

Introduction to Computer Vision

Instructors: Jean Ponce and Matthew Trager
jean.ponce@inria.fr, matthew.trager@cims.nyu.edu

TAs: Jiachen (Jason) Zhu and Sahar Siddiqui
jiachen.zhu@nyu.edu, ss12414@nyu.edu

Outline

- Training neural networks
- Object detection with CNNs

Parametric supervised learning

- Training examples $(x_1, y_1), \dots, (x_n, y_n) \in \mathcal{X} \times \mathcal{Y}$
- A function class $\mathcal{F} = \{f_\theta : \mathcal{X} \rightarrow \mathcal{Y} \mid \theta \in \mathbb{R}^d\}$
- A loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$

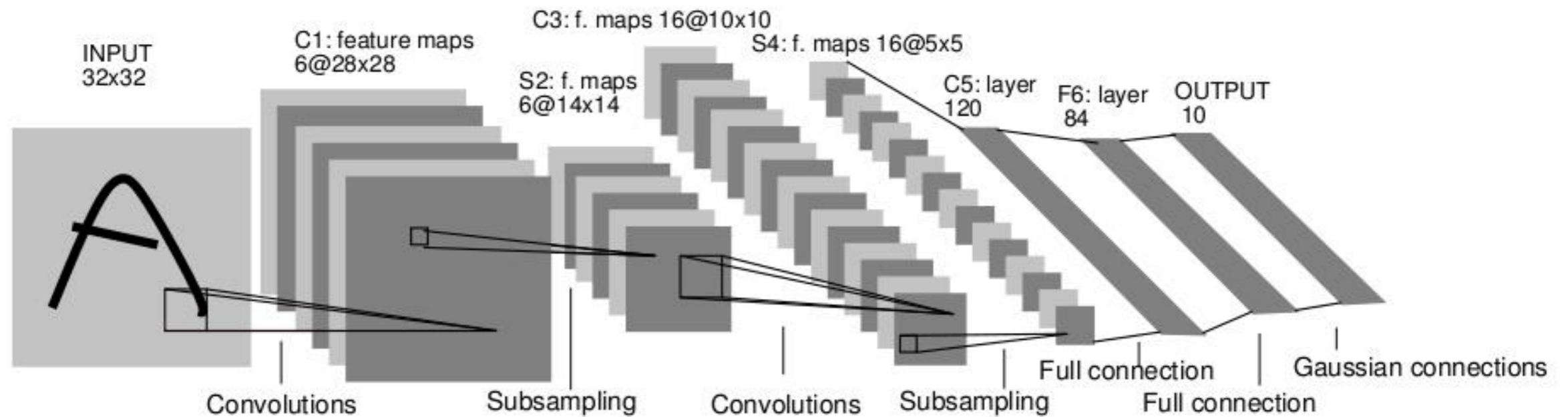
- Goal: minimize the empirical risk

$$\hat{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(f_\theta(x_i), y_i)$$

- Hope that this “generalizes”: if (X, Y) is a r.v., we would like to minimize

$$L(\theta) = \mathbb{E}[\ell(f_\theta(X), Y)]$$

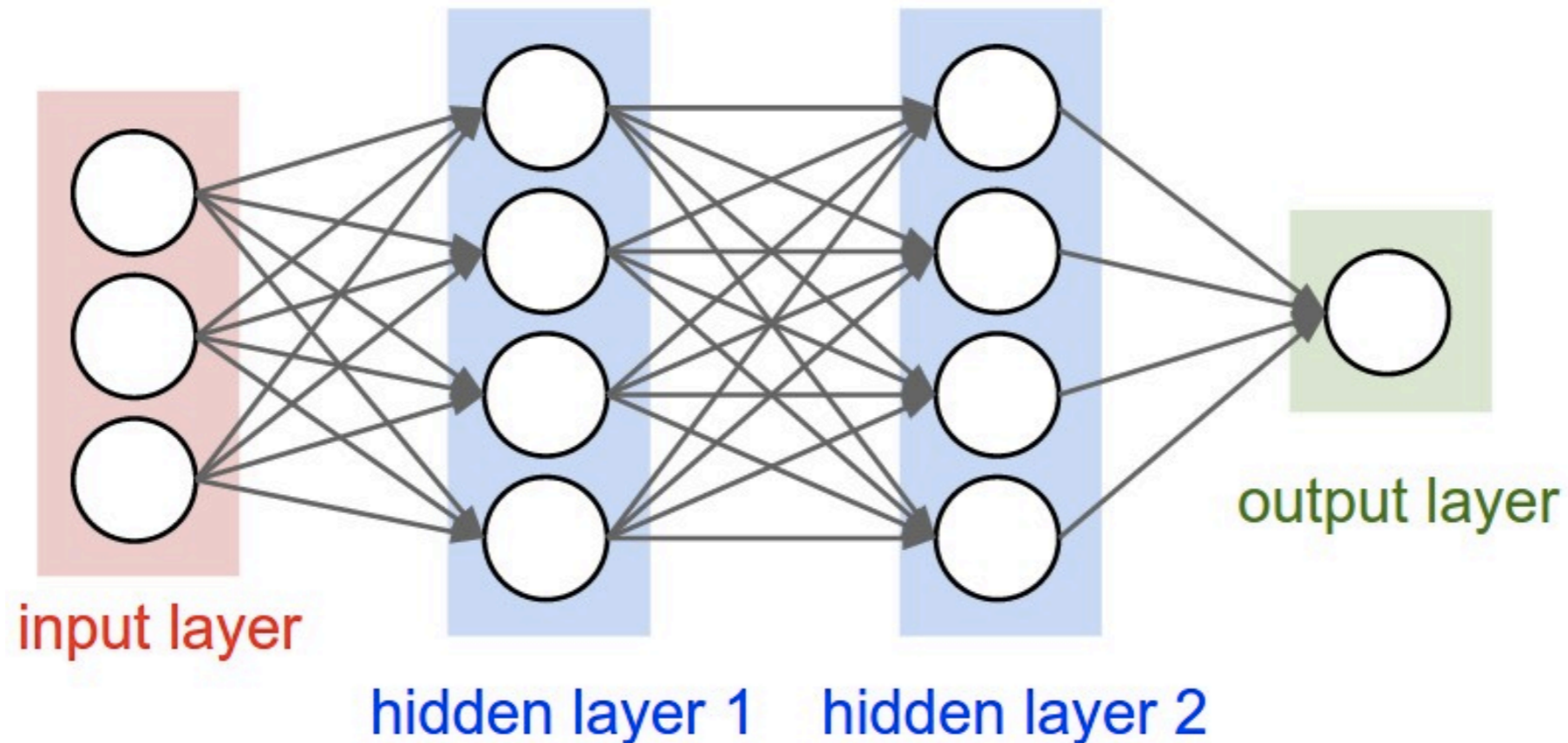
Neural networks



A neural network architecture describes a particular family of functions: $f_{\theta} : \mathcal{X} \rightarrow \mathcal{Y}$.

The parameters $\theta \in \mathbb{R}^d$ are the network's *weights*.

Feedforward NN (MLP)



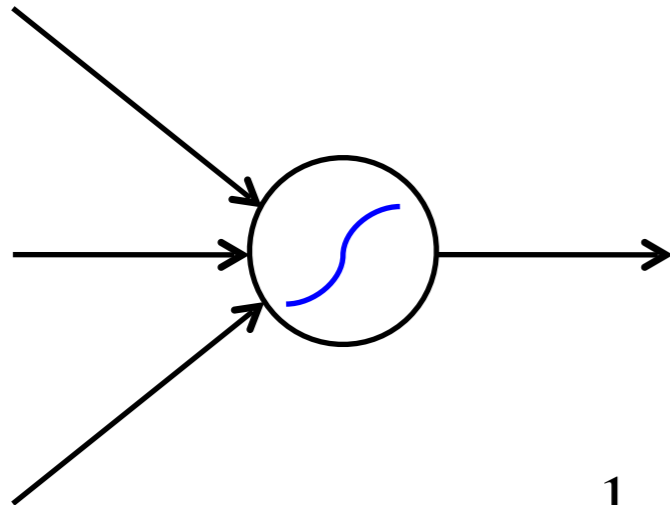
A feedforward NN is a composition of linear and non-linear functions, for example:

$$f_{\theta}(x) : \mathbb{R}^{d_0} \rightarrow \mathbb{R}^{d_h}, \quad f_{\theta}(x) = W_h \rho W_{h-1} \rho \dots W_2 \rho W_1 x,$$

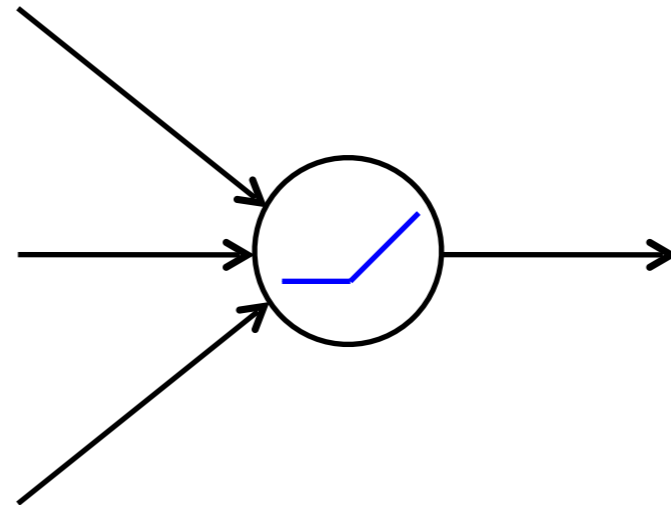
where $W_i \in \mathbb{R}^{d_i \times d_{i-1}}$, $\theta = (W_h, \dots, W_1)$, and ρ is a non-linear map acting coordinate-wise.

Non-linearity

- The nonlinearity ρ should be (almost everywhere) differentiable.

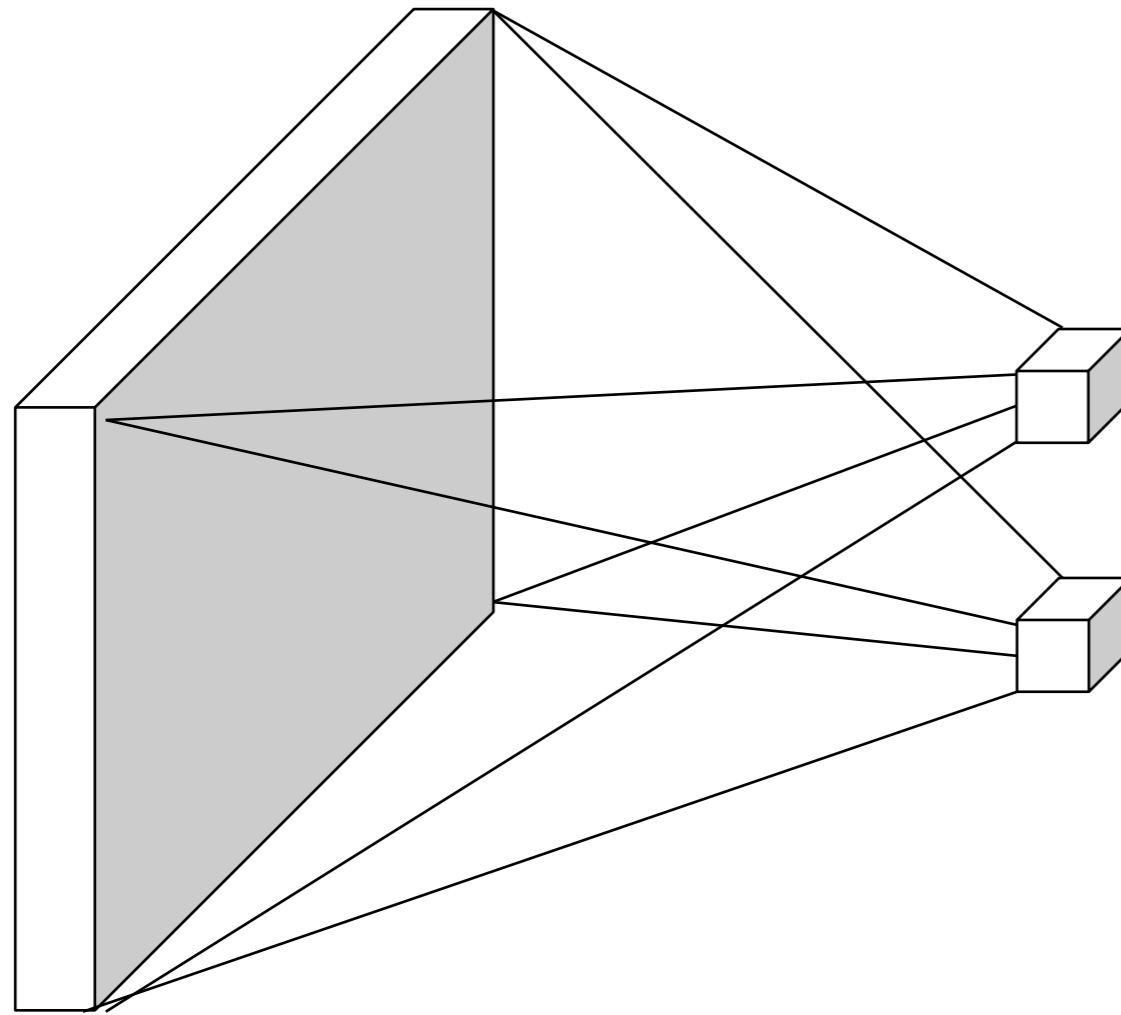


Sigmoid: $\rho(t) = \frac{1}{1 + e^{-t}}$



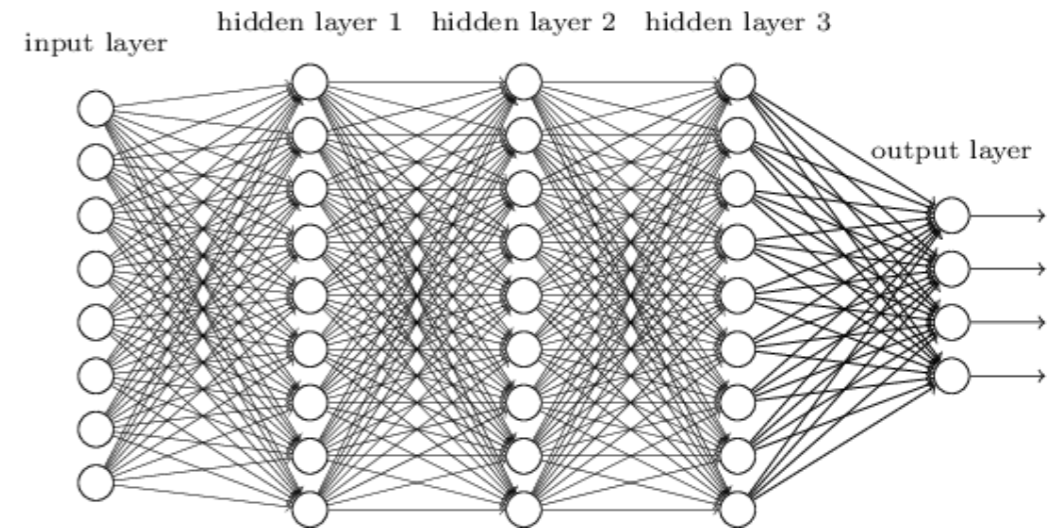
Rectified linear unit (ReLU): $\rho(t) = \max(0, t)$

