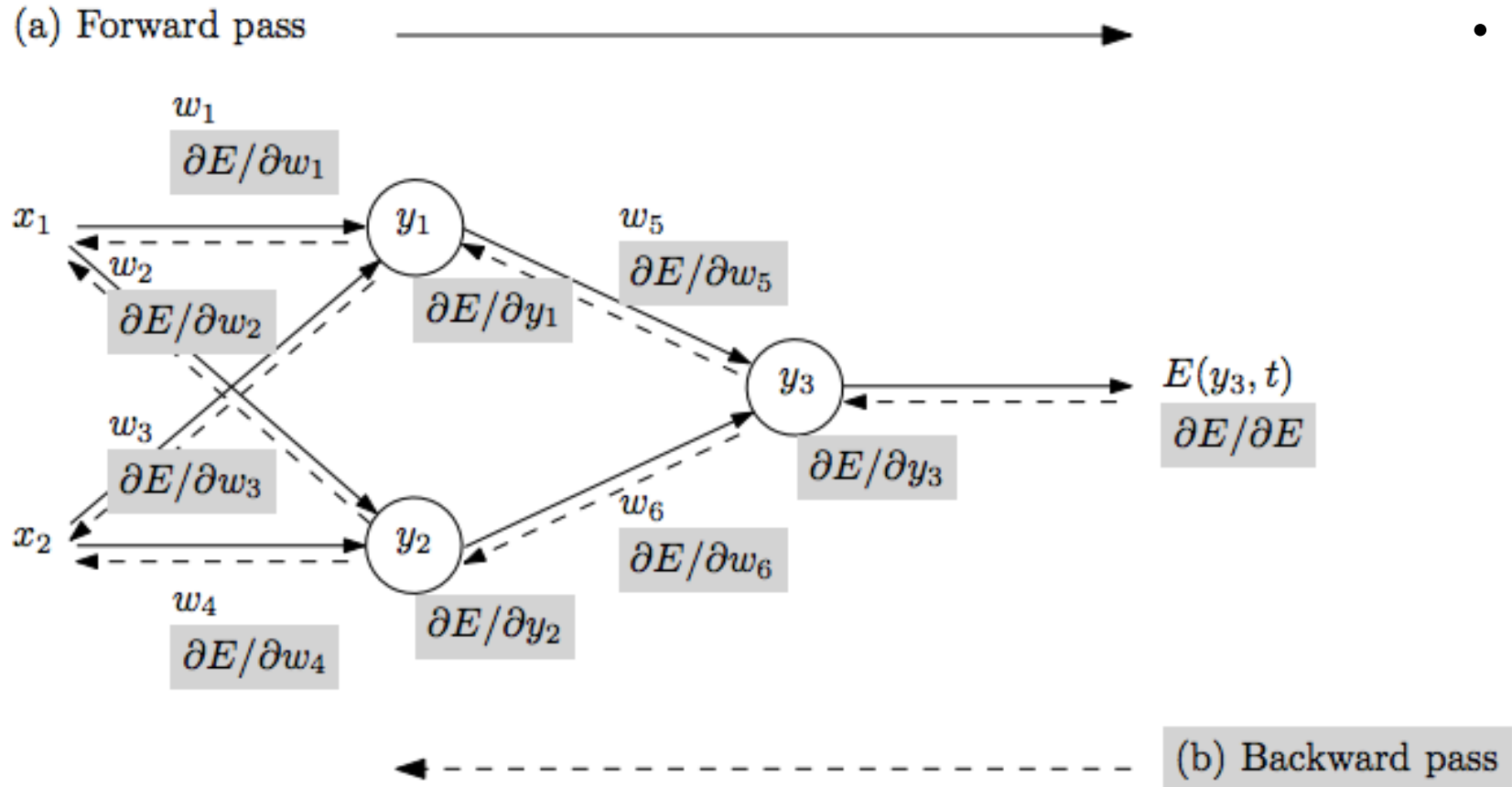
(a) Standard Neural Net



(b) After applying dropout.

Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. "Dropout: a simple way to prevent neural networks from overfitting." JMLR, 2014.

Algorithm Training: 90 Epochs on ImageNet



- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make predictions
 2. **Error quantification:** measure dissatisfaction with a model's predictions on training data
 3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
 4. Update each parameter using calculated gradients

AlexNet Analysis

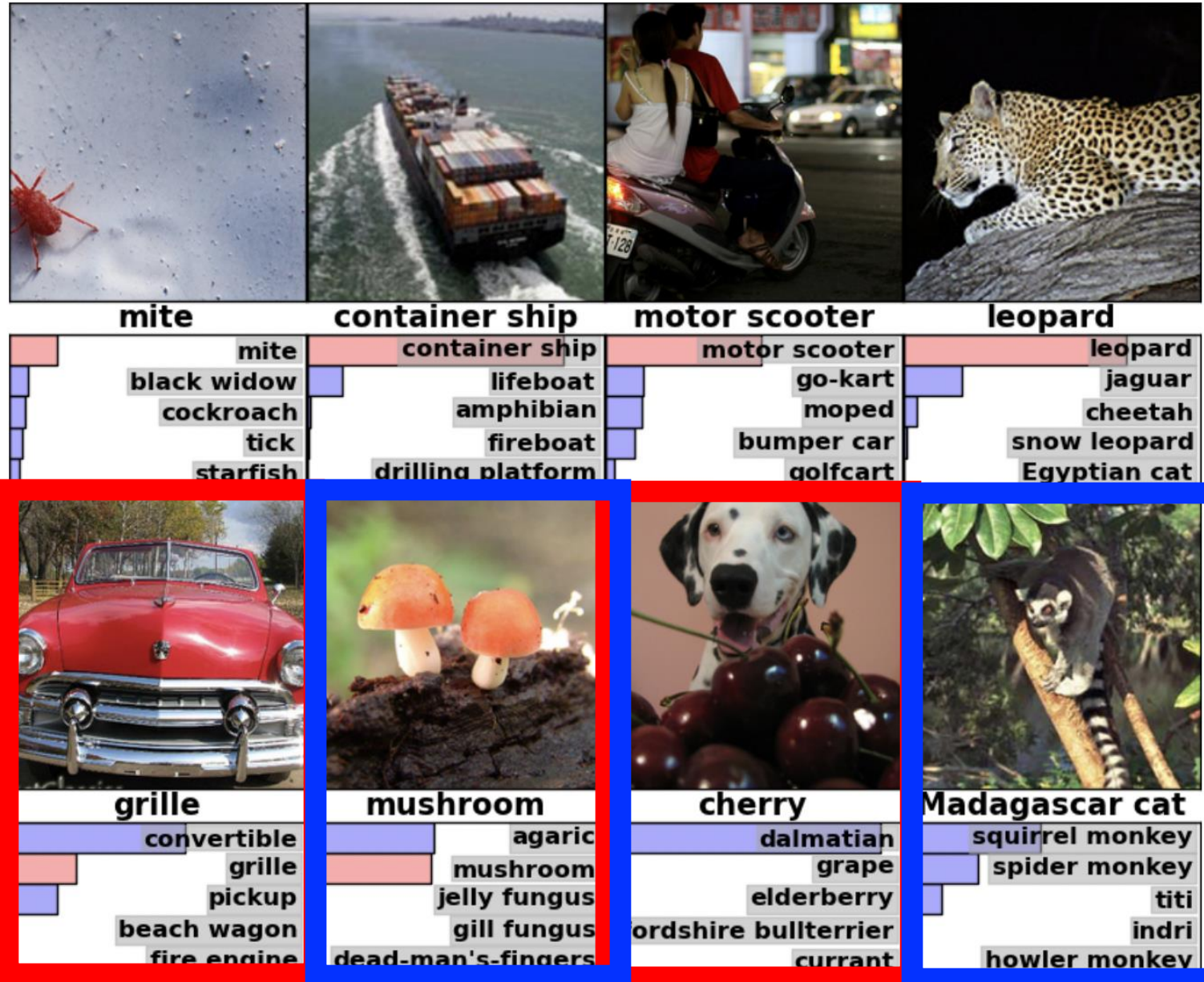
8 examples of predictions, correct and incorrect

When/why might the model succeed?

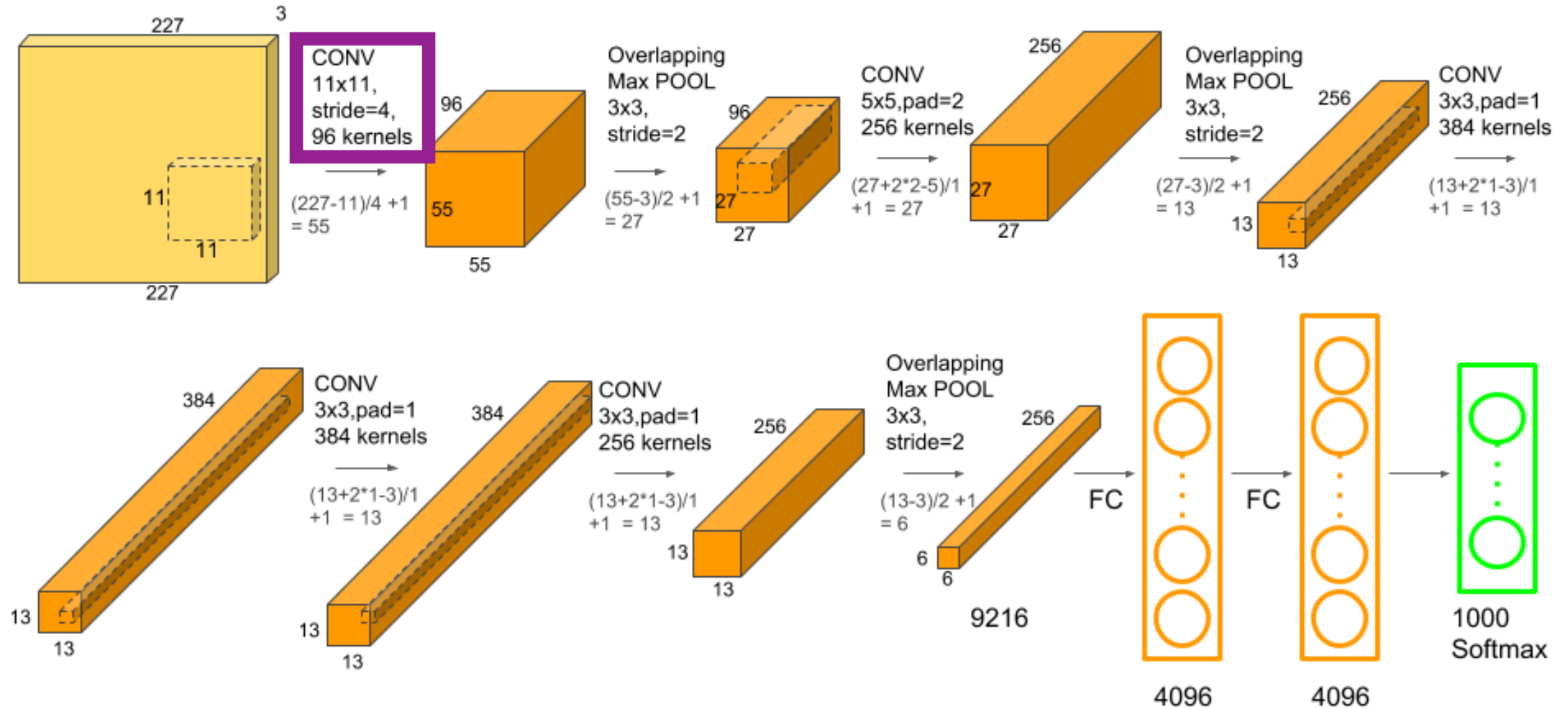
- Single well-defined object (even if off-centered)

When/why might the model fail?