Neural networks for images

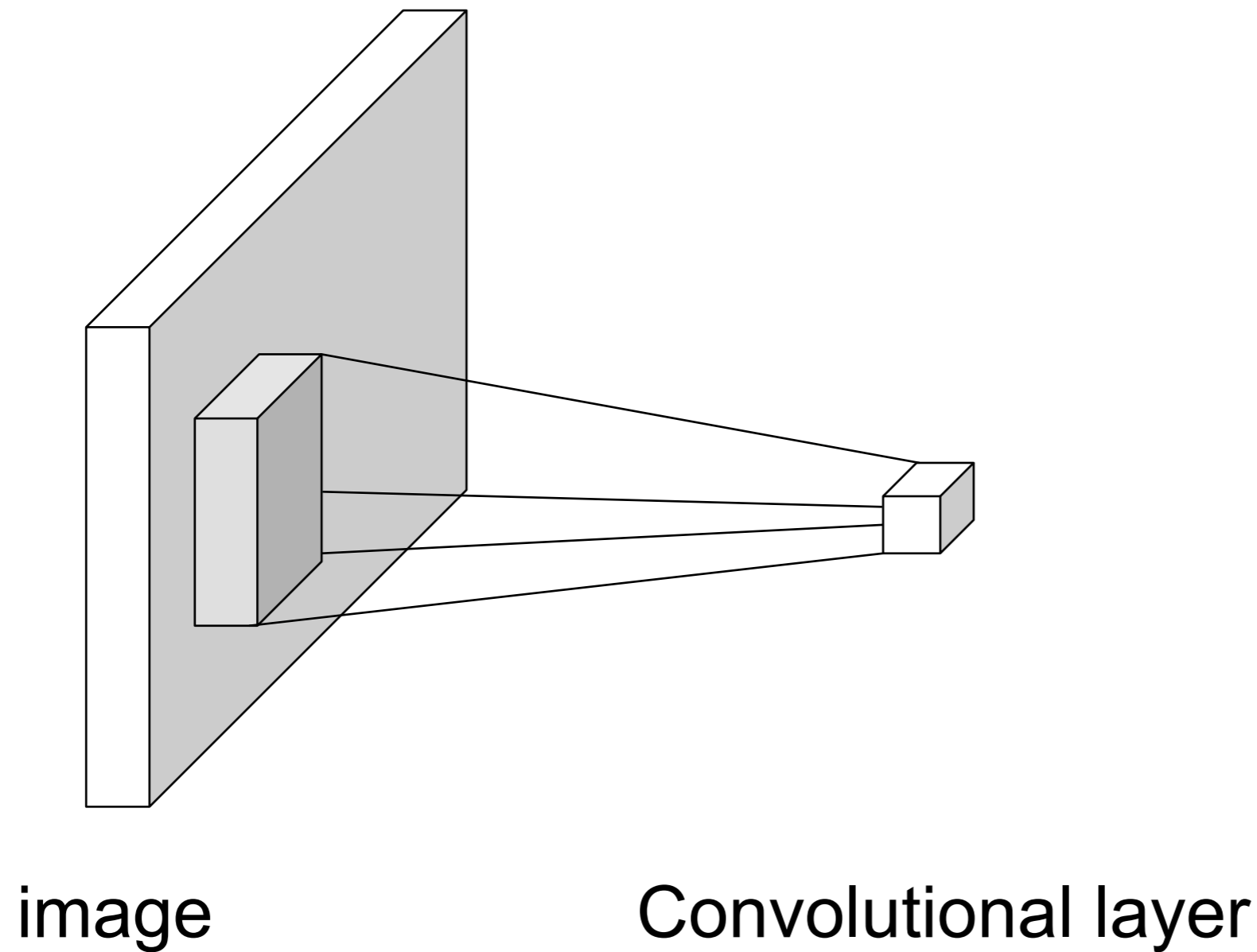


image

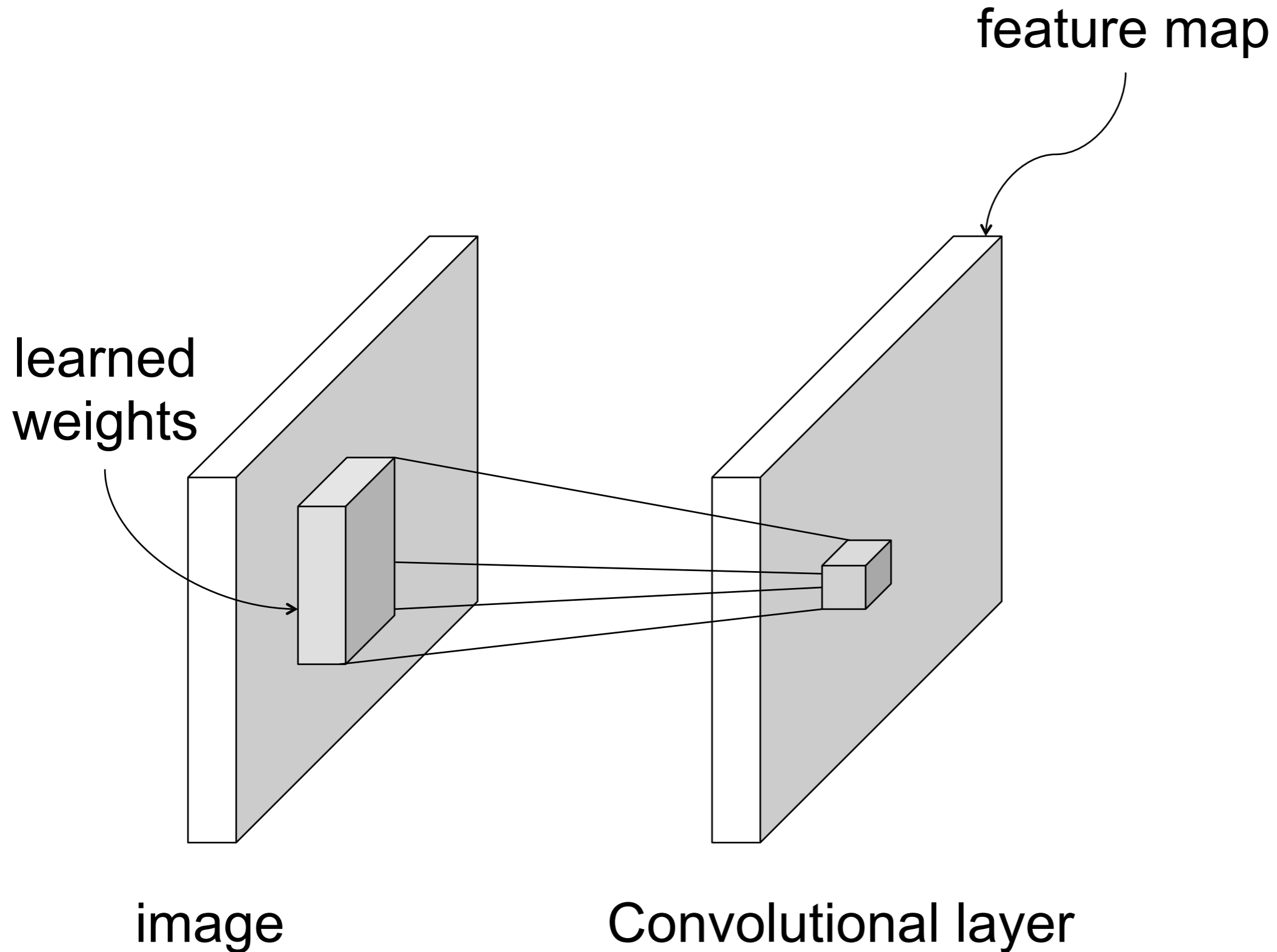
Fully connected layer



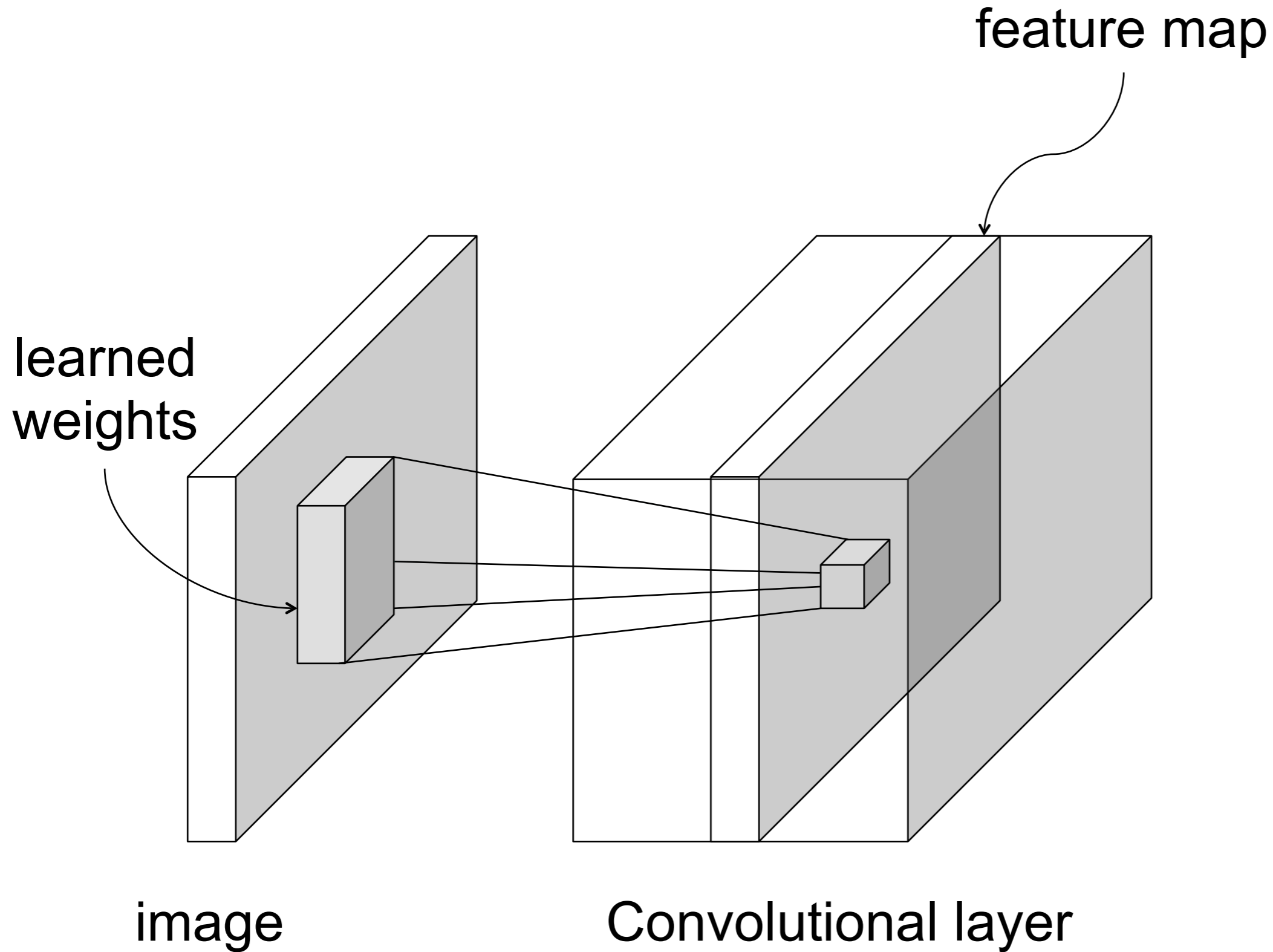
Neural networks for images



Neural networks for images



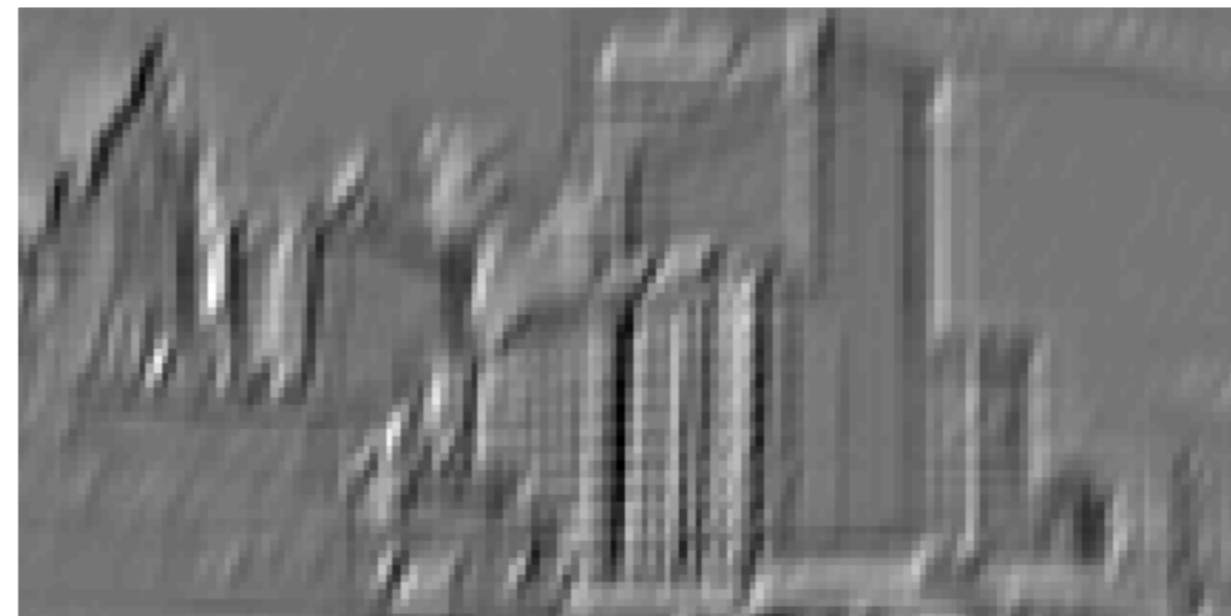
Neural networks for images



Convolution as feature extraction

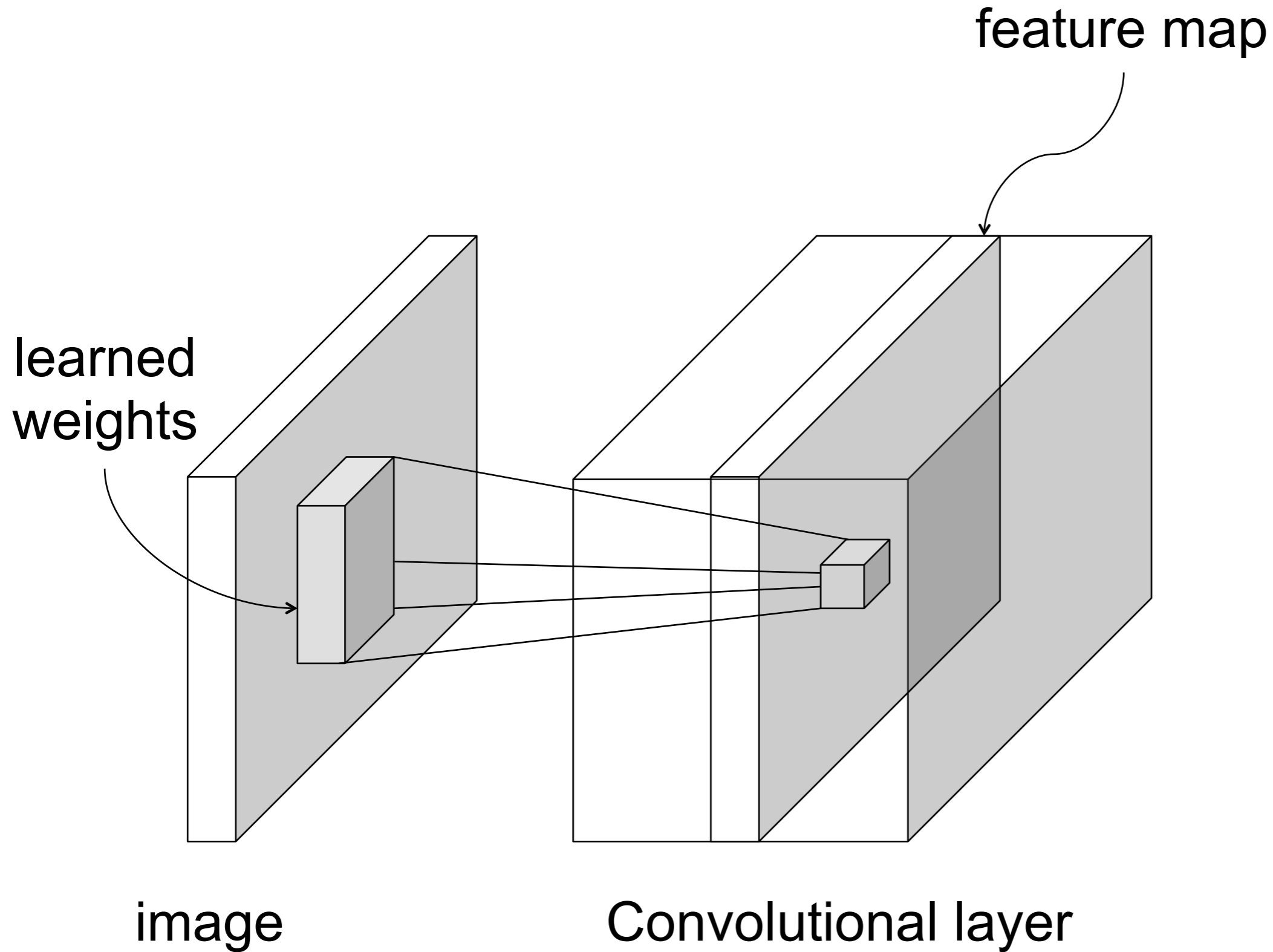


Input

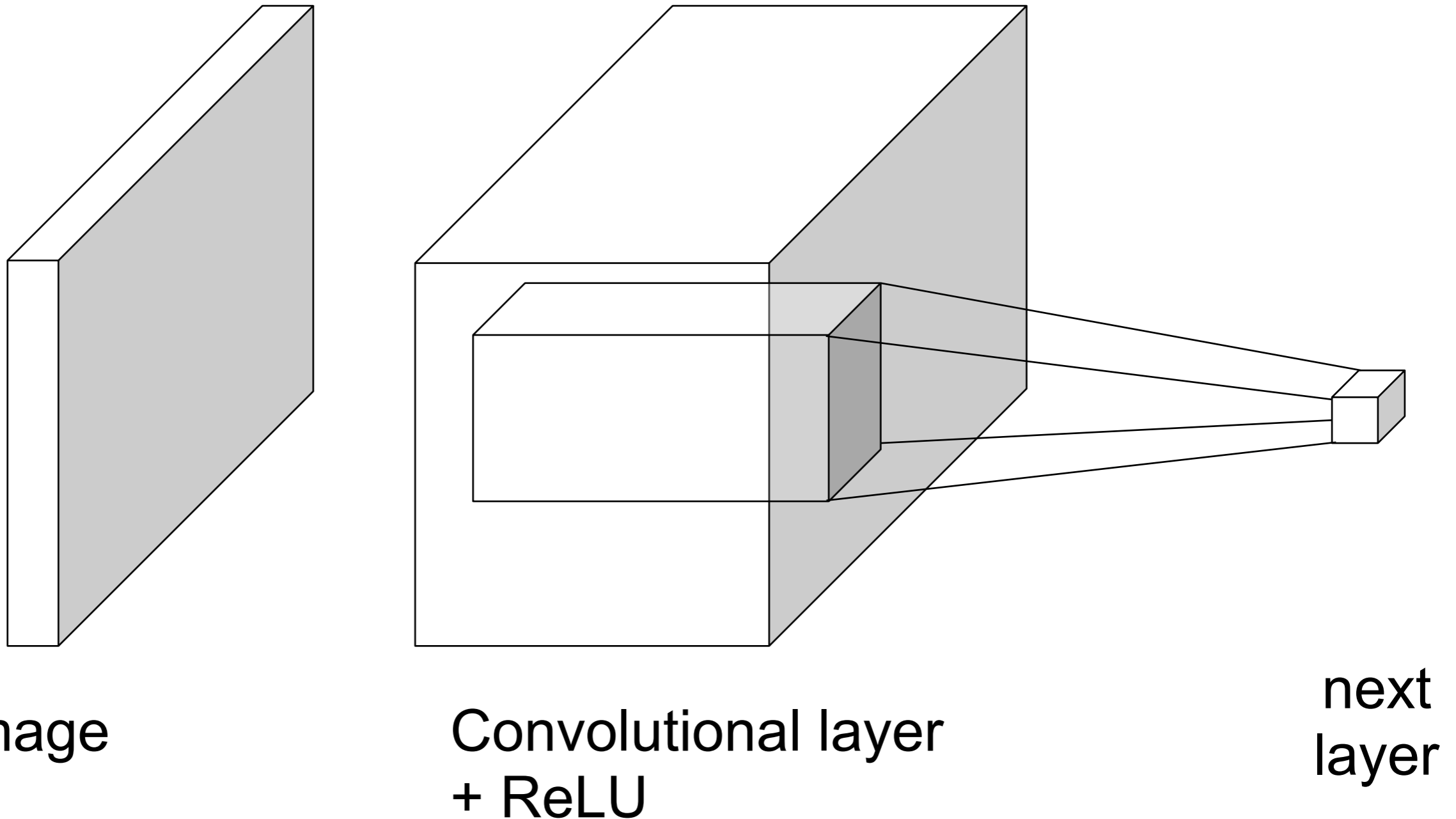


Feature Map

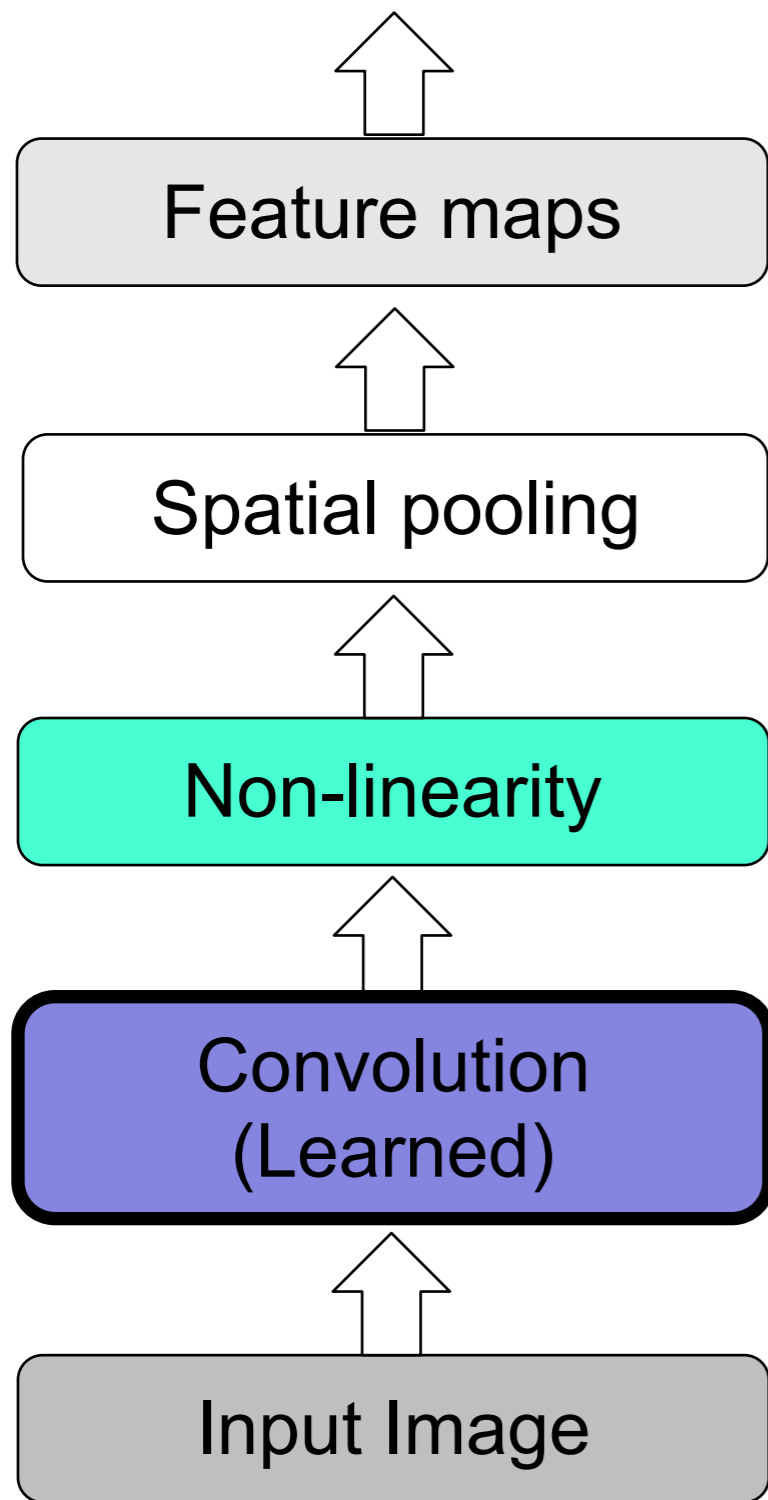
Neural networks for images



Neural networks for images



Key operations in a CNN

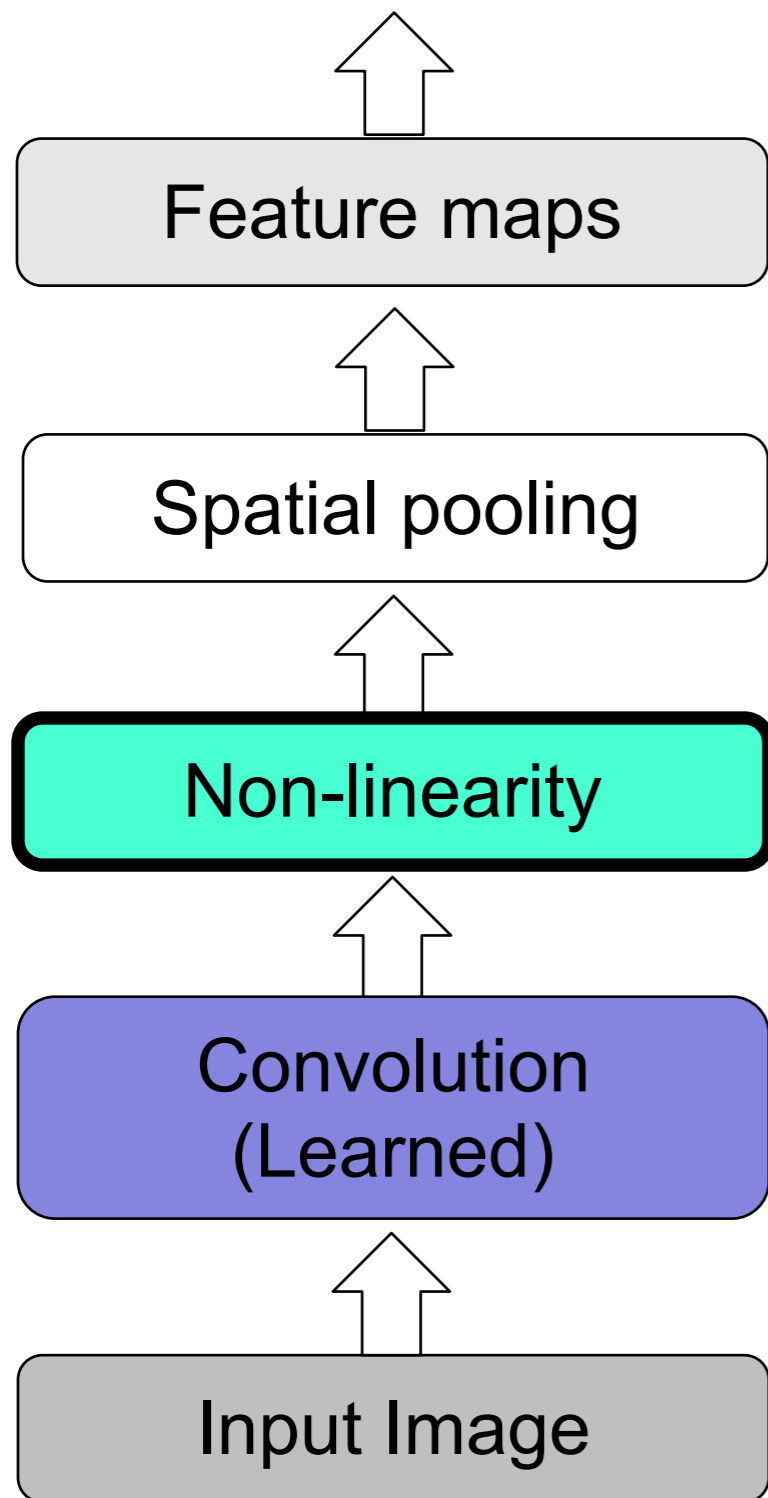


Input

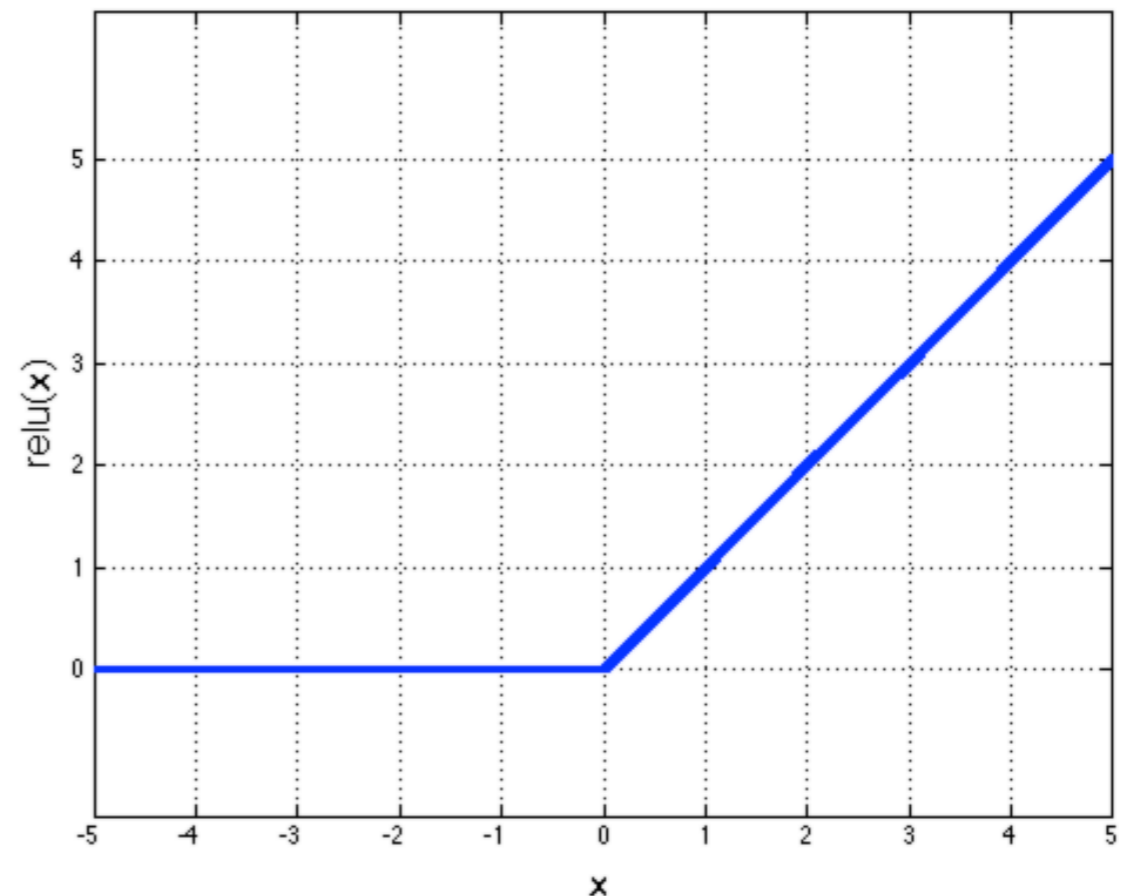


Feature Map

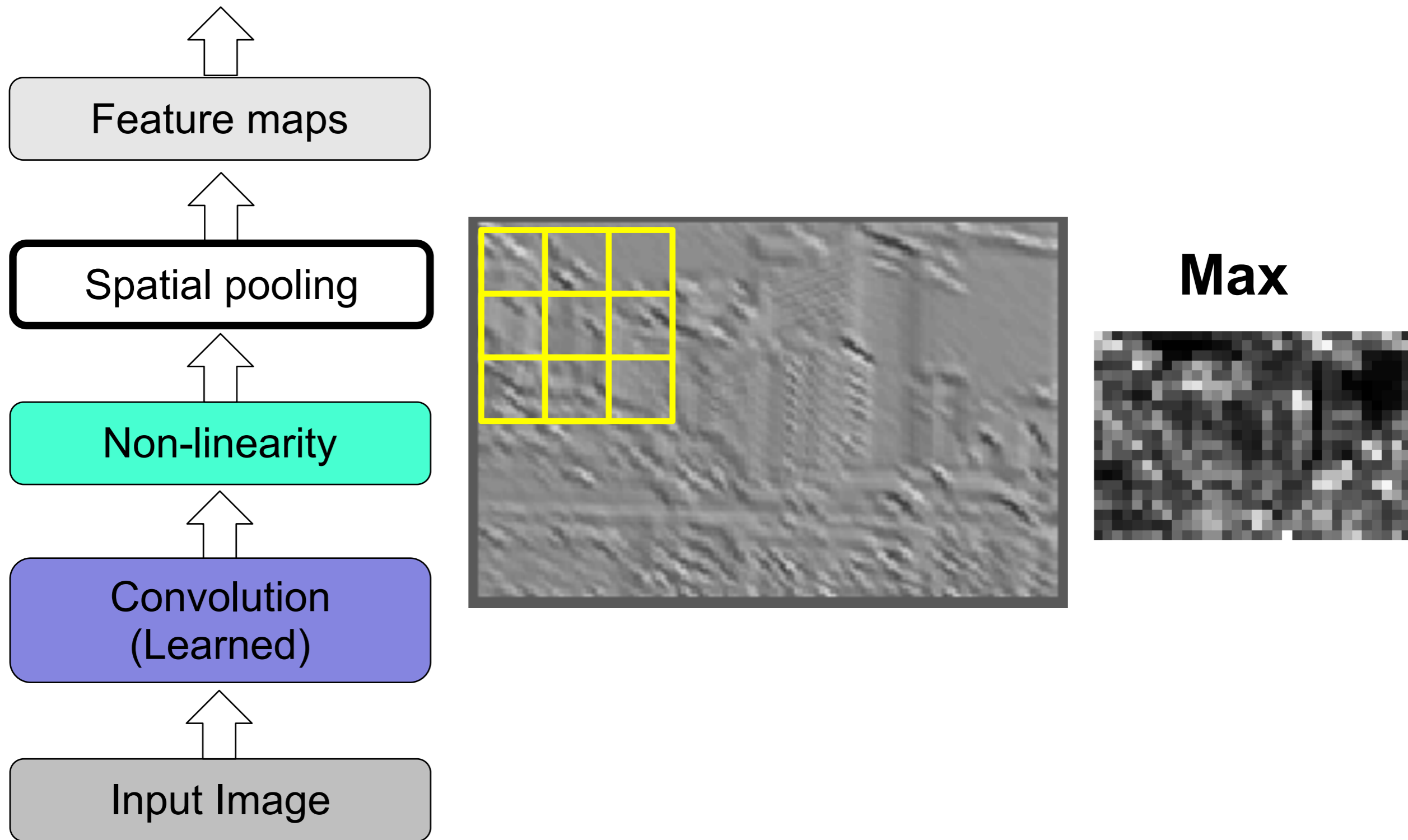
Key operations in a CNN



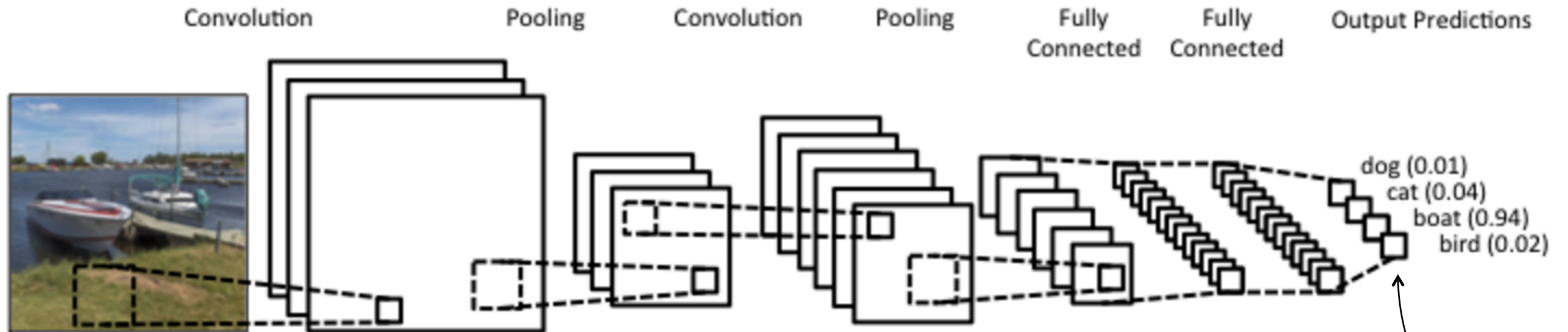
Rectified Linear Unit (ReLU)



Key operations in a CNN



Key operations in a CNN



Softmax layer:

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

Loss functions

The objective to minimize is of the form

$$L(\theta) = \sum_{i=1}^n \ell(f_{\theta}(x_i), y_i)$$

Examples for ℓ :