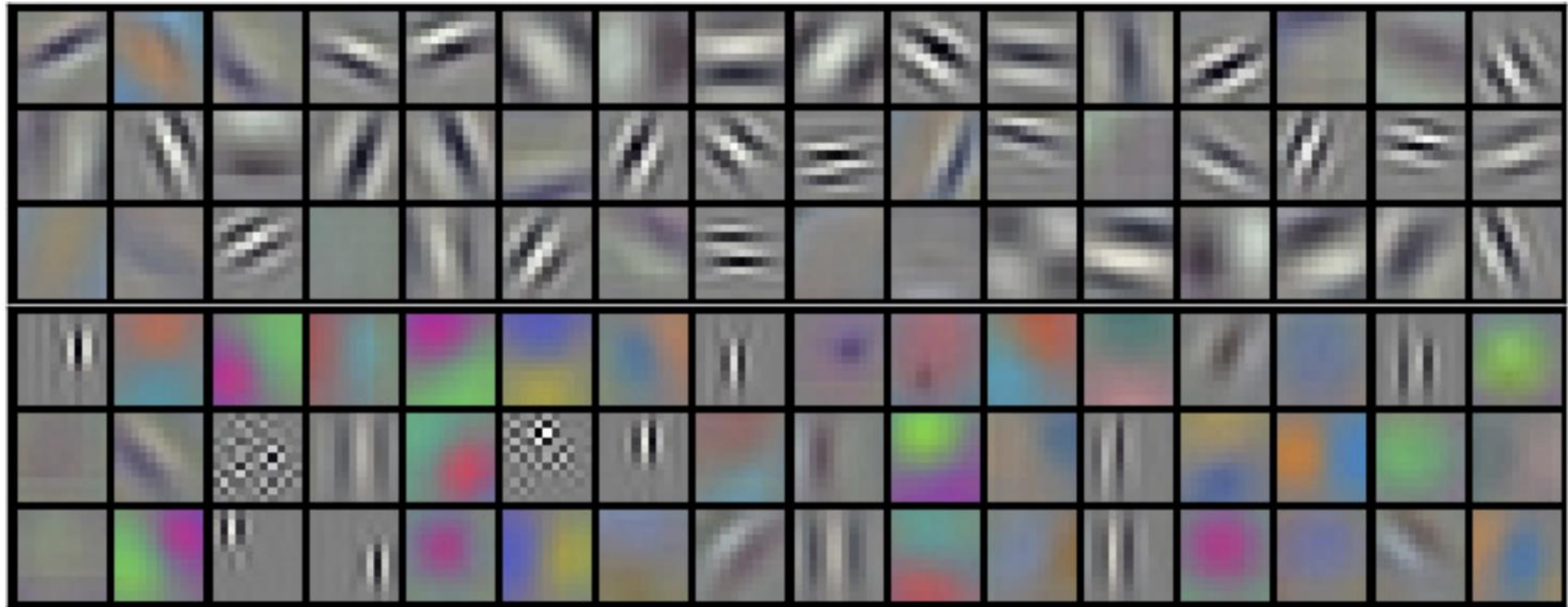
- **Ambiguity**
- **Similar/superset categories**



AlexNet: Inspecting What It Learned



AlexNet: Inspecting What It Learned (96 Filters)



Model learned filters that select based on frequency, orientation, and color!
(Aligns with Hubel & Weisel's findings for how vision systems work)

AlexNet: Key Tricks for Going Deeper

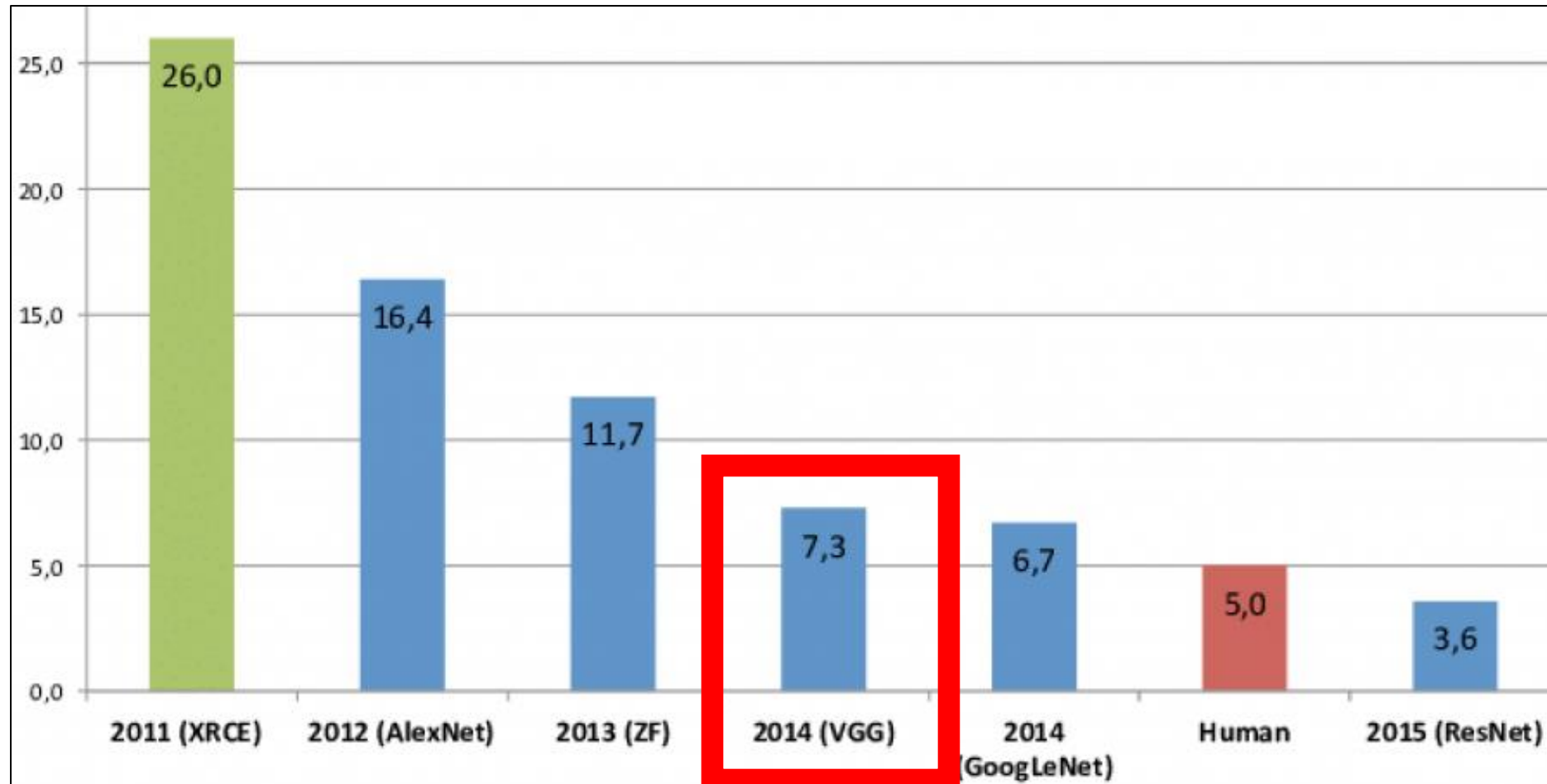
- ReLU instead of sigmoid or tanh activation functions
- Regularization techniques; e.g.,
 1. Data augmentation
 2. Dropout in fully connected layers
- Trained across two GPUs

Object Recognition: Today's Topics

- ImageNet Challenge Top Performers
- Baseline Model: AlexNet
- **VGG**
- ResNet
- Summary of CNN Era

VGG: A Deeper CNN

Progress of models on ImageNet (Top 5 Error)



Why VGGNet?

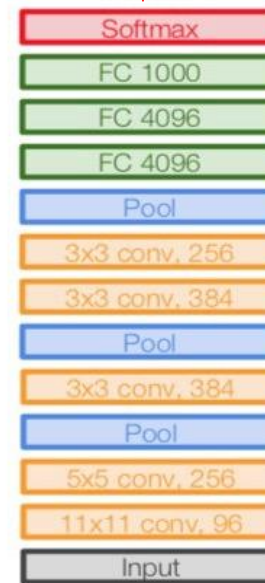
VGG stands for the **Visual Geometry Group (VGG)** at University of Oxford where the authors were based 😊

Karen Simonyan and Andrew Zisserman. "Very deep convolutional networks for large-scale image recognition." International Conference on Learning Representations (ICLR), 2015.