- *Regression:*

- Quadratic loss: $\ell(y, \hat{y}) = \|y - \hat{y}\|_2^2, \quad y, \hat{y} \in \mathbb{R}^d$

- *Classification:*

- *Cross-entropy:* $\ell(y, \hat{y}) = - \sum_i y_i \log(\hat{y}_i), \quad y \in \{0,1\}^d, \hat{y} \in [0,1]^d.$

Classification is usually better! (easier problem)

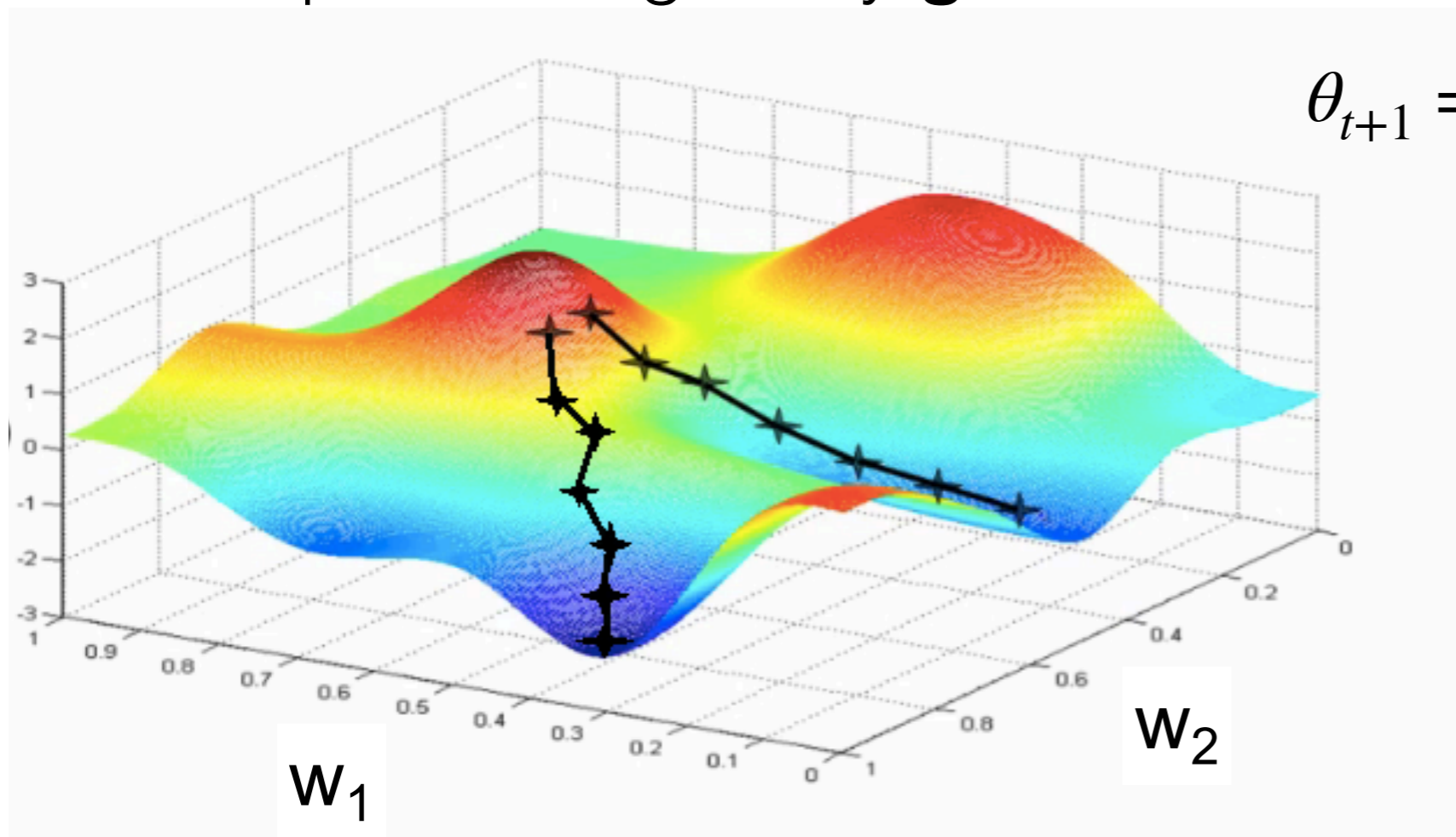
Training of multi-layer networks

- Find network weights to minimize the error between true and estimated labels of training examples:

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \ell(f_{\theta}(x_i), y_i)$$

- Update weights by **gradient descent**:

$$\theta_{t+1} = \theta_t - \gamma_t \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \ell(f_{\theta}(x_i), y_i)$$



Gradient descent

$$\theta_{t+1} = \theta_t - \gamma_t \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \ell(f_{\theta}(x_i), y_i)$$

- Need to choose the learning rate policy γ_t
- Can get stuck in a local minima ($L(\theta)$ is not convex)
- Each step can be expensive to compute if the dataset is large

Stochastic gradient descent

- Idea: instead of computing the gradient, compute an approximation

$$\theta_{t+1} = \theta_t - \gamma_t \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \ell(f_{\theta}(x_i), y_i)$$

↓

$$\theta_{t+1} = \theta_t - \gamma_t \frac{\partial}{\partial \theta} \ell(f_{\theta}(x_{i_t}), y_{i_t})$$

Note that in expectation

$$\mathbb{E}(\theta_{t+1} | \theta_t) = \theta_{t+1} = \theta_t - \gamma_t \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \ell(f_{\theta}(x_i), y_i)$$

Batch SGD

- It's faster to compute several gradients in parallel
 - Some variance is good, too much can be bad
- compute gradient estimates on a batch instead of a single sample.

$$\theta_{t+1} = \theta_t - \gamma_t \frac{1}{K} \sum_{k=1}^K \frac{\partial}{\partial \theta} \ell(f_{\theta}(x_{i_k,t}), y_{i_k,t})$$

- In practice, using batches as large as possible so that the network fits in the GPU memory

Beyond SGD

- Many other algorithms: see <https://ruder.io/optimizing-gradient-descent/>

- GD with momentum: encourage directions that are coherent:

$$v_t = \eta v_{t-1} + \gamma \nabla_{\theta_t}$$
$$\theta_{t+1} = \theta_t - v_t$$

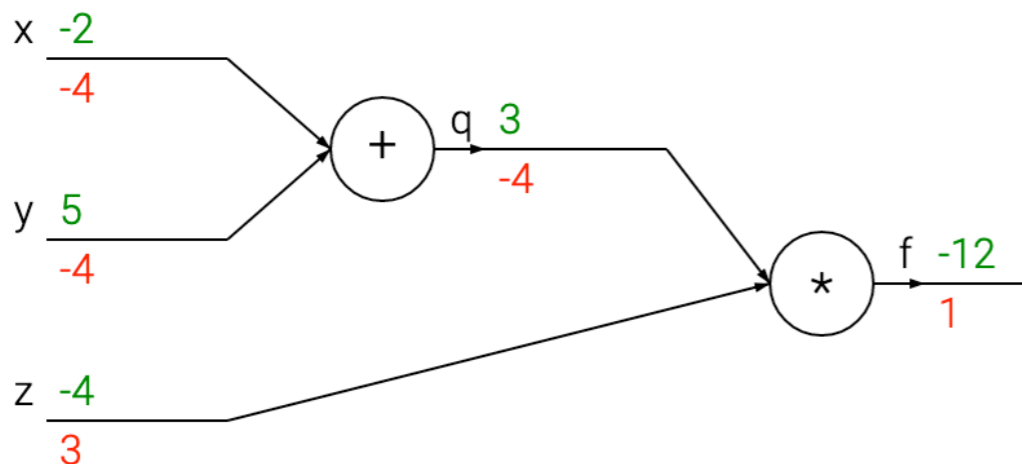
- Adagrad, Adadelata, RMSProp, ADAM...
- All these optimizers are coded in standard deep learning libraries

Backpropagation

- AKA “reverse-mode automatic differentiation” (“chain rule”)
 - *not* finite differences
 - *not* symbolic differentiation

Simple example: if $y = f_h(f_{h-1}(\dots f_1(w_0)))$, $w_i = f_i(w_{i-1})$, then compute

$$\frac{dy}{dw_i} = \frac{dy}{dw_{i+1}} \frac{dw_{i+1}}{dw_i}$$



Other example: $f(x, y, z) = (x + y)z$,
with $x = -2, y = 5, z = -4$

Initialization

- Since loss is not convex, initialization is really important!
- If two neurons on the same layer are initialized with the same weight, they will stay the same.
- For ReLU activations, PyTorch uses “Kaiming Uniform”:

$$w_i \sim \mathcal{U} \left[-\sqrt{k}, \sqrt{k} \right] \quad \text{where } k = \frac{1}{n_{in}}$$

Regularization

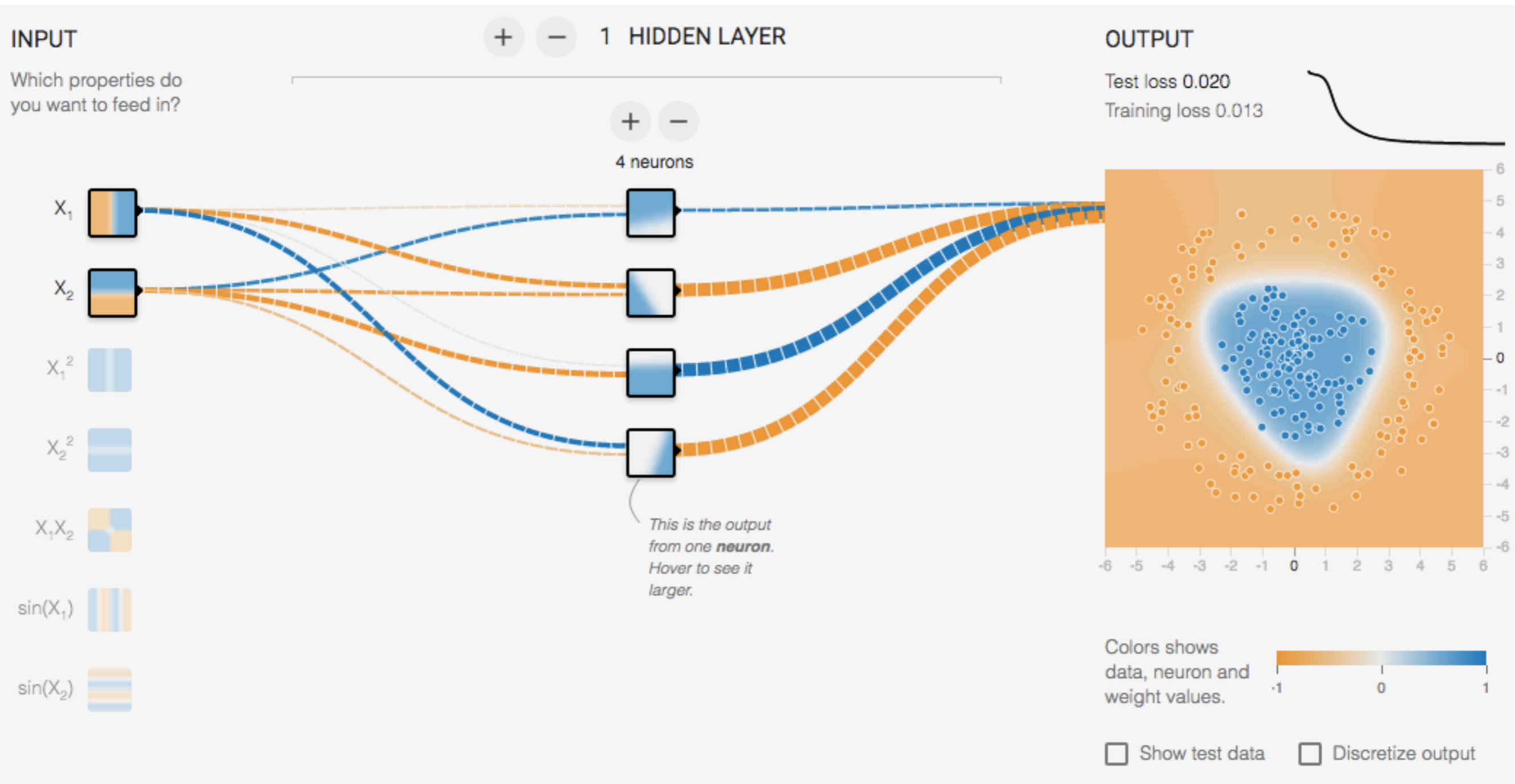
Several techniques for regularization (to avoid overfitting):

- L2 penalization on the weights (also called weight decay)

$$L(\theta) = \sum_{i=1}^n \ell(f_{\theta}(x_i), y_i) + \lambda \|\theta\|^2$$

- Dropout: only keep a neuron active with some probability p , or set it to zero otherwise.
- Early stopping
- Data augmentation

Interactive Demo



<http://playground.tensorflow.org/>

Common problem: vanishing/exploding gradients

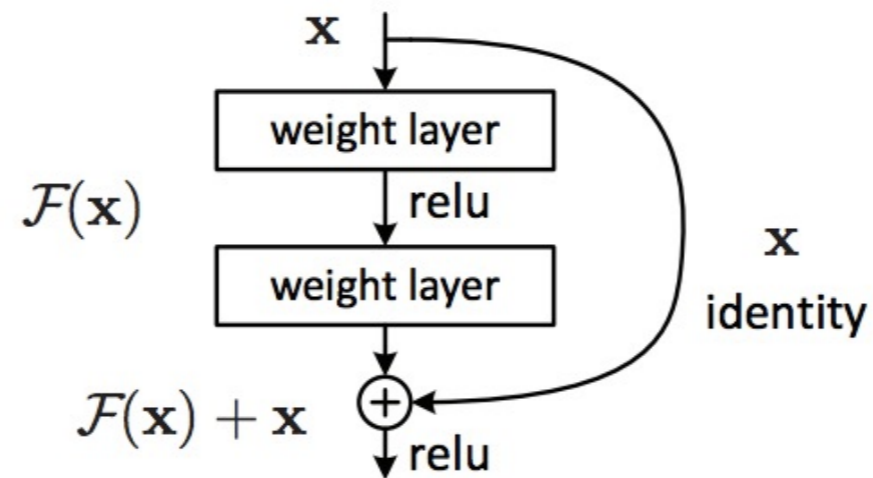
- If all linear layers have width 1, and are initialized at α :

$$\frac{\partial L}{\partial w_1} = \alpha^{N-1} \frac{\partial L}{\partial w_N}$$

- Depending on the value of α , the gradient with respect to w_1 will be huge or very small.
- Similar effects happen for more complex deep networks (even worse if non-linearities have ~ 0 gradients).

Vanishing/exploding gradients: solutions

- Use ReLU (non-saturating)
- Use skip-connections



- Use batch-normalization

Batch normalization

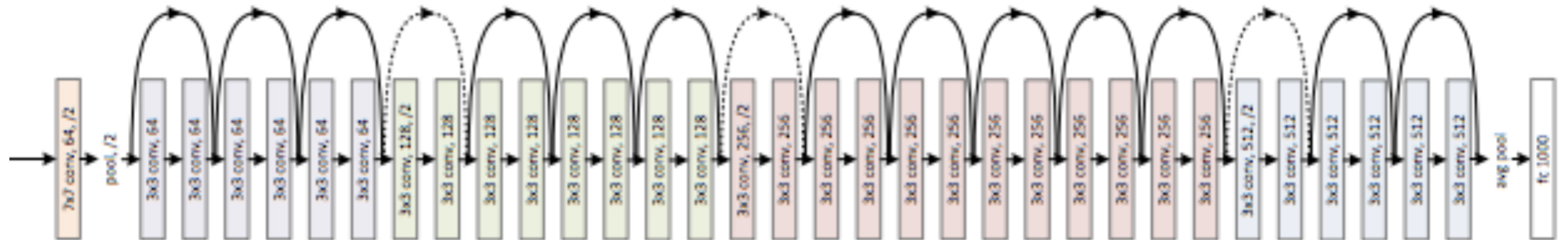
- Additional layer to avoid vanishing or exploding signal.
- Idea: normalize the data everywhere in the network using estimates of the mean/variance

$$BN_{\alpha, \beta}(x, \mu, \sigma) = \alpha \frac{x - \mu}{\sigma} + \beta$$

- α, β are learned, μ, σ are estimated over a mini-batch.
- Batch-norm layers are typically placed just before non-linearities

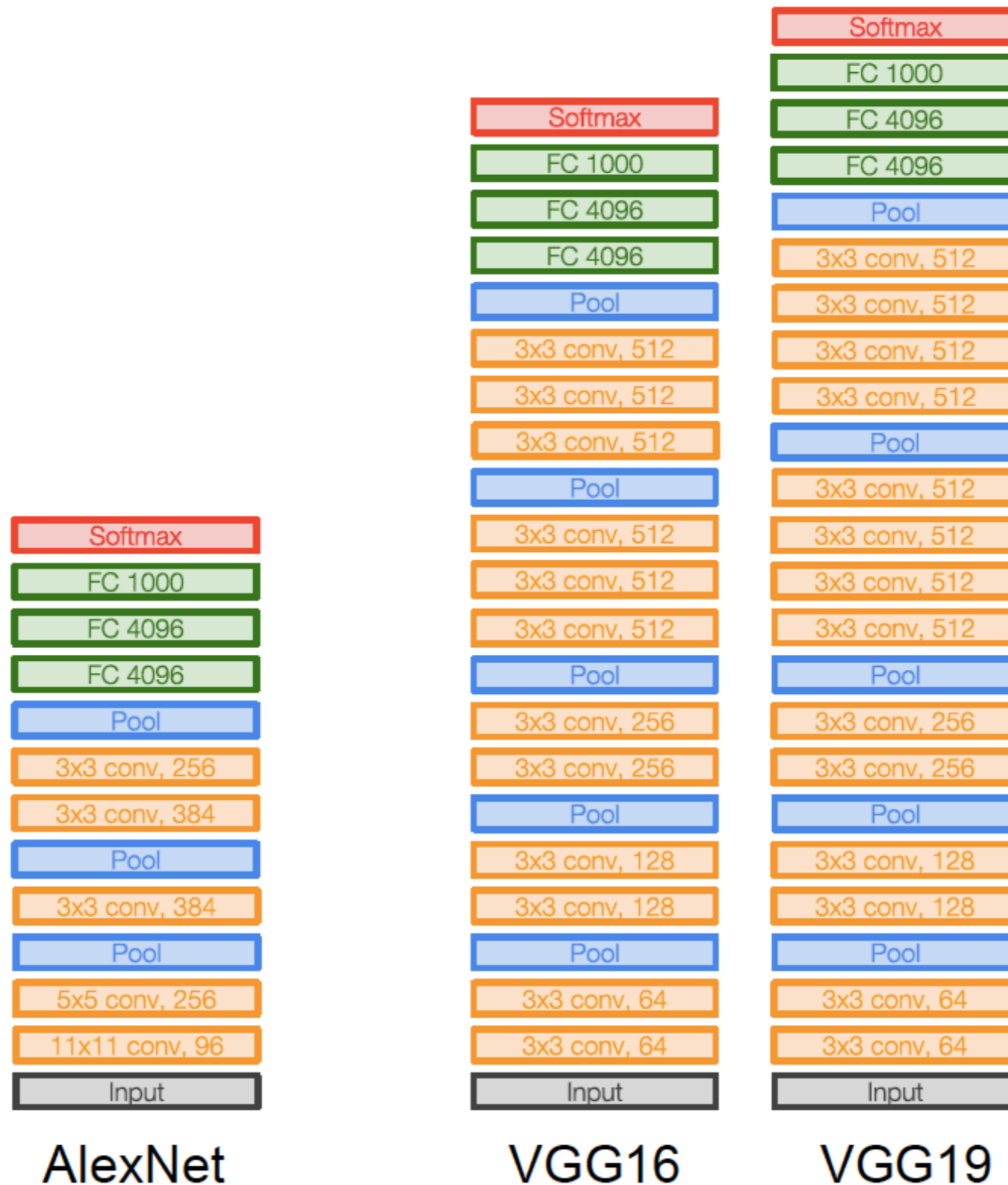
Typical modern network: Resnet

- Identical layers are repeated many times



- 3x3 convolution kernels + BN + ReLUs.
- Skip connections.
- Only one fully connected layer at the end.