Key Novelty: Deeper Does Better

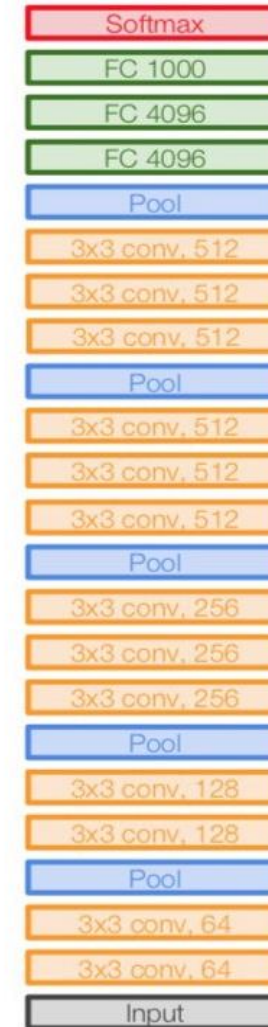
** Number of layers with learnable model parameters between input and output layer (i.e., excludes pooling layers)*

8 layers



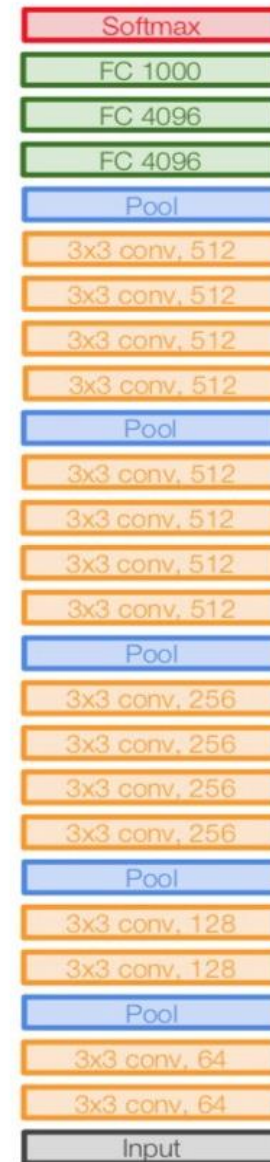
AlexNet

16 layers



VGG16

19 layers

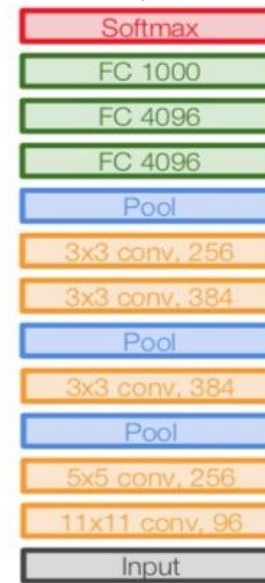


VGG19

Figure Source (edited to fix mistakes): <https://medium.com/deep-learning-g/cnn-architectures-vggnet-e09d7fe79c45>

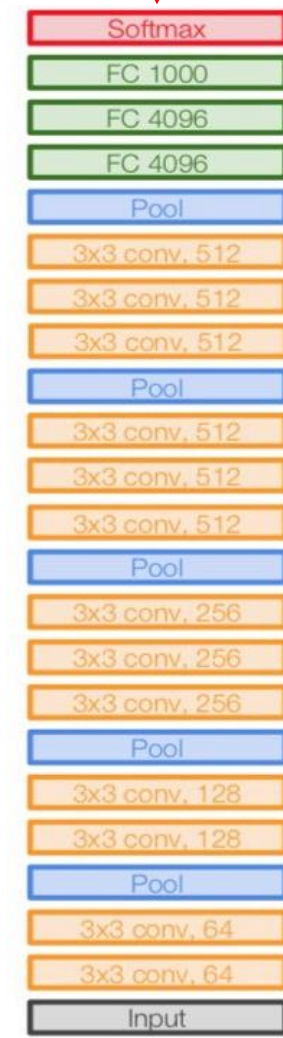
Key Novelty: Deeper Does Better

15.3% top-5 error



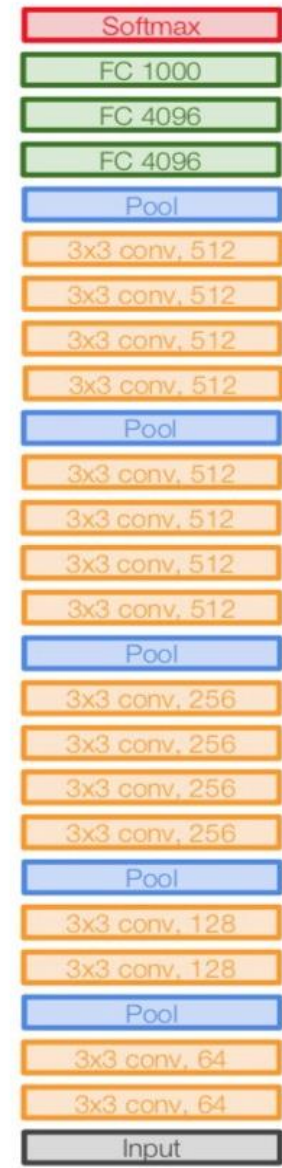
AlexNet

7.7% top-5 error



VGG16

7.3% top-5 error



VGG19

Figure Source (edited to fix mistakes): <https://medium.com/deep-learning-g/cnn-architectures-vggnet-e09d7fe79c45>

Key Novelty: Deeper Does Better

* Number of layers with learnable model parameters between input and output layer (i.e., exclude pooling layers)

Layers with differences

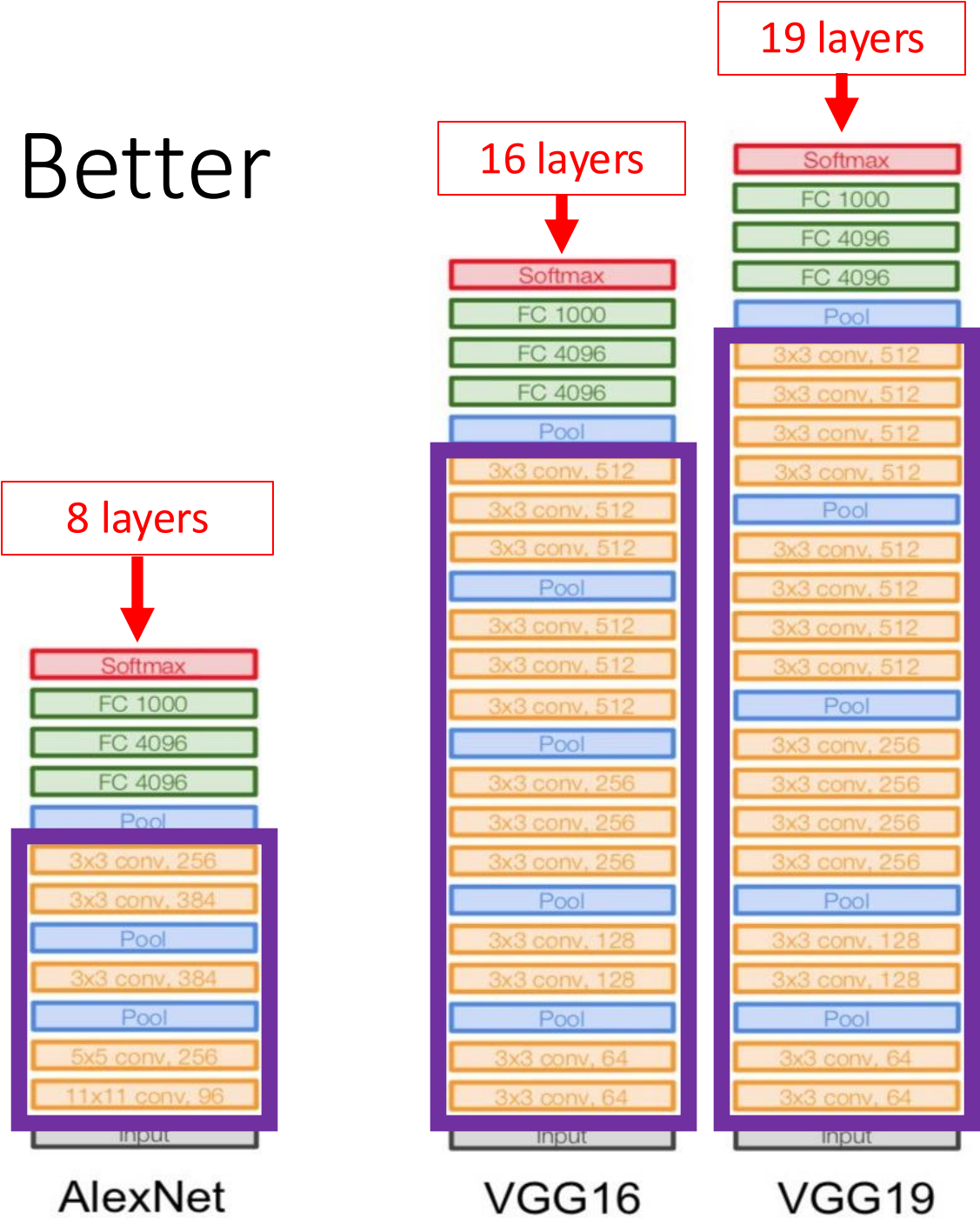


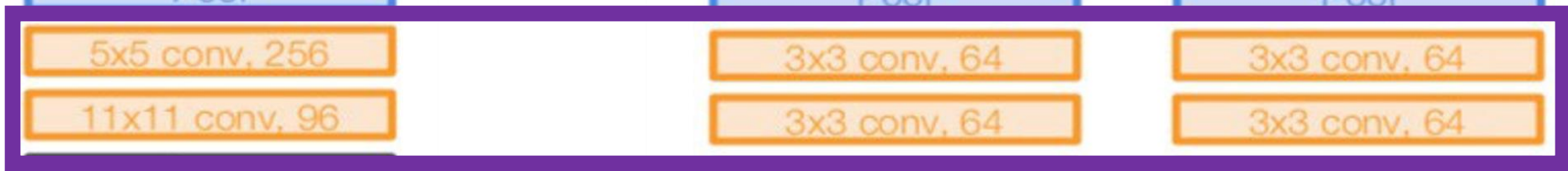
Figure Source (edited to fix mistakes): <https://medium.com/deep-learning-g/cnn-architectures-vggnet-e09d7fe79c45>