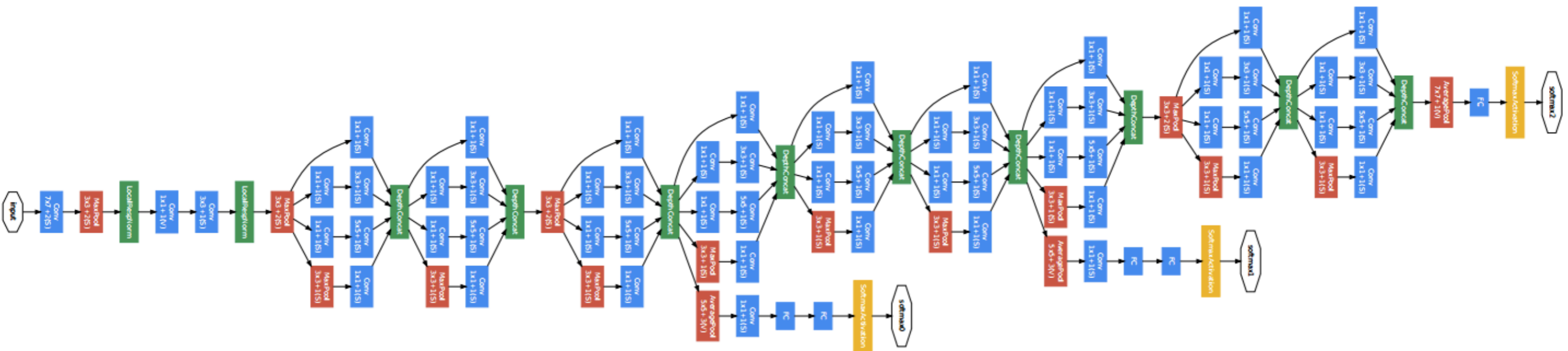
AlexNet and VGGNet



[Image source](#)

K. Simonyan and A. Zisserman, [Very Deep Convolutional Networks for Large-Scale Image Recognition](#), ICLR 2015

GoogLeNet



ResNet

AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)



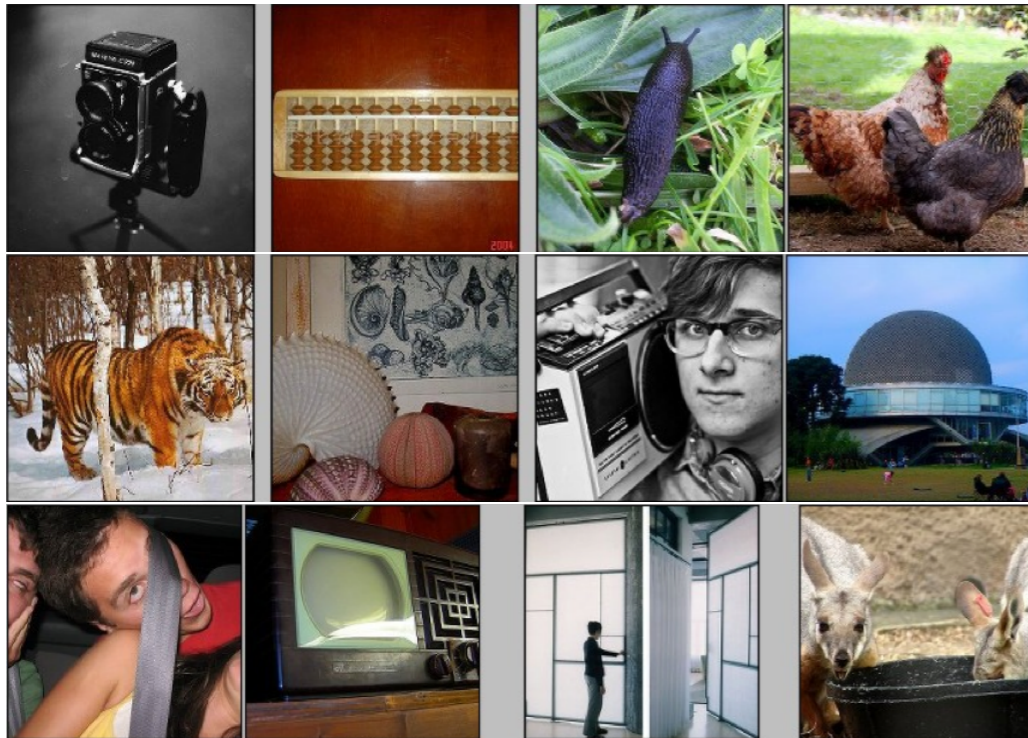
ResNet, **152 layers**
(ILSVRC 2015)

K. He, X. Zhang, S. Ren, and J. Sun, [Deep Residual Learning for Image Recognition](#), CVPR 2016 (Best Paper)



ImageNet Challenge

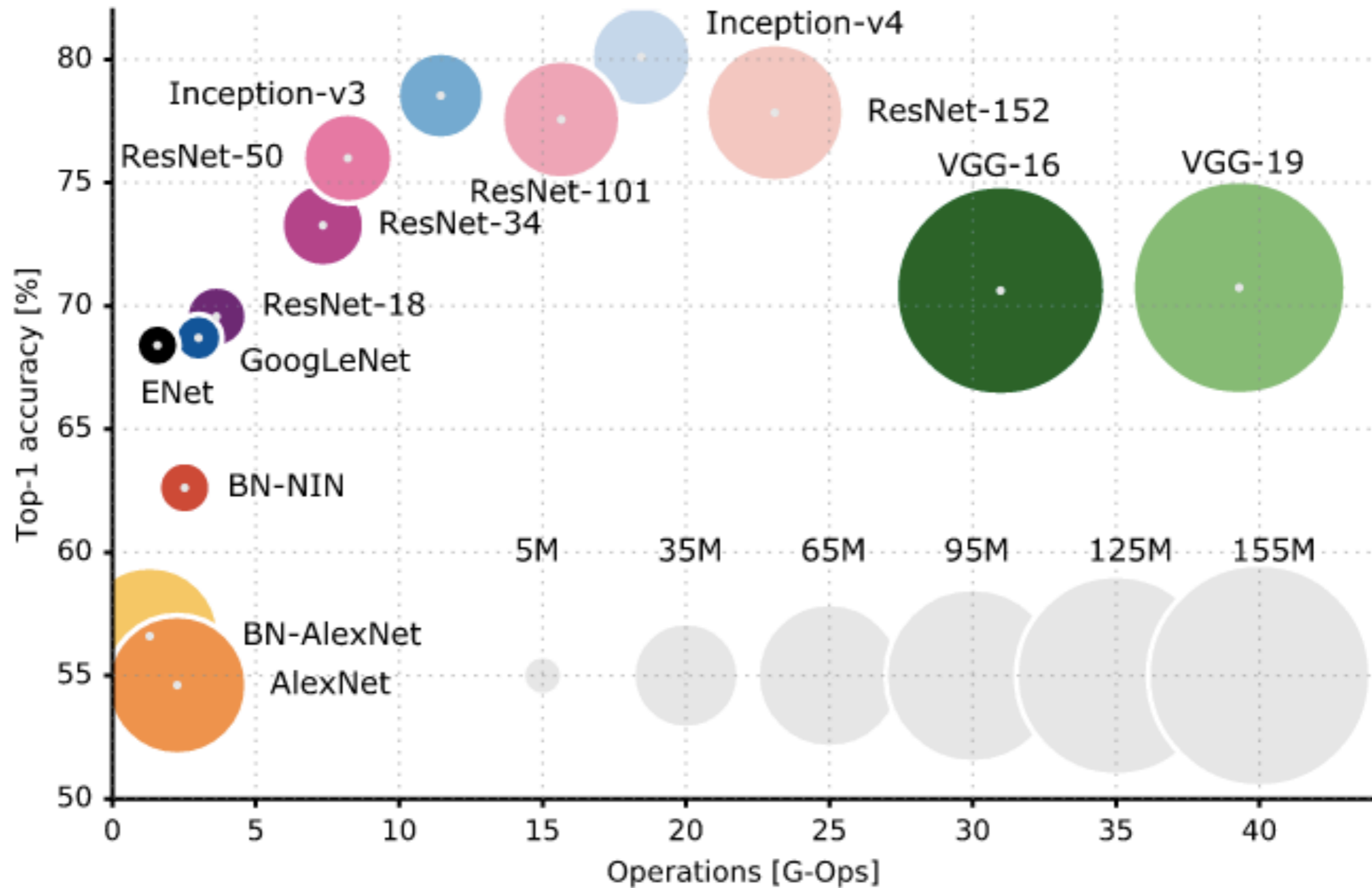
IM  GENET



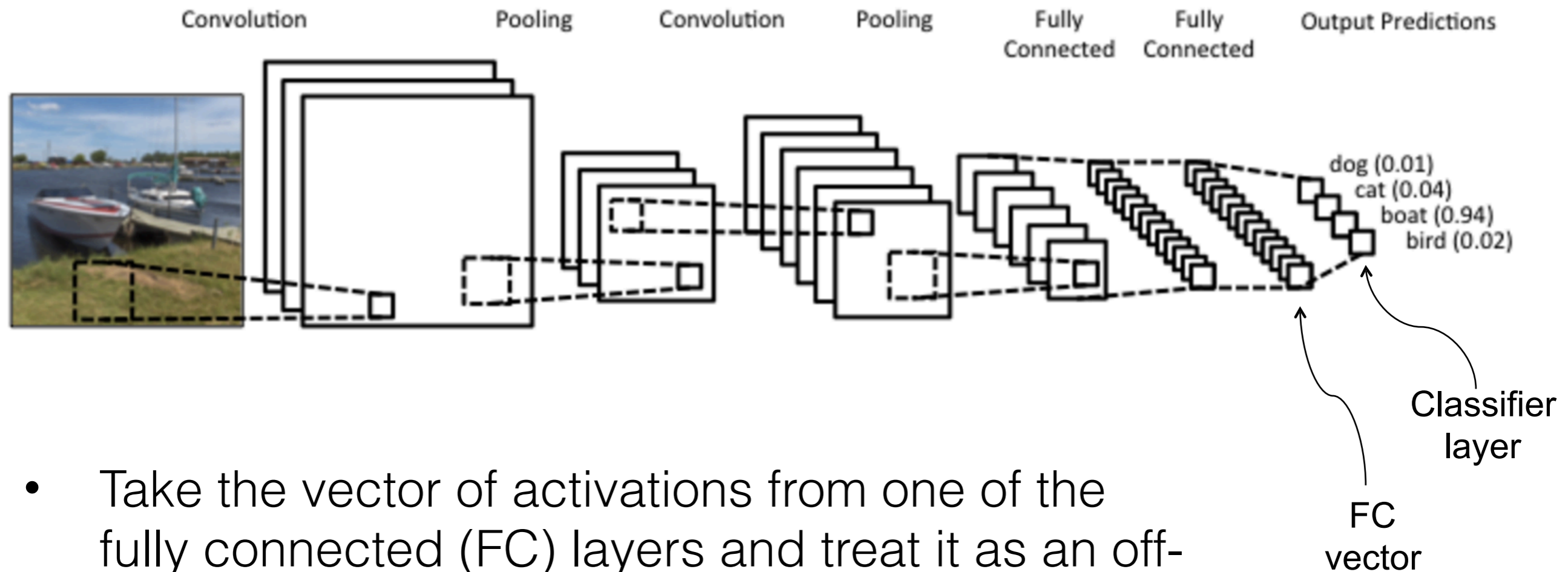
- ~14 million labeled images, 20k classes
- Images gathered from Internet
- Human labels via Amazon MTurk
- ImageNet Large-Scale Visual Recognition Challenge (ILSVRC):
1.2 million training images, 1000 classes
- Starting with 2015, image classification is not part of ILSVRC challenge (but people continue to benchmark on the data)

www.image-net.org/challenges/LSVRC/

Comparing architectures



How to use a trained network for a new task?



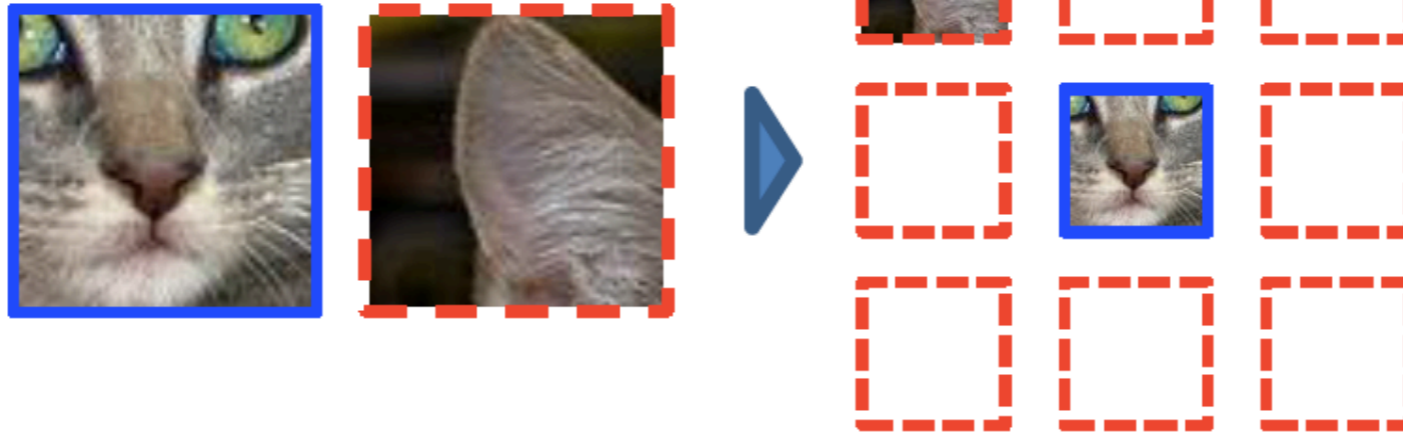
- Take the vector of activations from one of the fully connected (FC) layers and treat it as an off-the-shelf feature
- Train a new classifier layer on top of the FC layer
- Fine-tune the whole network.

More general pre-training than ImageNet?