Key Idea: Smaller Convolutional Filters

- Replace larger filter with stack of smaller filters



11x11 filter, 96 filters



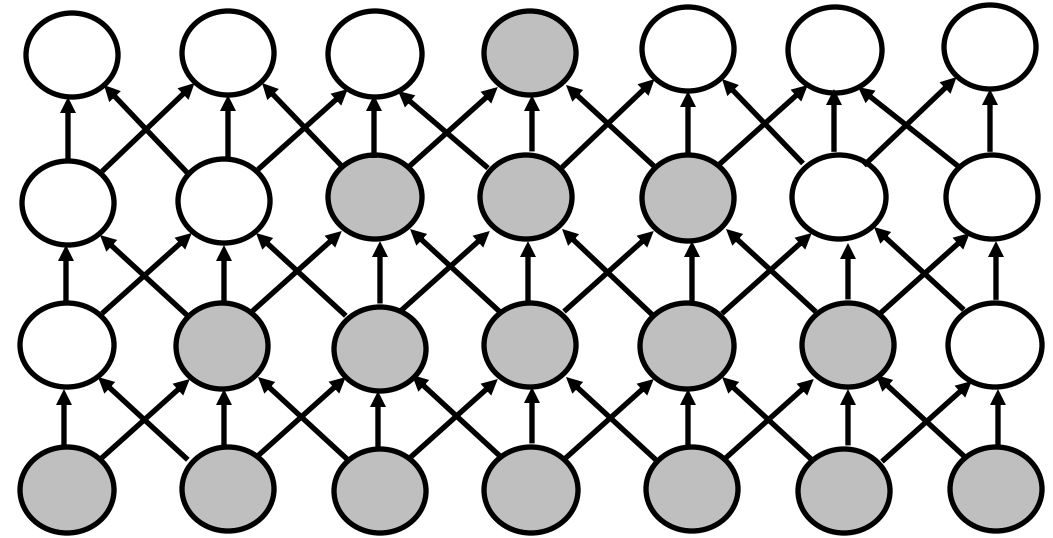
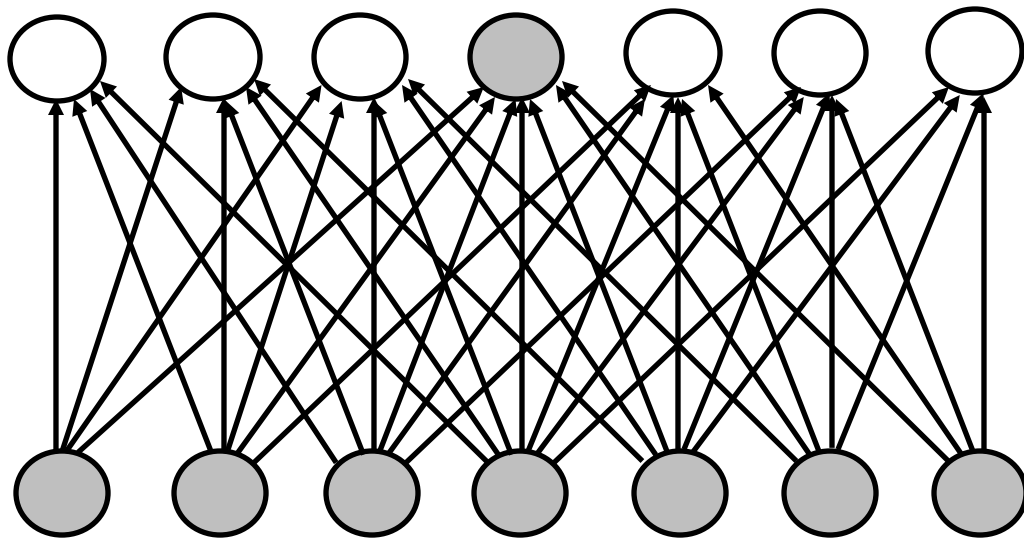
AlexNet

VGG16

VGG19

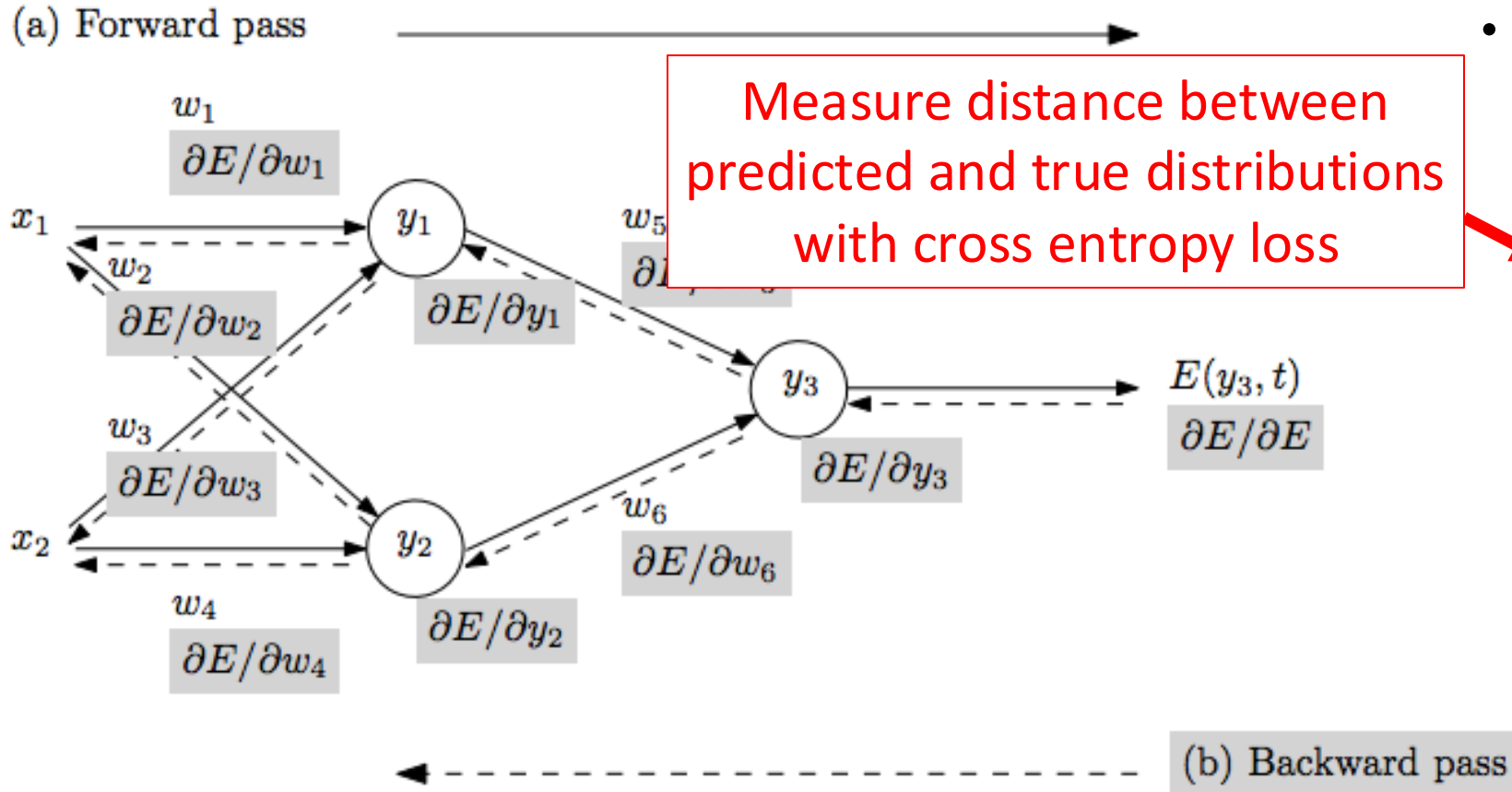
Key Idea: Smaller Convolutional Filters

- Replace larger filter with stack of smaller filters; e.g., replace 7x7 with three 3x3s



- Benefits:
 - More discriminative classifier since more non-linear rectifications: 3 vs 1
 - Reduces # of parameters: multiple of 27 (3×3^2) parameters vs 49 (7×7) parameters

VGG Training (follows AlexNet): 74 Epochs



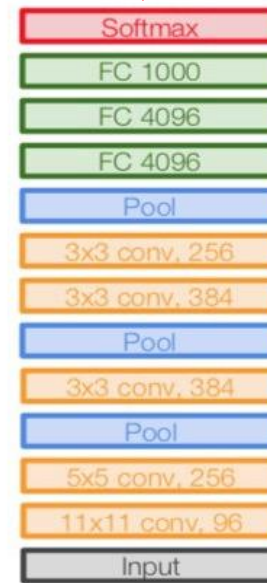
- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make predictions
 2. **Error quantification:** measure dissatisfaction with a model's predictions on training data
 3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
 4. Update each parameter using calculated gradients

Algorithm Training (follows AlexNet)

- Strategies to mitigate overfitting
 1. Data augmentation
 1. Random patches and their mirror images (2048x more data)
 2. Adjust RGB channels (using PCA to add multiples of principal components)
 2. Dropout (50% of nodes for first two fully connected layers); mimics ensembles by learning to solve same problem with different subnetworks

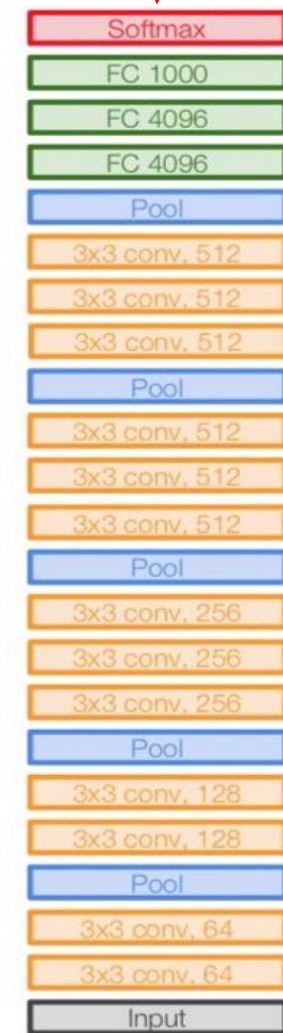
VGG Limitation: Models Are Large!

60 million parameters



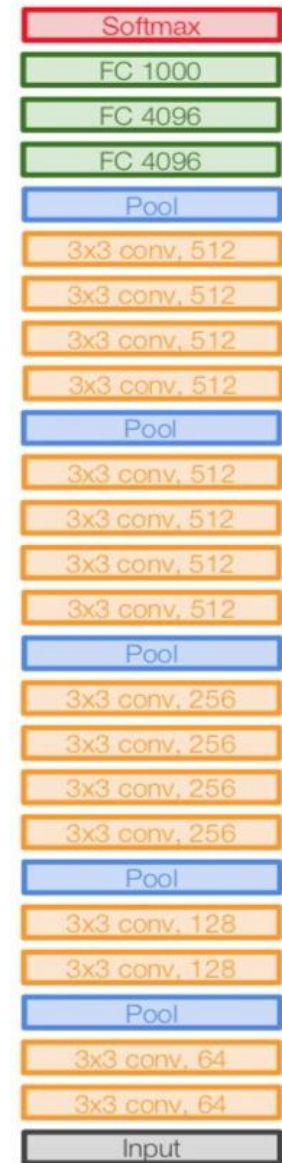
AlexNet

138 million parameters



VGG16

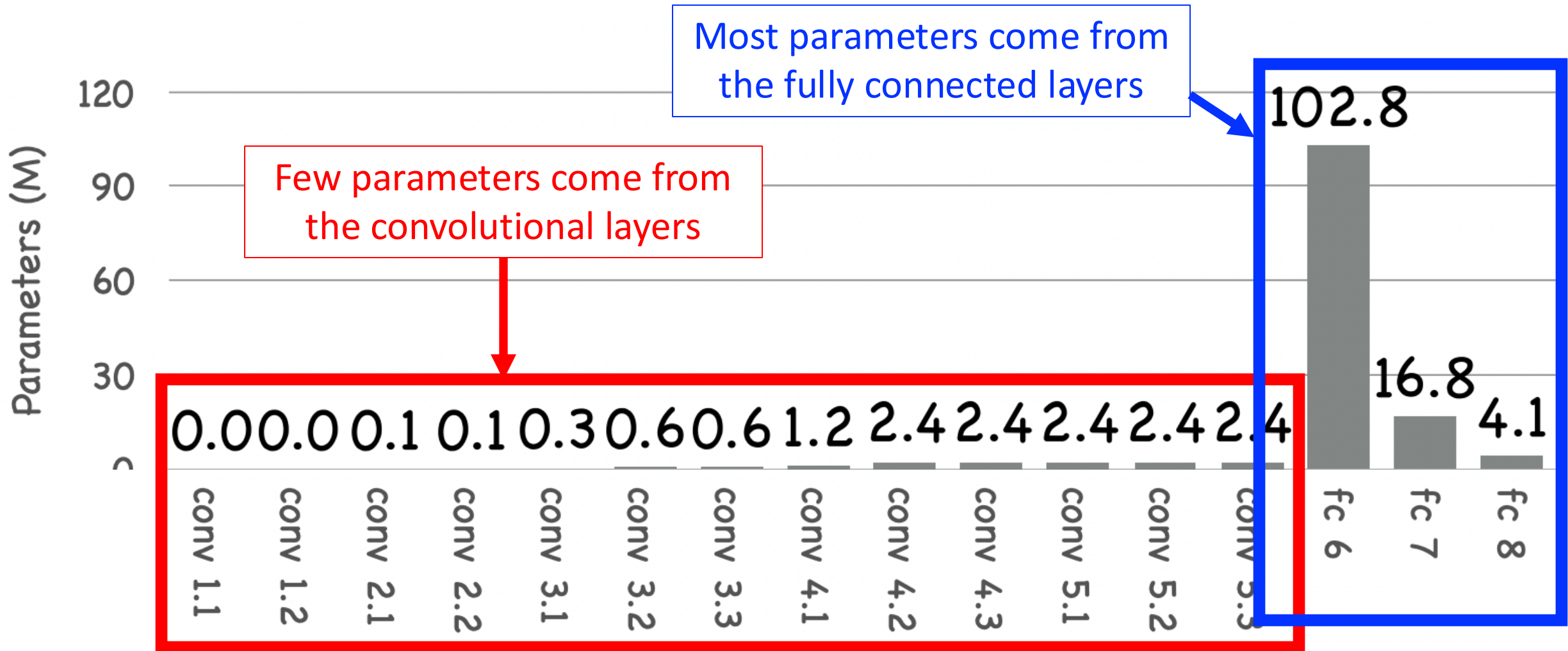
144 million parameters



VGG19

Figure Source (edited to fix mistakes): <https://medium.com/deep-learning-g/cnn-architectures-vggnet-e09d7fe79c45>

VGG Limitation: Models Are Large (e.g., VGG16)



VGG: Key Tricks for Going Deeper

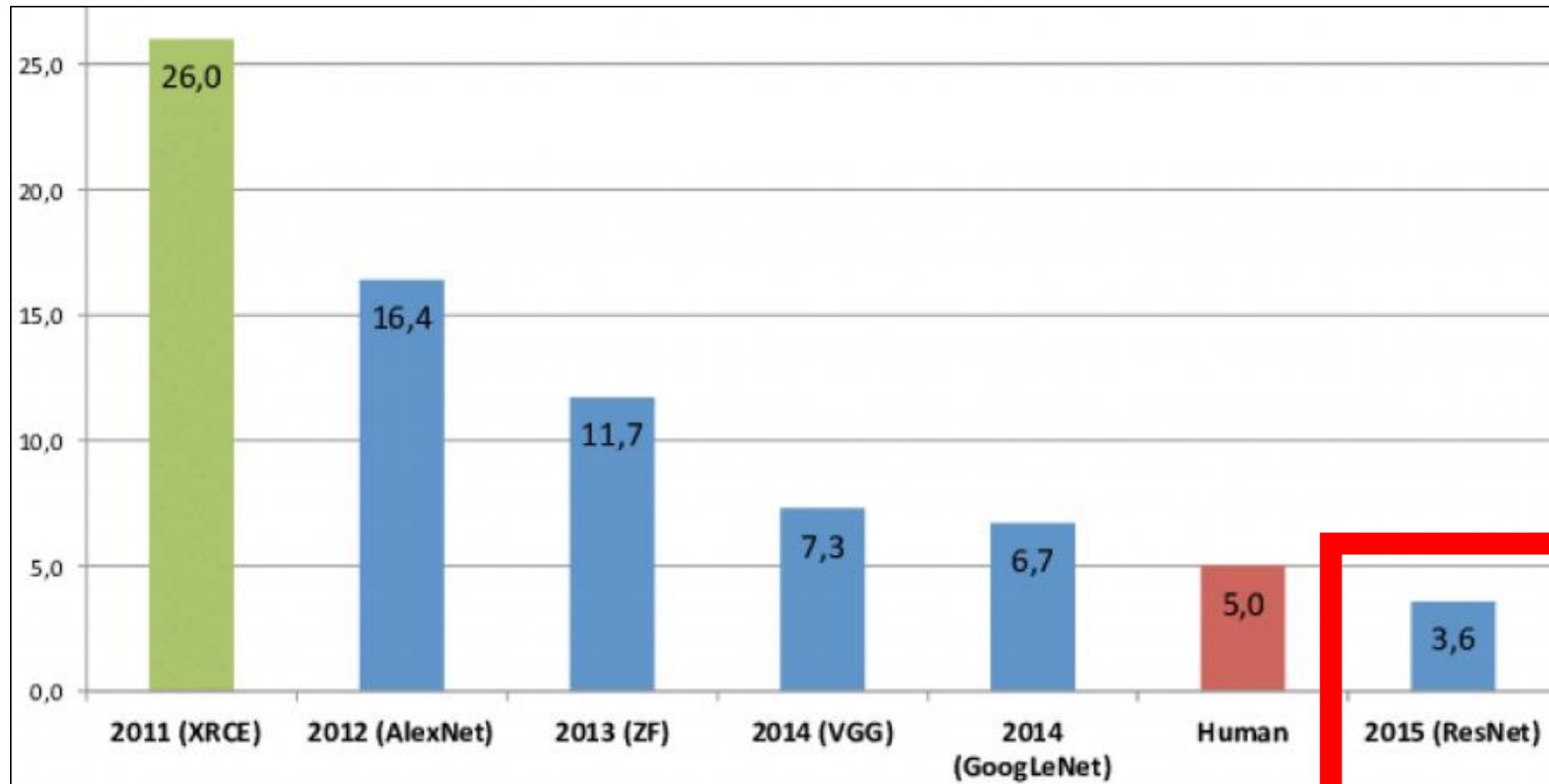
- 3x3 filters instead of larger filters
- Regularization techniques; e.g.,
 1. Data augmentation
 2. Dropout in fully connected layers
- Trained across multiple GPUs

Object Recognition: Today's Topics

- ImageNet Challenge Top Performers
- Baseline Model: AlexNet
- VGG
- **ResNet**
- Summary of CNN Era

Secret Sauce for State-of-Art: Deeper CNNs

Progress of models on ImageNet (Top 5 Error)



Why ResNet?

“Res” stands for residuals, which is the novel proposed idea.

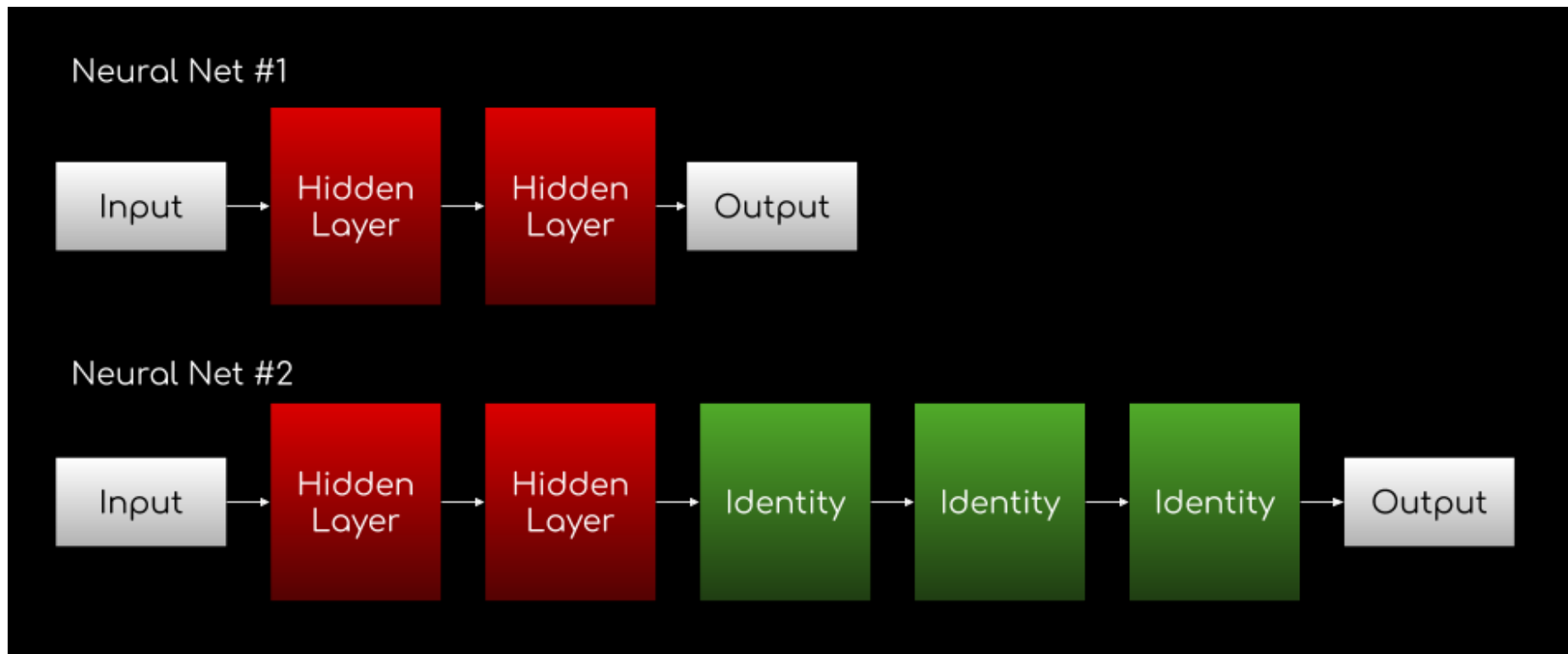
Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. “Deep Residual Learning for Image Recognition.” CVPR, 2016.

Motivating Observation

Idea: deeper networks should perform comparable or better than shallower networks since they can learn the shallower function by simply learning “identity” functions for later layers

Observation: adding more layers leads to WORSE results!

Is the problem overfitting?