- Goal: learn generic features, that can be useful for other tasks.
- Idea of “self supervised” learning: find auxiliary task which allows to learn useful features

Puzzle solving

Example:



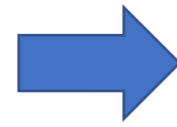
Question 1:



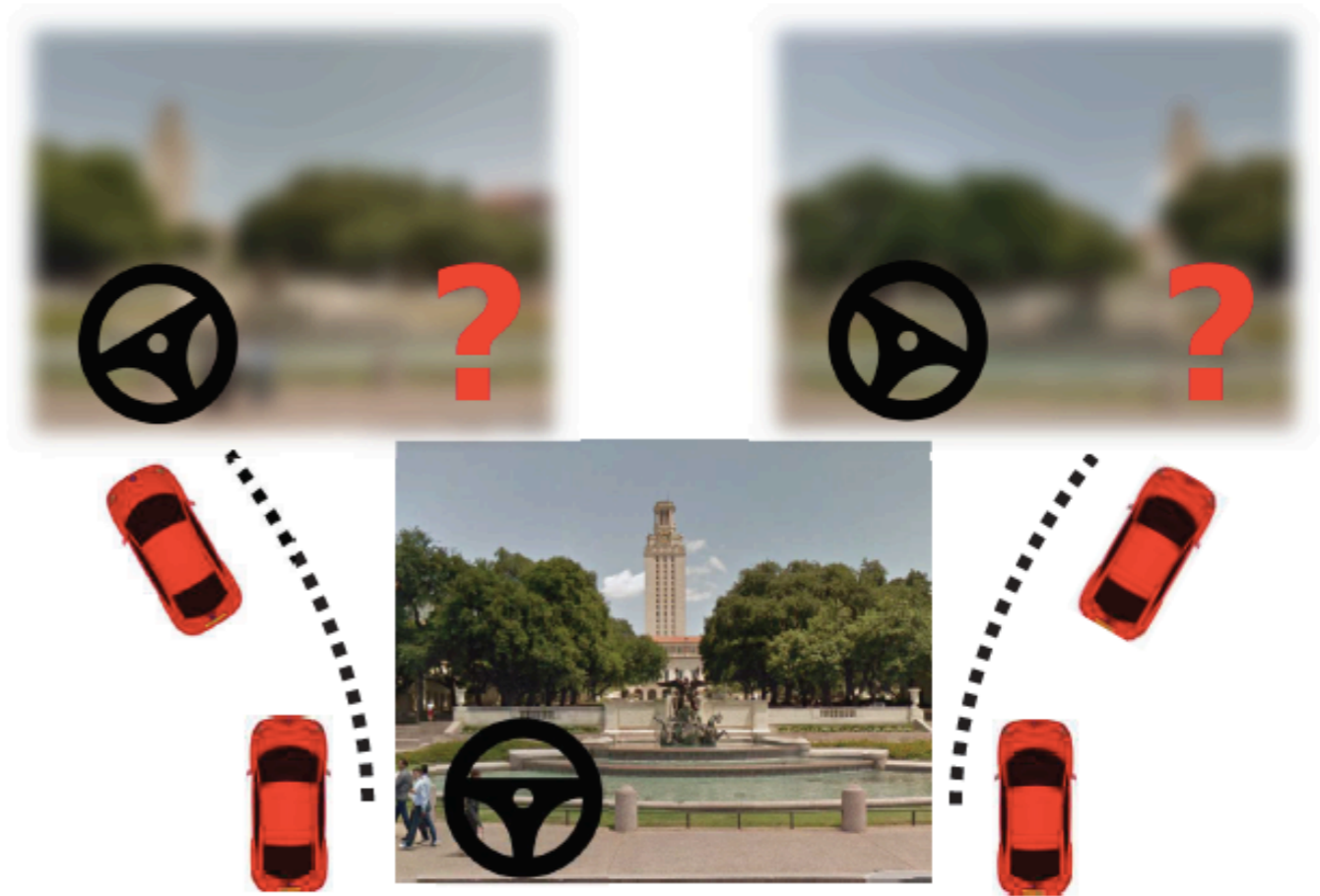
Question 2:



Image completion

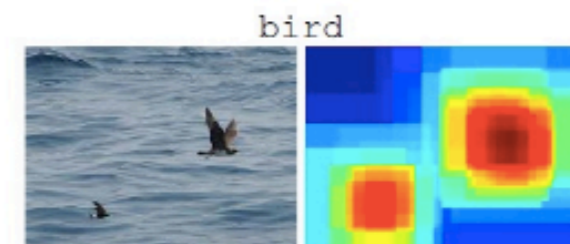
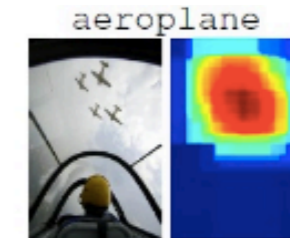
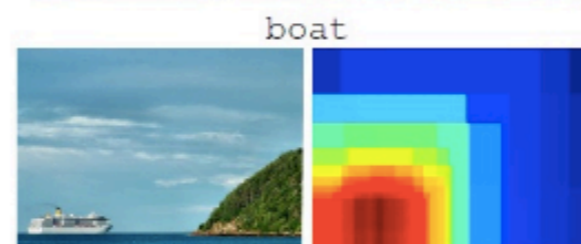
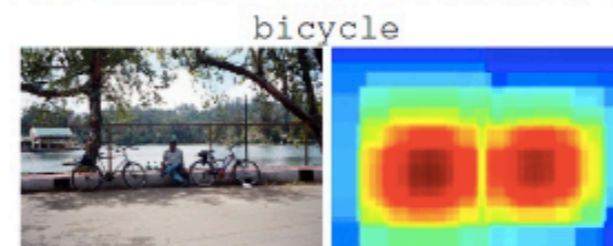
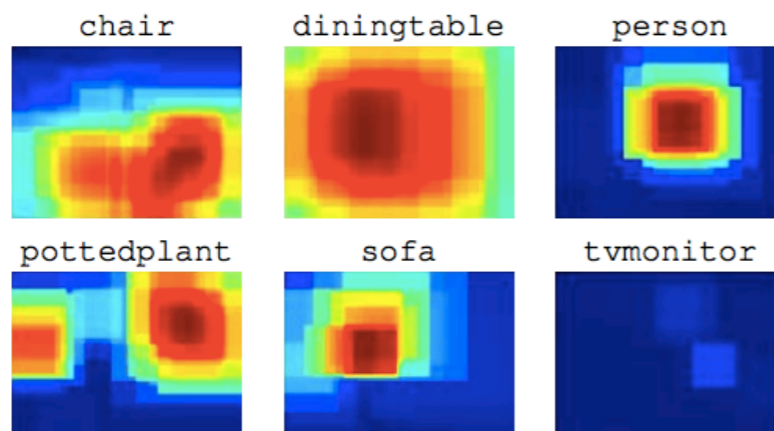


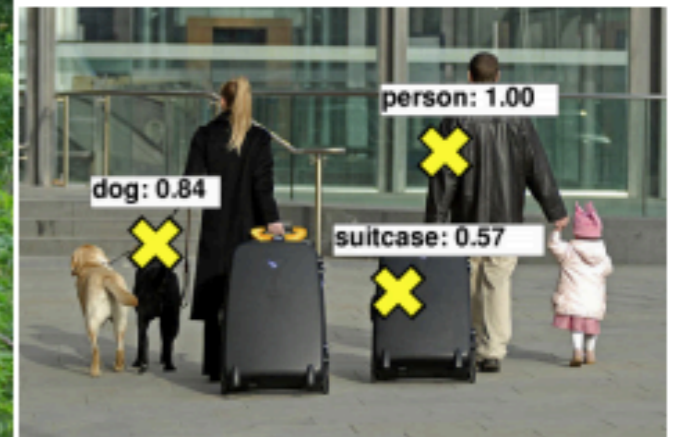
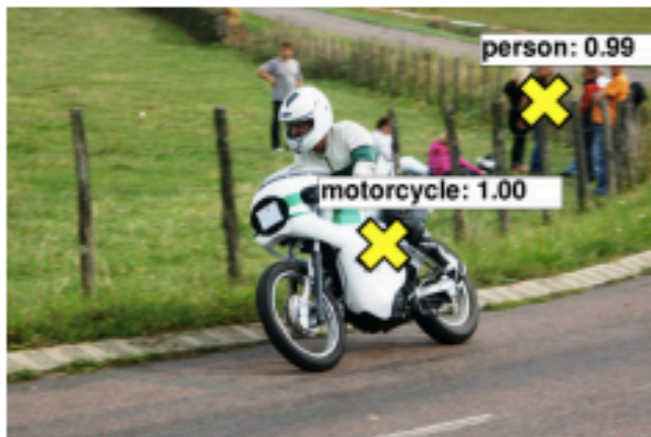
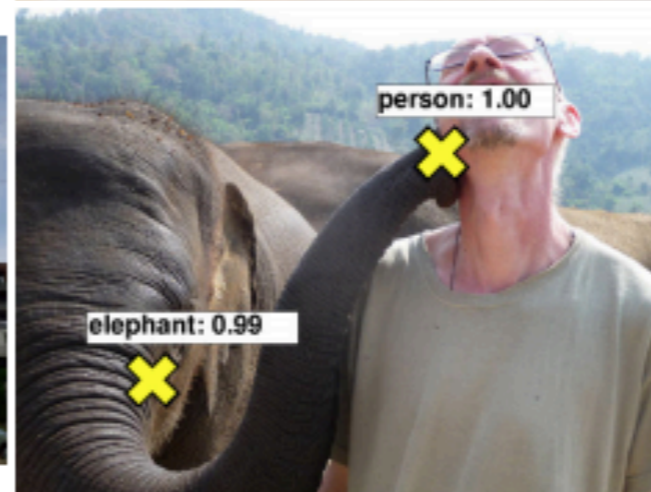
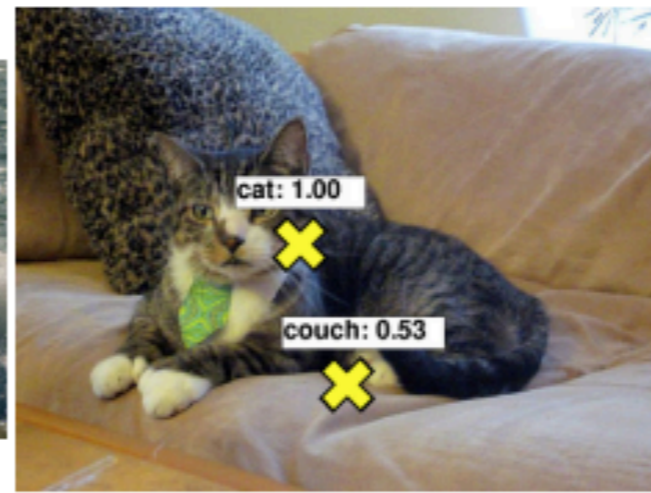
Predicting the future



Free localization from classification

- Apply densely a classification network and look at the results:

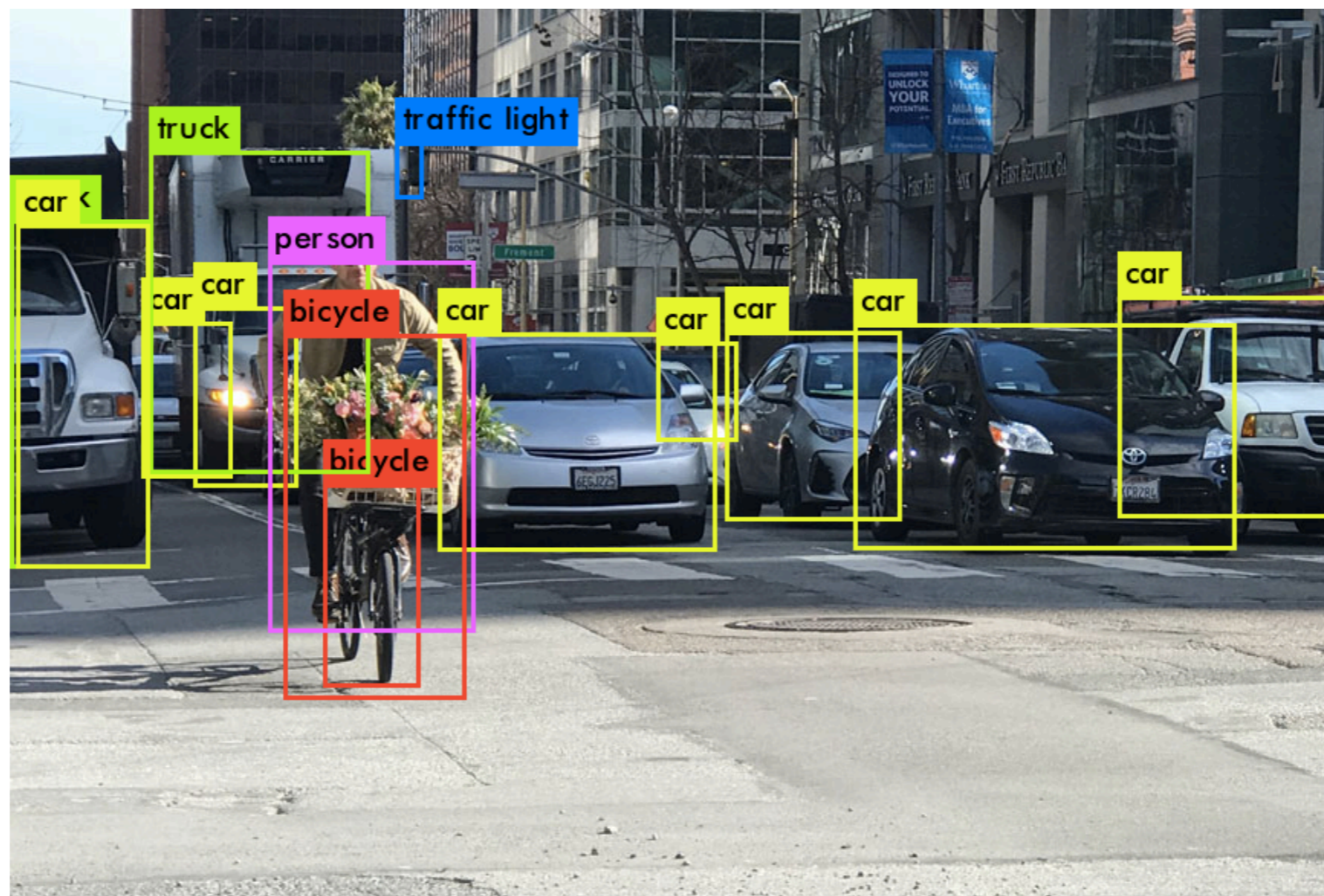




Neural networks for object detection

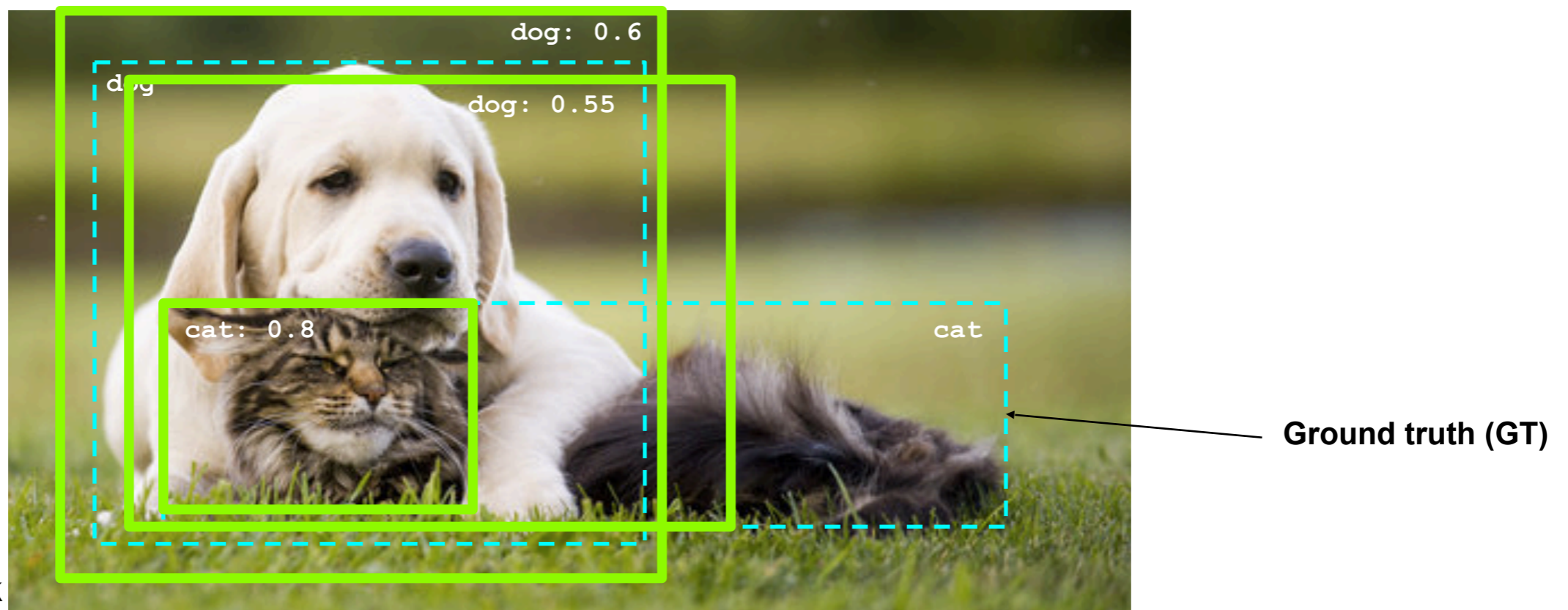
What are the challenges of object detection?