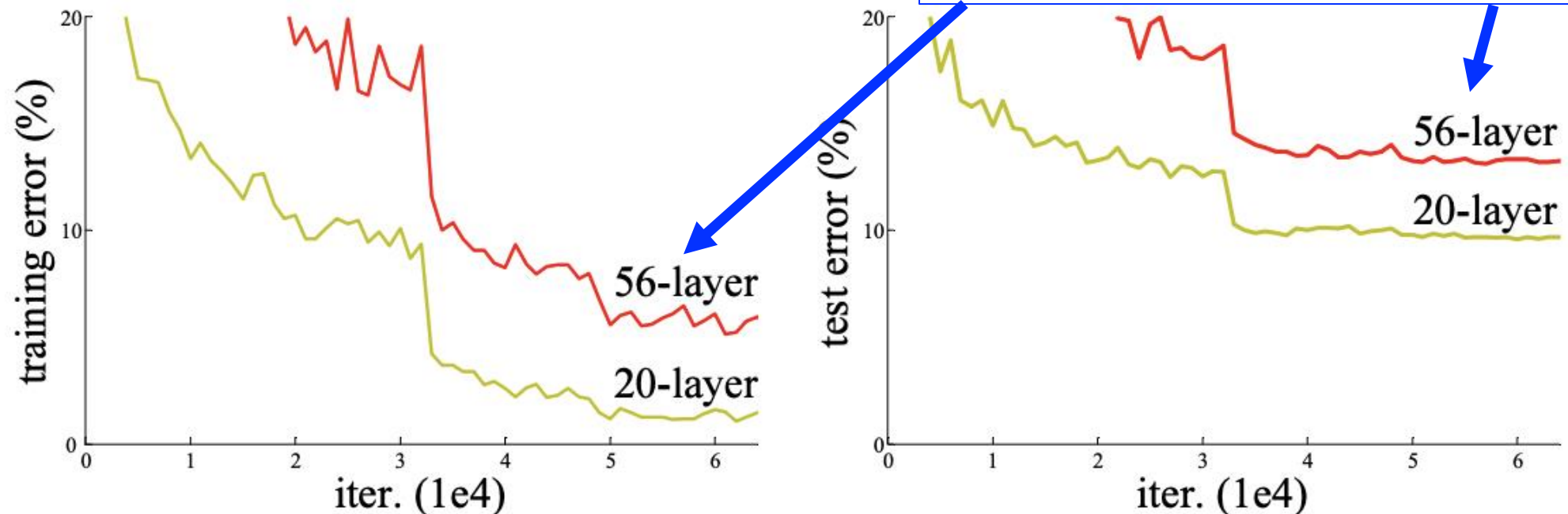
Motivating Observation

Idea: deeper networks should perform comparable or better than shallower networks since they can learn the shallower function by simply learning “identity” functions for later layers

Observation: adding more layers leads to WORSE results!

Is the problem overfitting? **NO**

Training data error (and test error) is greater with more layers



Motivating Observation

Idea: a deeper network should perform as good if not better than shallower networks since they can learn the shallower function by simply learning “identity” functions for later layers

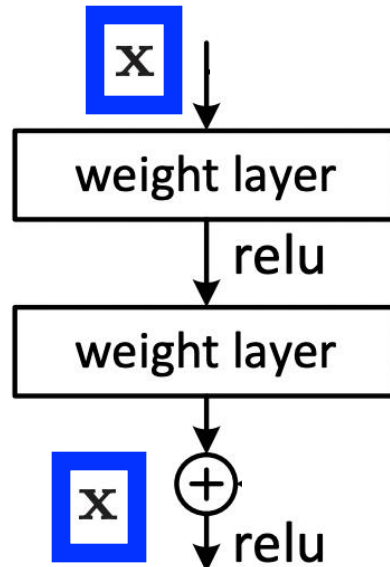
Observation: adding more layers leads to WORSE results!

Is the problem overfitting? NO

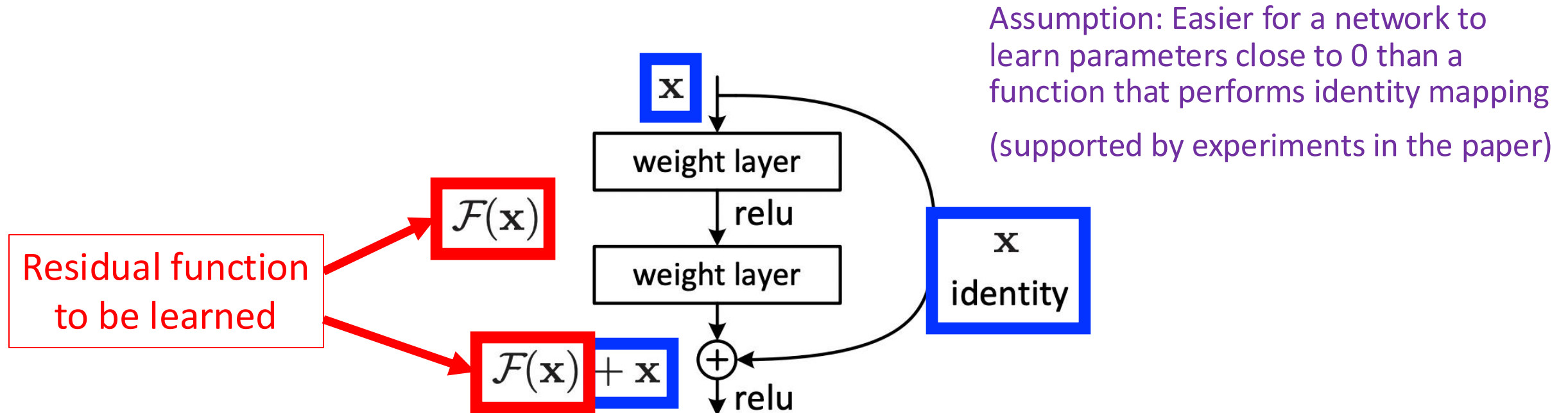
Problem: It is difficult for the algorithm to learn layers of identity mappings