- Images may contain more than one class, multiple instances from the same class
- Bounding box localization
- Evaluation



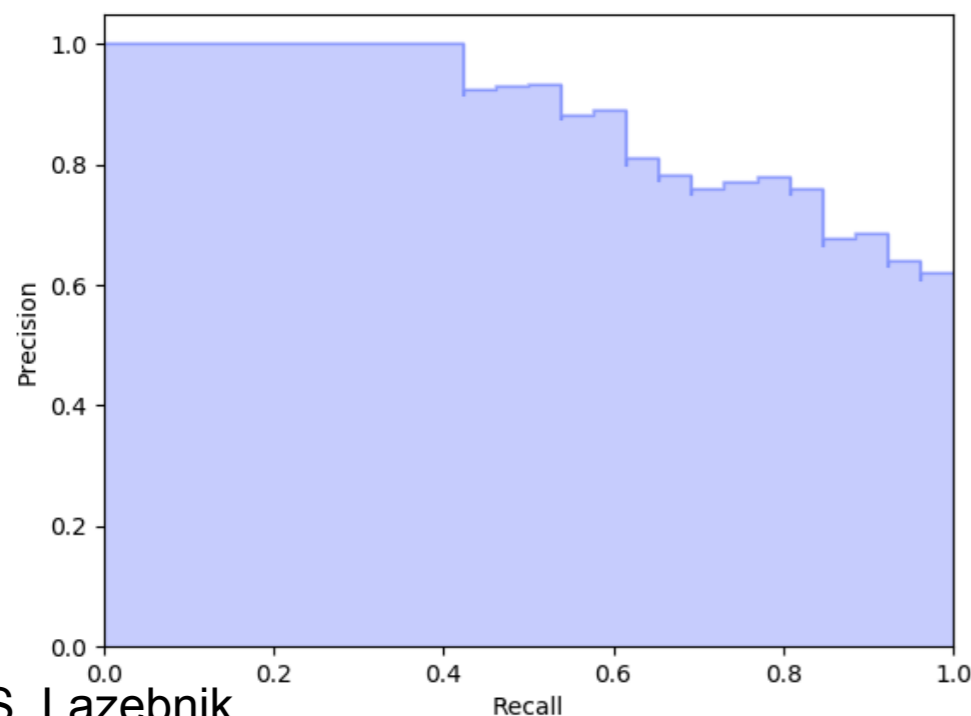
Object detection evaluation

- At test time, predict bounding boxes, class labels, and confidence scores
- For each detection, determine whether it is a true or false positive
 - PASCAL criterion: $\text{Area}(\text{GT} \cap \text{Det}) / \text{Area}(\text{GT} \cup \text{Det}) > 0.5$
 - For multiple detections of the same ground truth box, only one considered a true positive



Object detection evaluation

- At test time, predict bounding boxes, class labels, and confidence scores
- For each detection, determine whether it is a true or false positive
- For each class, plot **Recall-Precision curve** and compute **Average Precision** (area under the curve)
- Take mean of AP over classes to get **mAP**



Precision:

true positive detections /
total detections

Recall:

true positive detections /
total positive test instances

Simple approach: Sliding window detection



- Slide a window across the image and evaluate a detection model at each location
 - Thousands of windows to evaluate: efficiency and low false positive rates are essential
 - Difficult to extend to a large range of scales, aspect ratios

Histograms of oriented gradients (HOG)

- Partition image into blocks and compute histogram of gradient orientations in each block

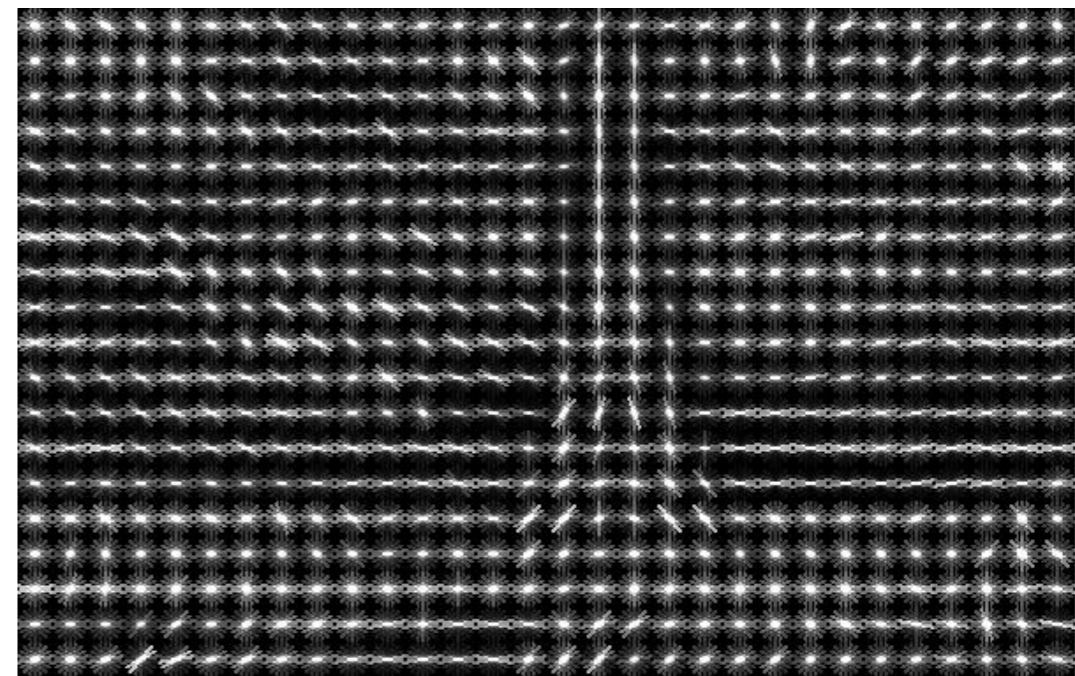
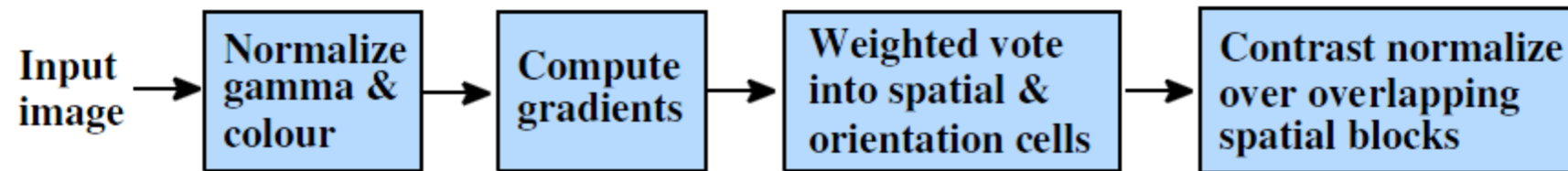


Image credit: N. Snavely

N. Dalal and B. Triggs, [Histograms of Oriented Gradients for Human Detection](#), CVPR 2005

Pedestrian detection with HOG

- Train a pedestrian template using a linear support vector machine

positive training examples



negative training examples

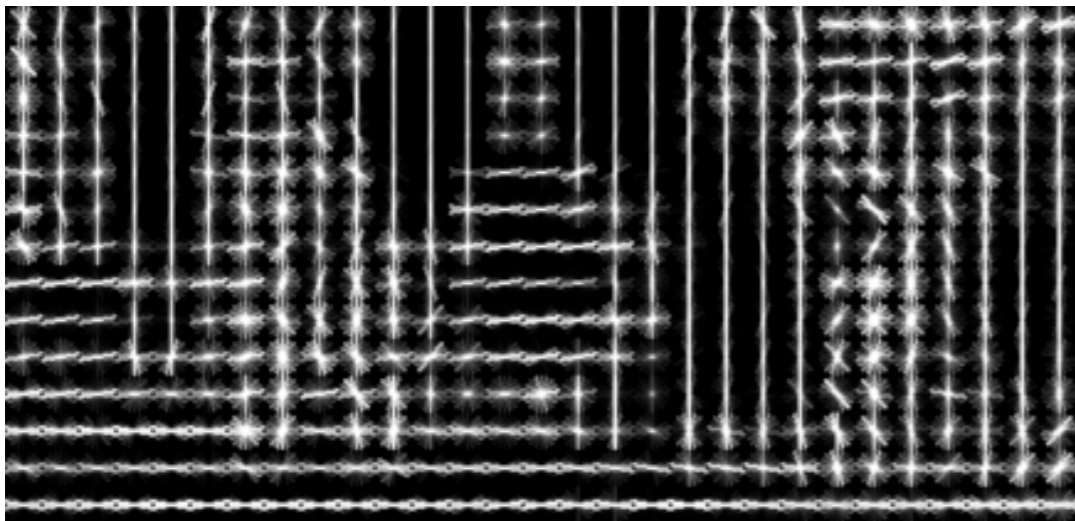


N. Dalal and B. Triggs, [Histograms of Oriented Gradients for Human Detection](#), CVPR 2005

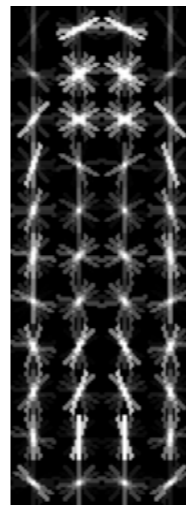
Pedestrian detection with HOG

- Train a pedestrian template using a linear support vector machine
- At test time, convolve feature map with template
- Find local maxima of response
- For multi-scale detection, repeat over multiple levels of a HOG *pyramid*

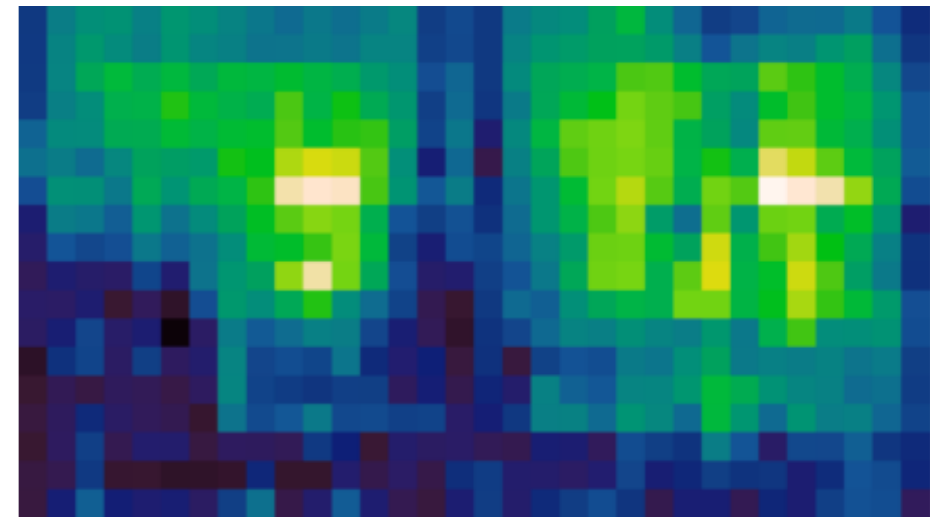
HOG feature map



Template



Detector response map



N. Dalal and B. Triggs, [Histograms of Oriented Gradients for Human Detection](#), CVPR 2005

Example detections



Discriminative part-based models

- Single rigid template usually not enough to represent a category
 - Many objects (e.g. humans) are articulated, or have parts that can vary in configuration

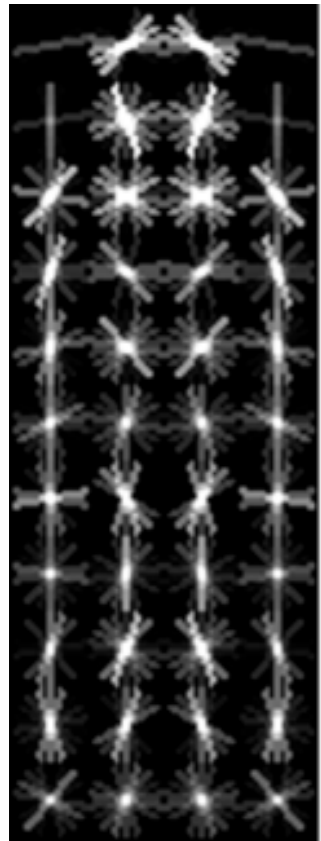


- Many object categories look very different from different viewpoints, or from instance to instance

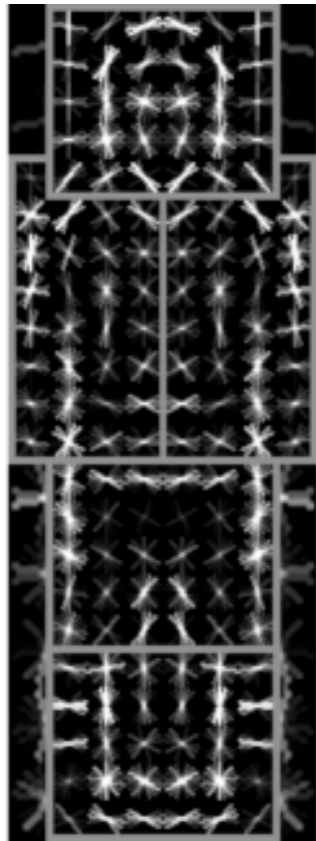


Discriminative part-based models

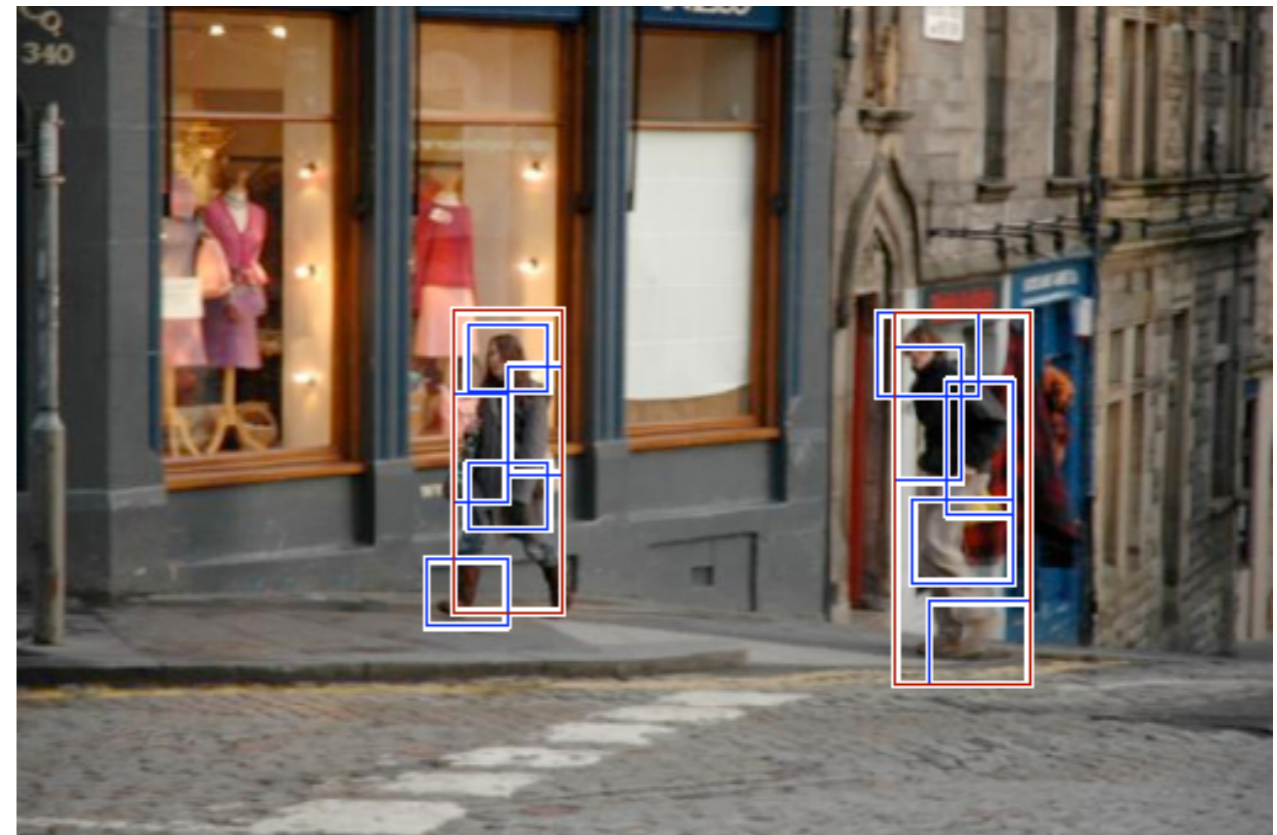
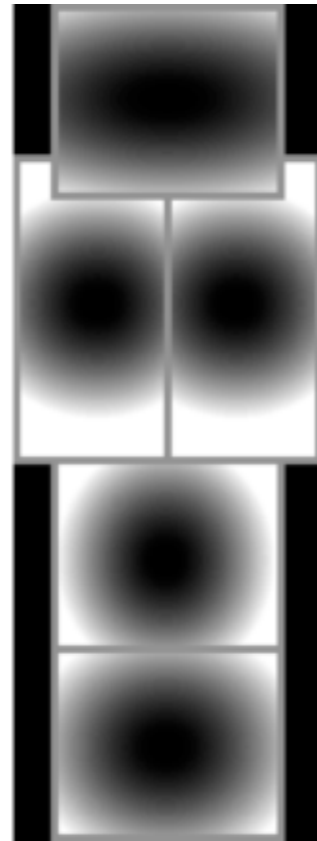
Root filter



Part filters



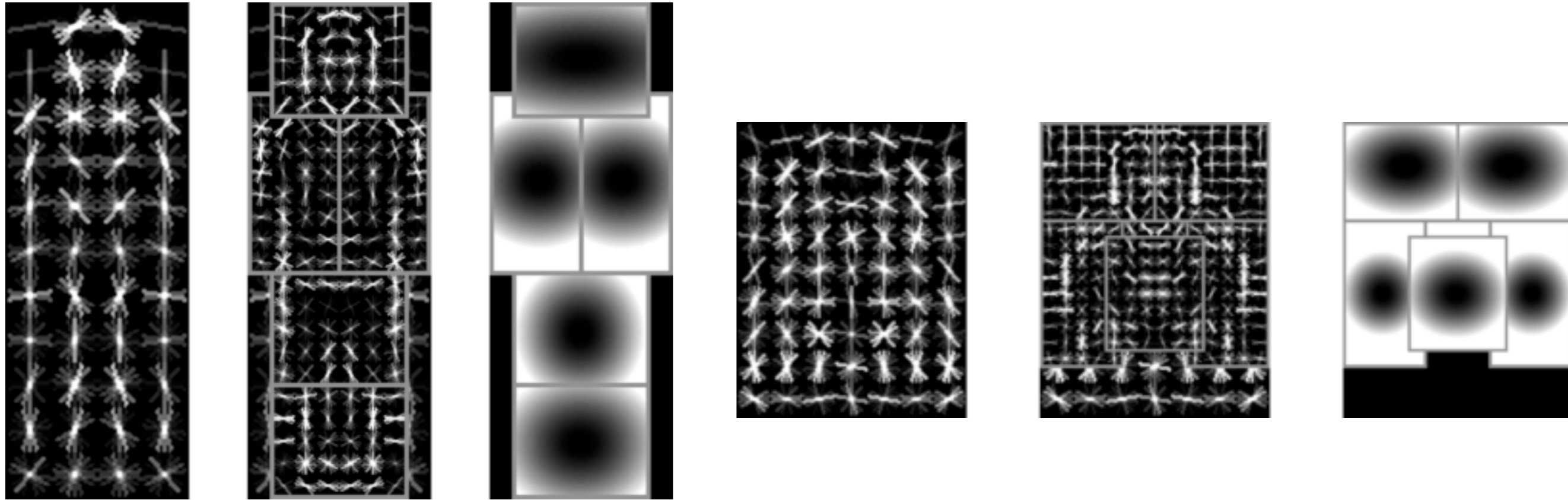
Deformation weights



P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, [Object Detection with Discriminatively Trained Part Based Models](#), PAMI 32(9), 2010