Problem: Difficult to Perform Identity Mapping

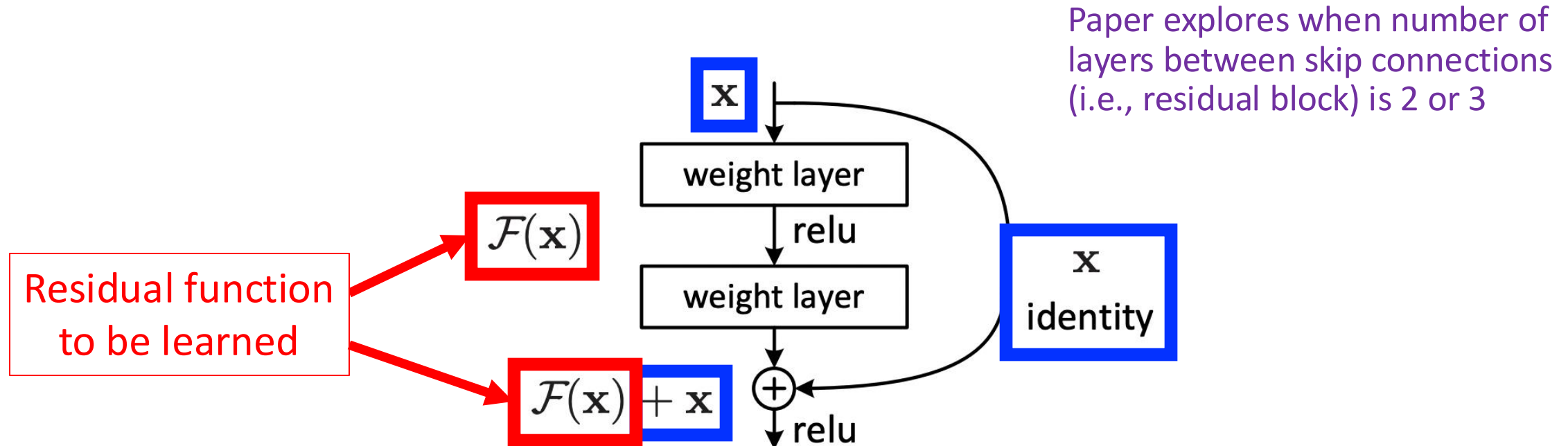
e.g.,



Key Idea: Skip Connections that Perform Identity Mapping



Key Idea: Skip Connections that Perform Identity Mapping



Key Idea: Skip Connections that Perform Identity Mapping

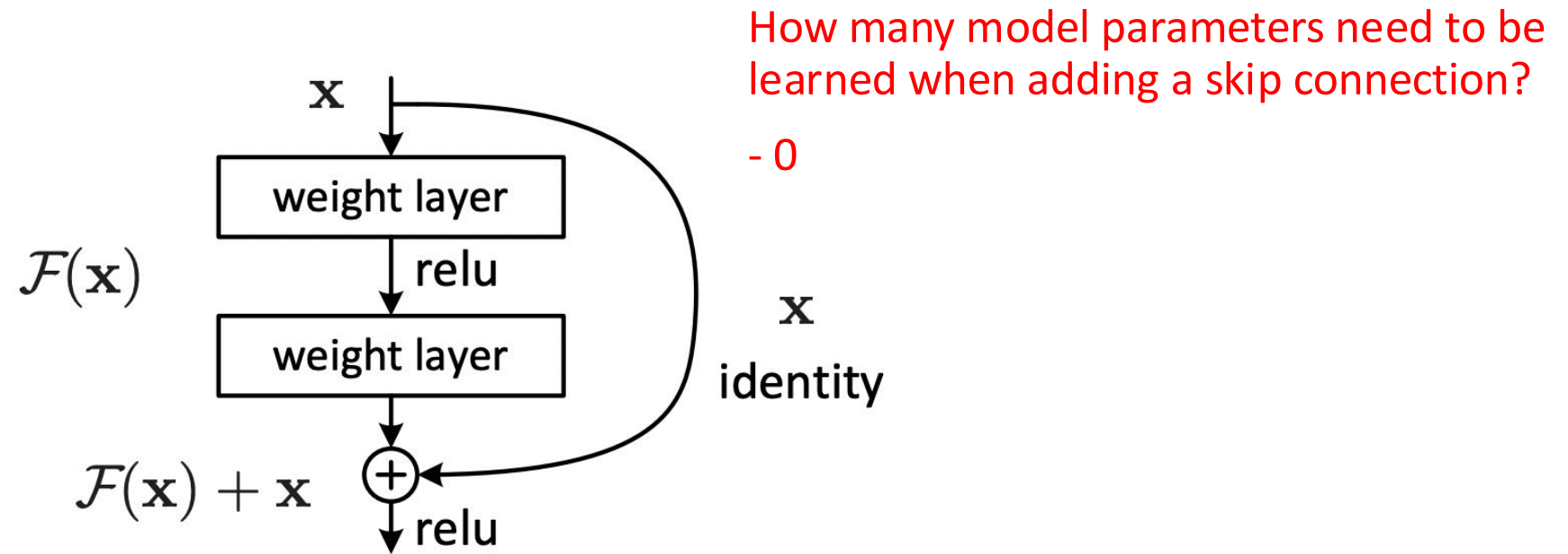
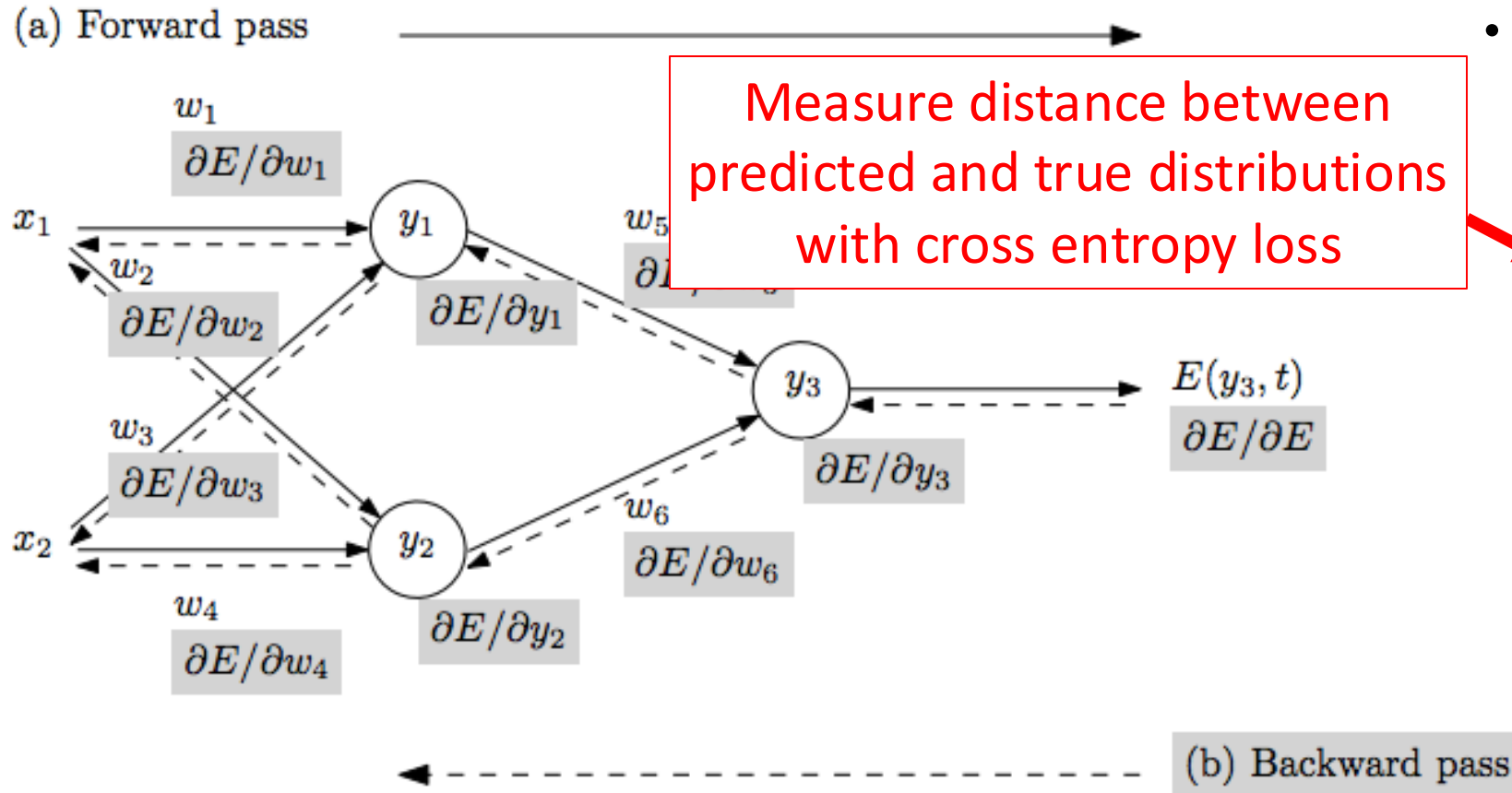


Figure 2. Residual learning: a building block.

ResNet Training (follows AlexNet)



- Repeat until stopping criterion met:
 1. **Forward pass:** propagate training data through model to make predictions
 2. **Error quantification:** measure dissatisfaction with a model's predictions on training data
 3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter
 4. Update each parameter using calculated gradients

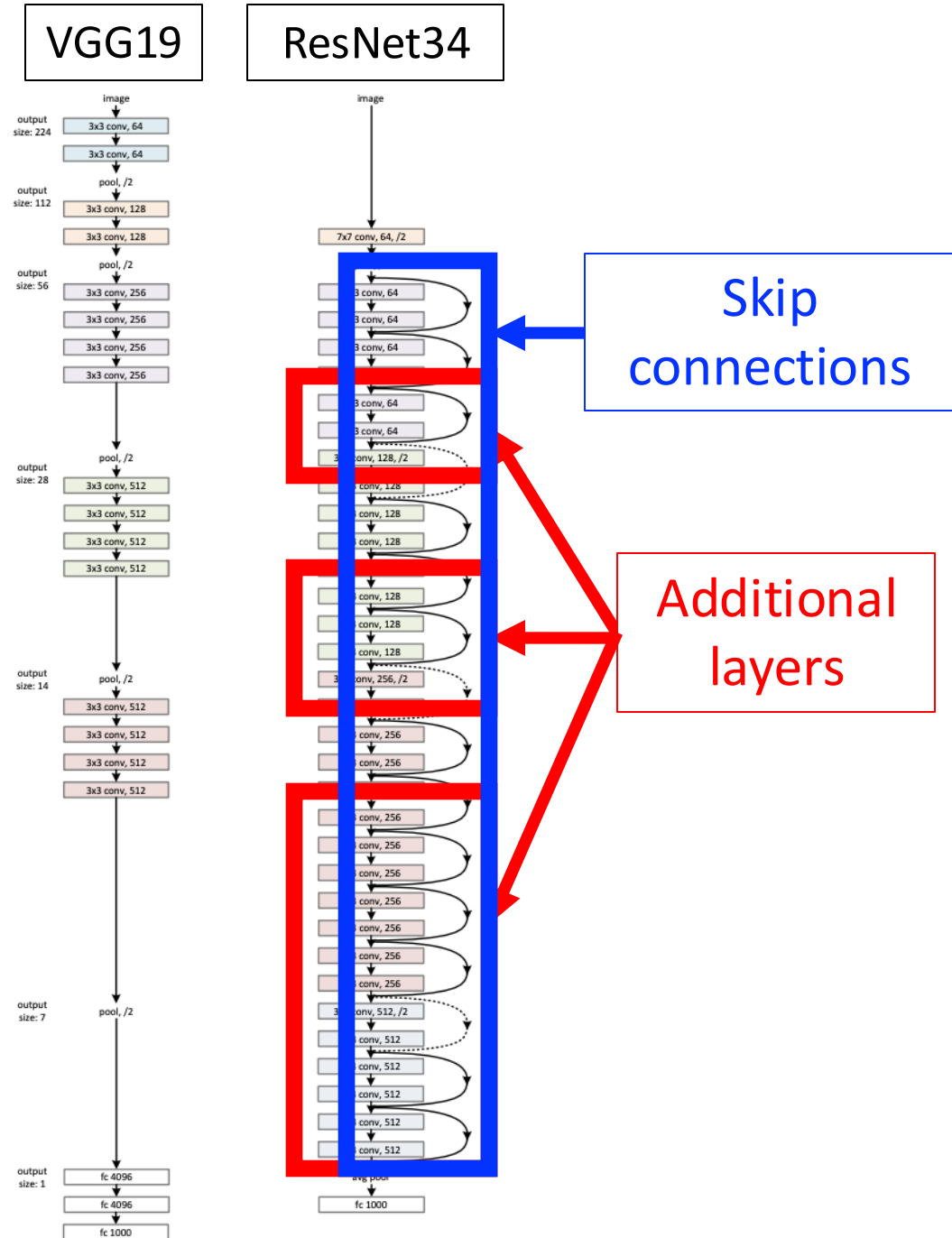
ResNet Training (follows AlexNet)

- Strategy to mitigate overfitting
 1. Data augmentation
 1. Random patches and their mirror images (2048x more data)
 2. Adjust RGB channels (using PCA to add multiples of principal components)

ResNet Implementations

Deep residual learning framework using **skip connections** enabled successfully learning **deeper models** than prior work

(18, 34, 50, 101, & 152 layers!)



Experimental Results on Validation Set

model	top-1 err.	top-5 err.
VGG-16 [40]	28.07	9.33
GoogLeNet [43]	-	9.15
PReLU-net [12]	24.27	7.38
ResNet-50	22.85	6.71
ResNet-101	21.75	6.05
ResNet-152	21.43	5.71

Performance
improves with
more layers



ResNet models outperform prior state-of-art models!

ResNet: Key Tricks for Going Deeper

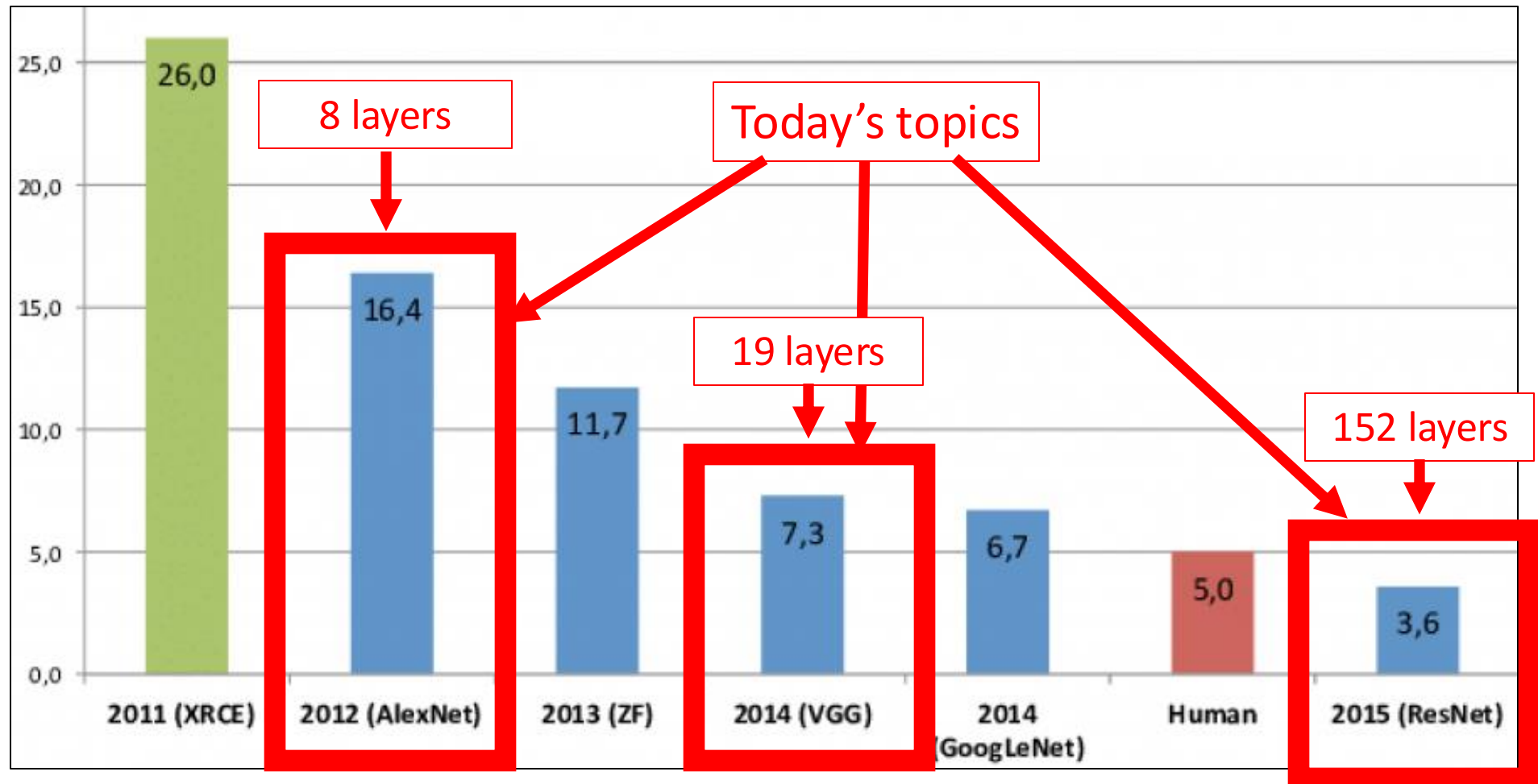
- Skip connections with residual learning

Object Recognition: Today's Topics

- ImageNet Challenge Top Performers
- Baseline Model: AlexNet
- VGG
- ResNet
- **Summary of CNN Era**

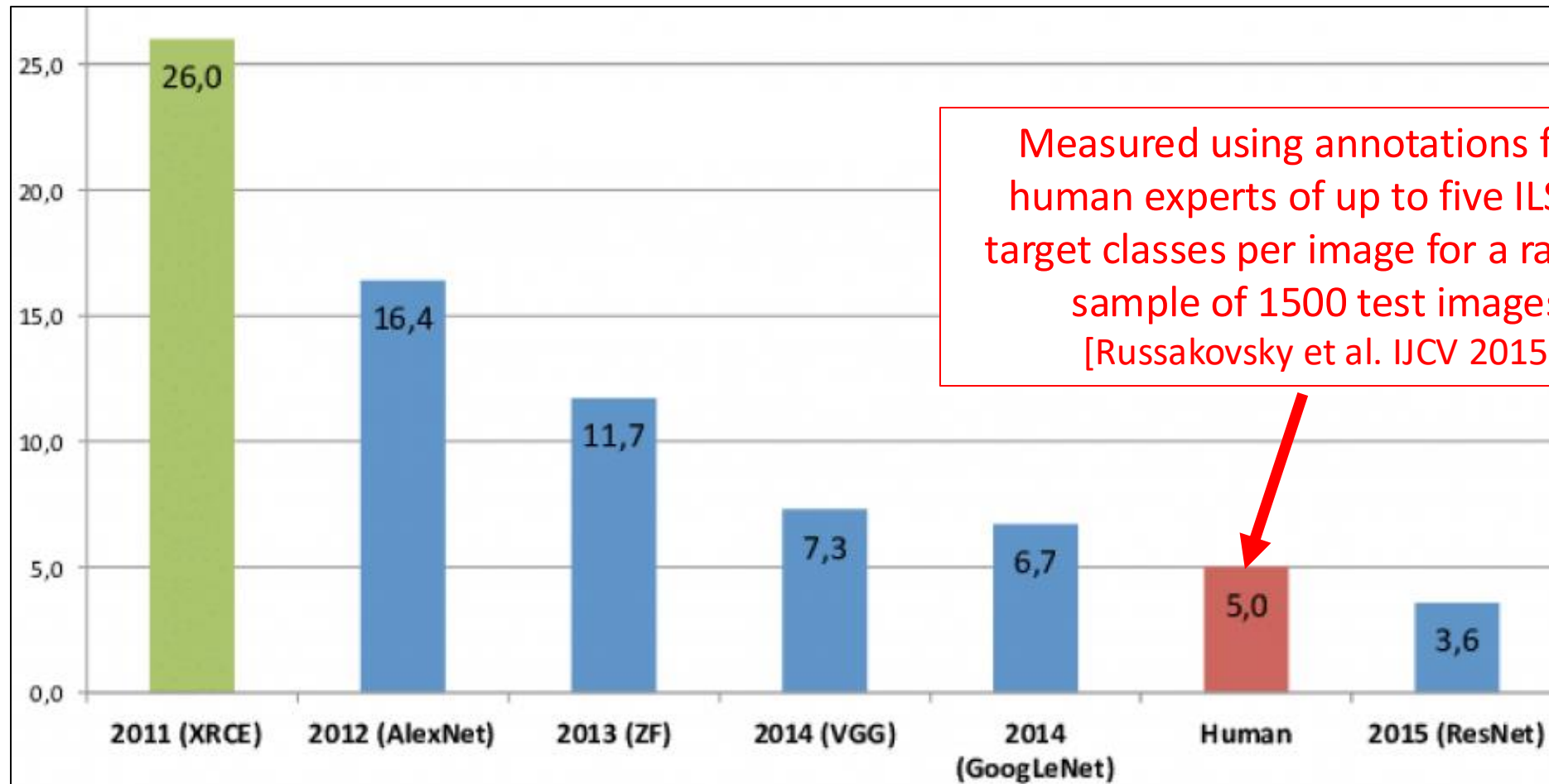
State-of-Art Design Models Go “Deeper”

Progress of models on ImageNet (Top 5 Error)



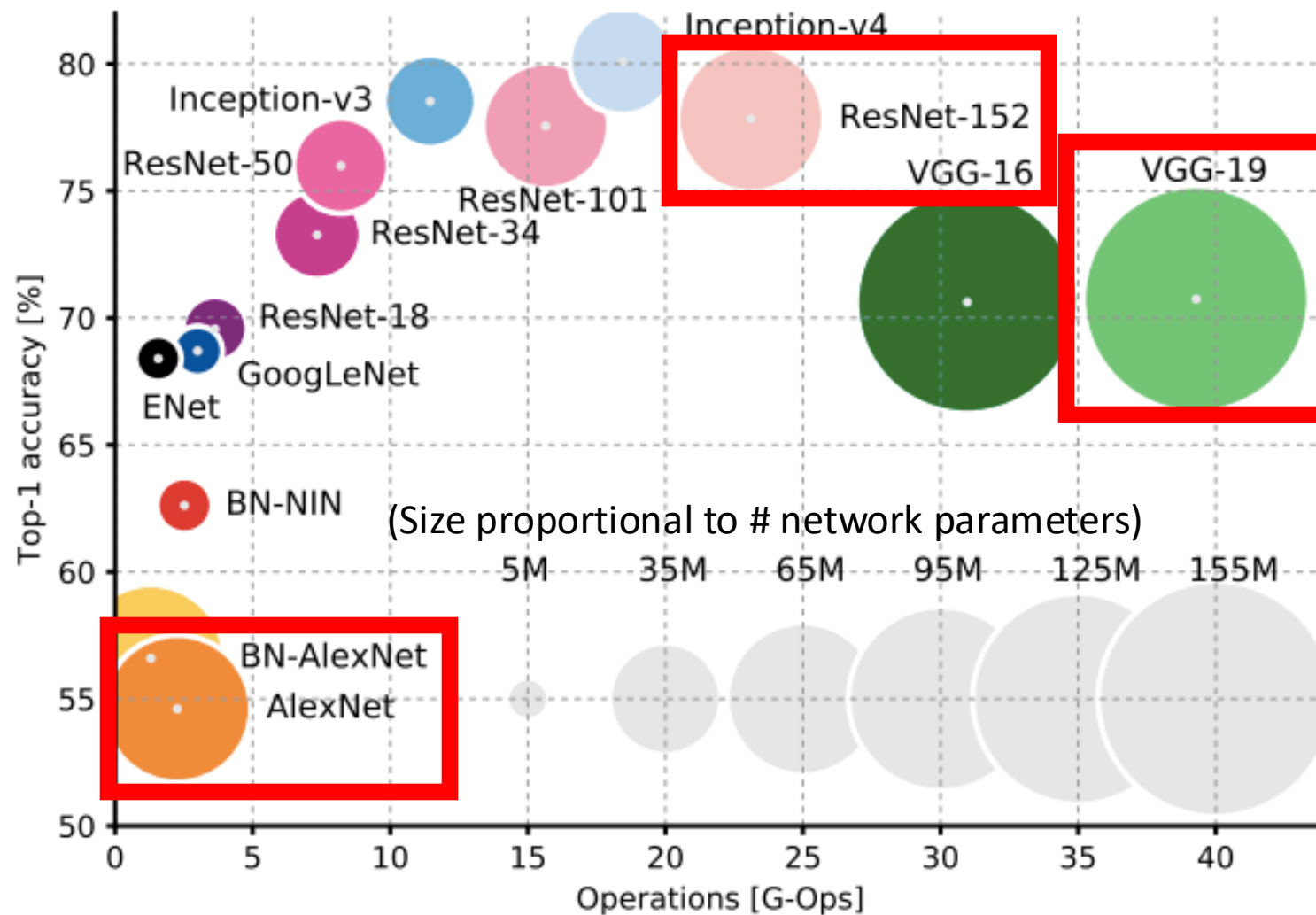
State-of-Art Model Exceeds Human Performance!

Progress of models on ImageNet (Top 5 Error)



Measured using annotations from human experts of up to five ILSVRC target classes per image for a random sample of 1500 test images.
[Russakovsky et al. IJCV 2015]

CNN Architectures Are a Great Start... Transformers to Follow



(required for a single forward pass)

Object Recognition: Today's Topics

- ImageNet Challenge Top Performers
- Baseline Model: AlexNet
- VGG
- ResNet
- Summary of CNN Era

The image features a central area with a radial gradient background, transitioning from a light center to a darker outer edge. This central area is framed by a dark grey border that mimics the appearance of a film strip, with white rectangular sprocket holes along the top and bottom edges. The text "The End" is centered within this frame.

The End