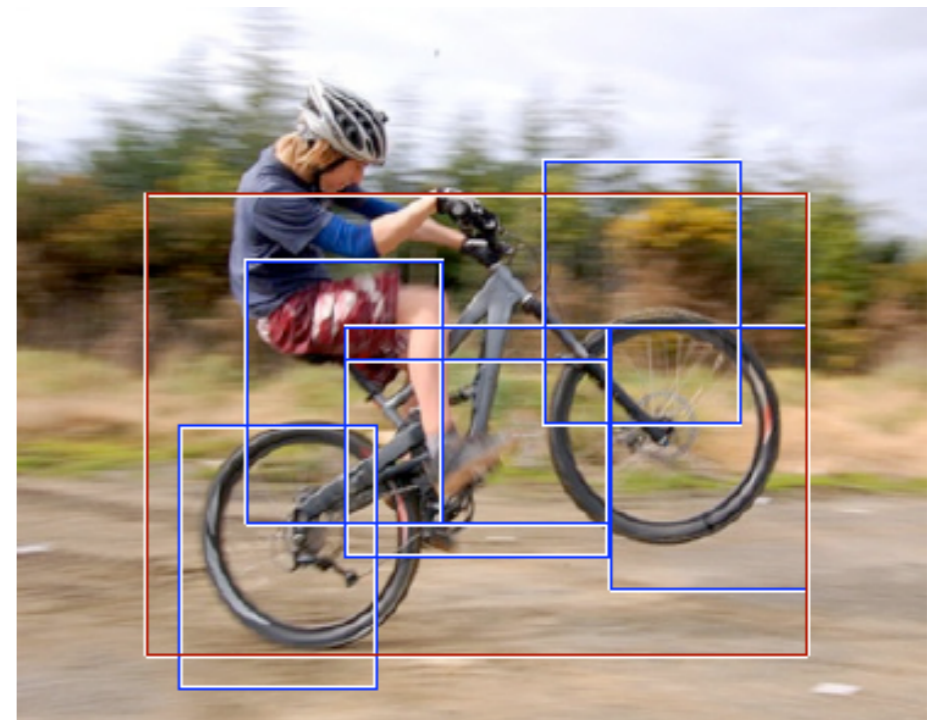
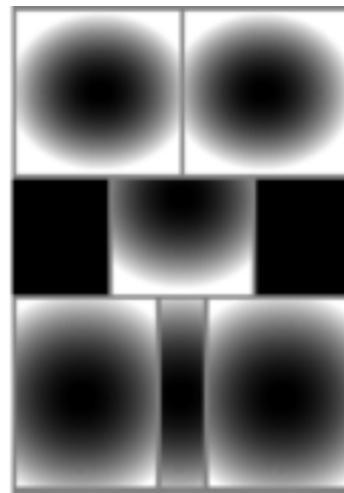
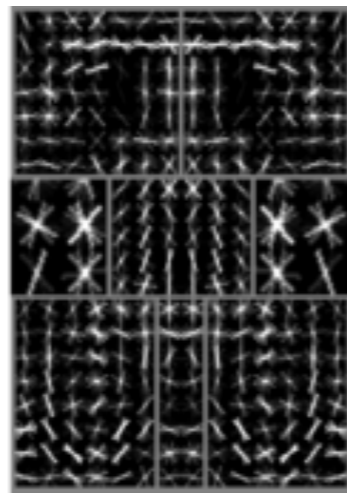
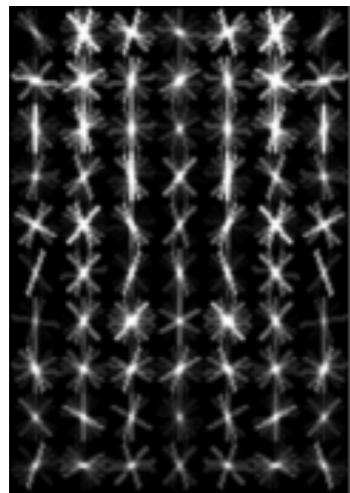
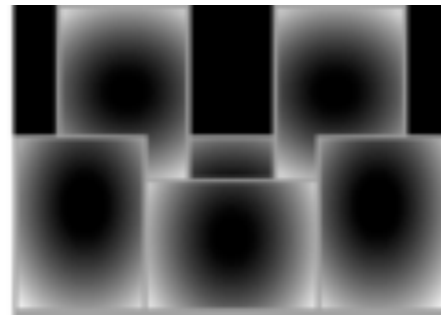
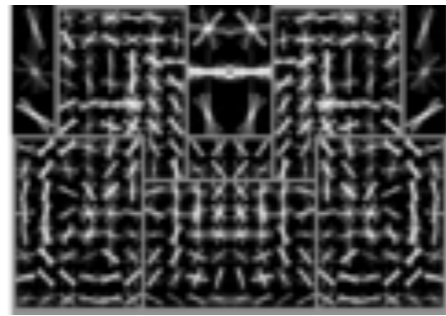
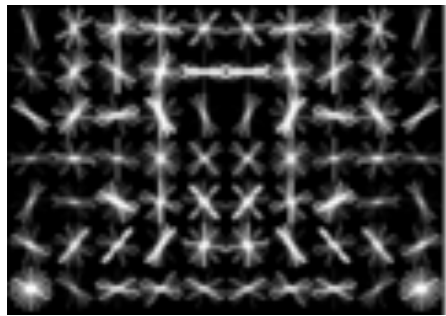
Discriminative part-based models

Multiple components



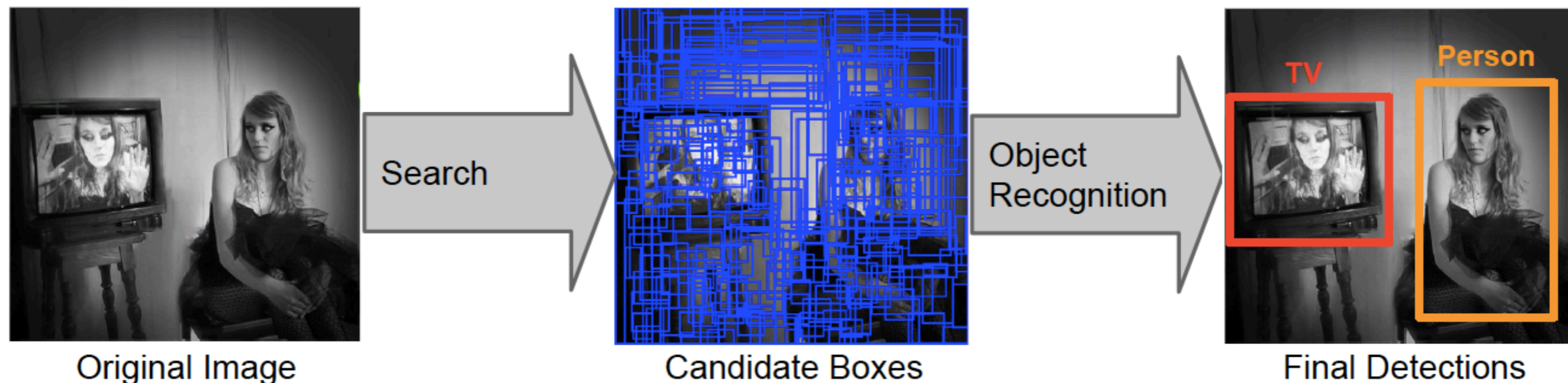
P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, [Object Detection with Discriminatively Trained Part Based Models](#), PAMI 32(9), 2010

Discriminative part-based models



P. Felzenszwalb, R. Girshick, D. McAllester, D. Ramanan, [Object Detection with Discriminatively Trained Part Based Models](#), PAMI 32(9), 2010

Conceptual approach: Proposal-driven detection



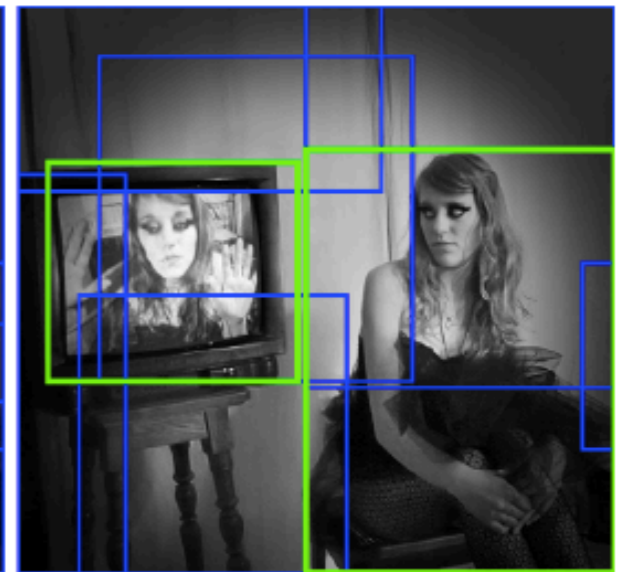
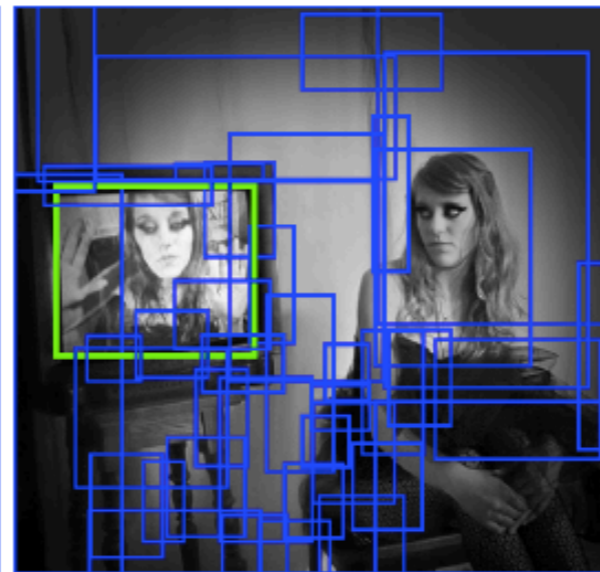
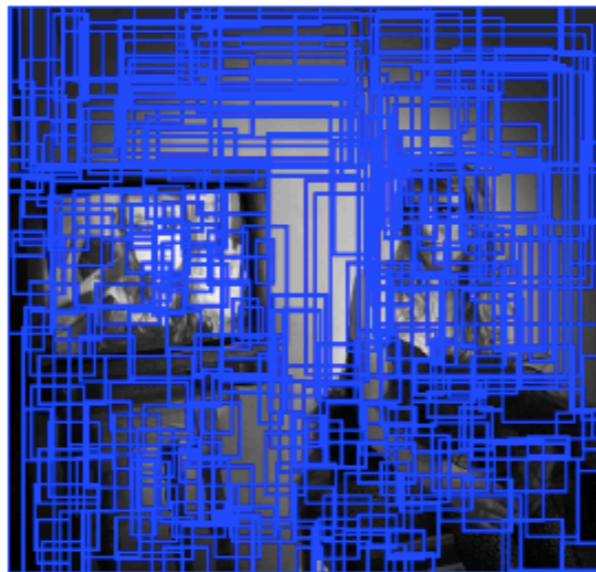
- Generate and evaluate a few hundred *region proposals*
 - Proposal mechanism can take advantage of low-level *perceptual organization* cues
 - Proposal mechanism can be category-specific or category-independent, hand-crafted or trained
 - Classifier can be slower but more powerful

Selective search for detection

- Use hierarchical segmentation: start with small *superpixels* and merge based on diverse cues

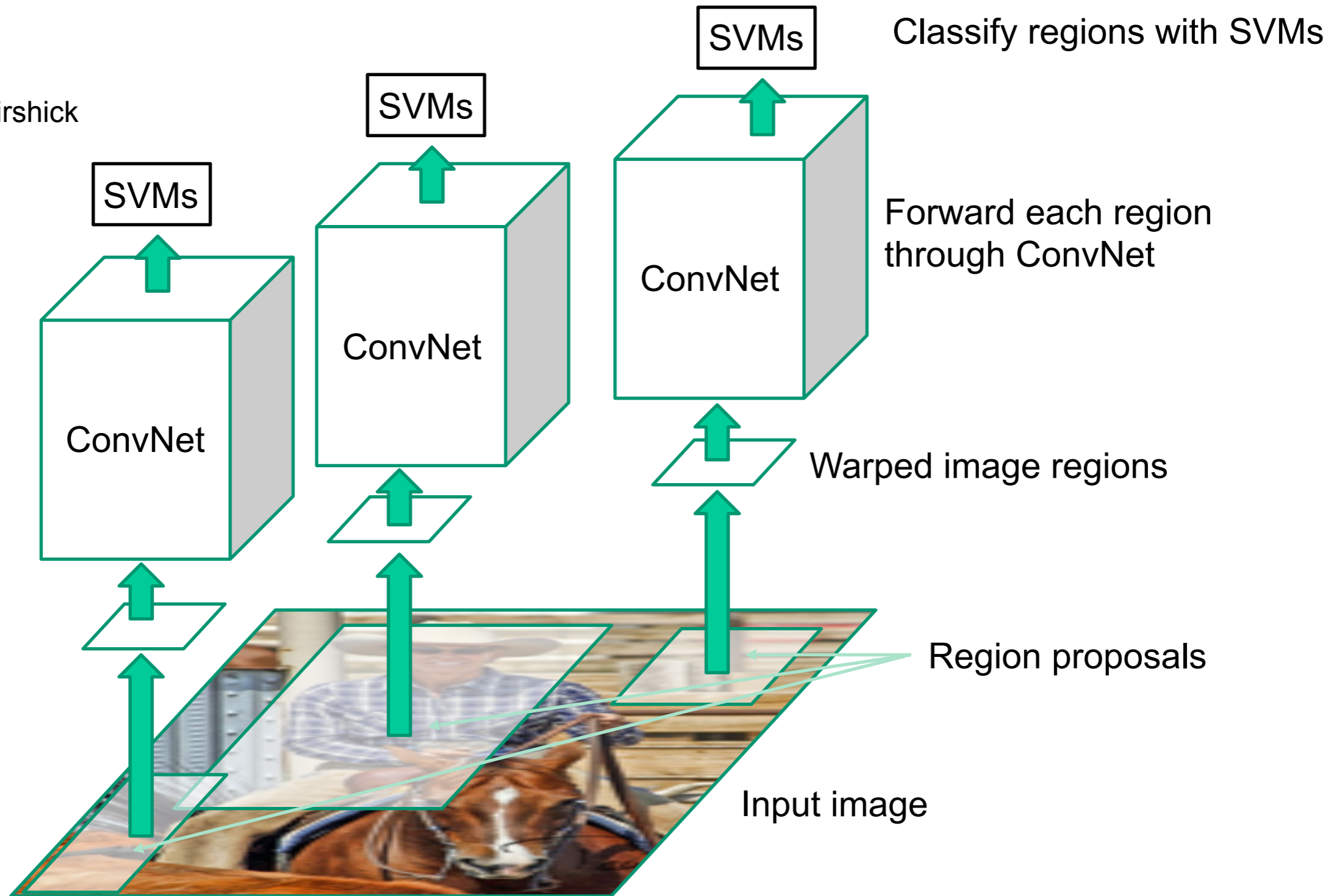


Input Image

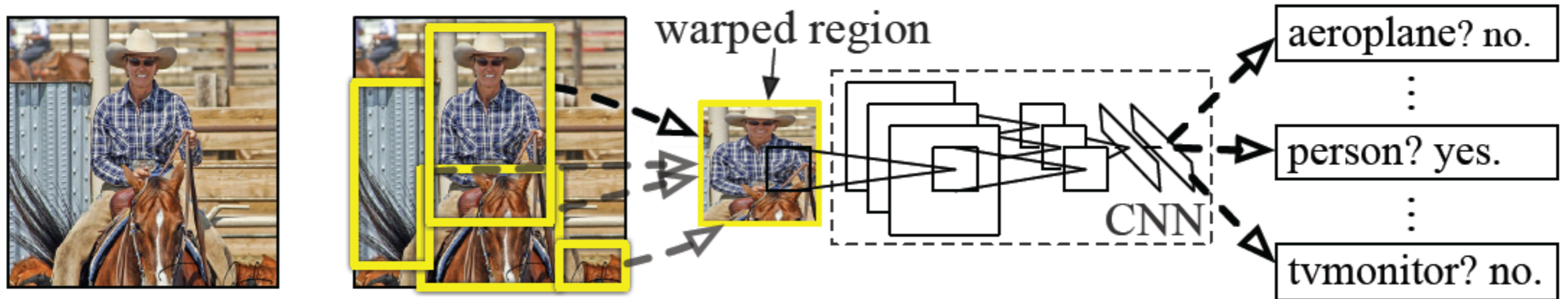


R-CNN: Region proposals + CNN features

Source: R. Girshick



R-CNN details

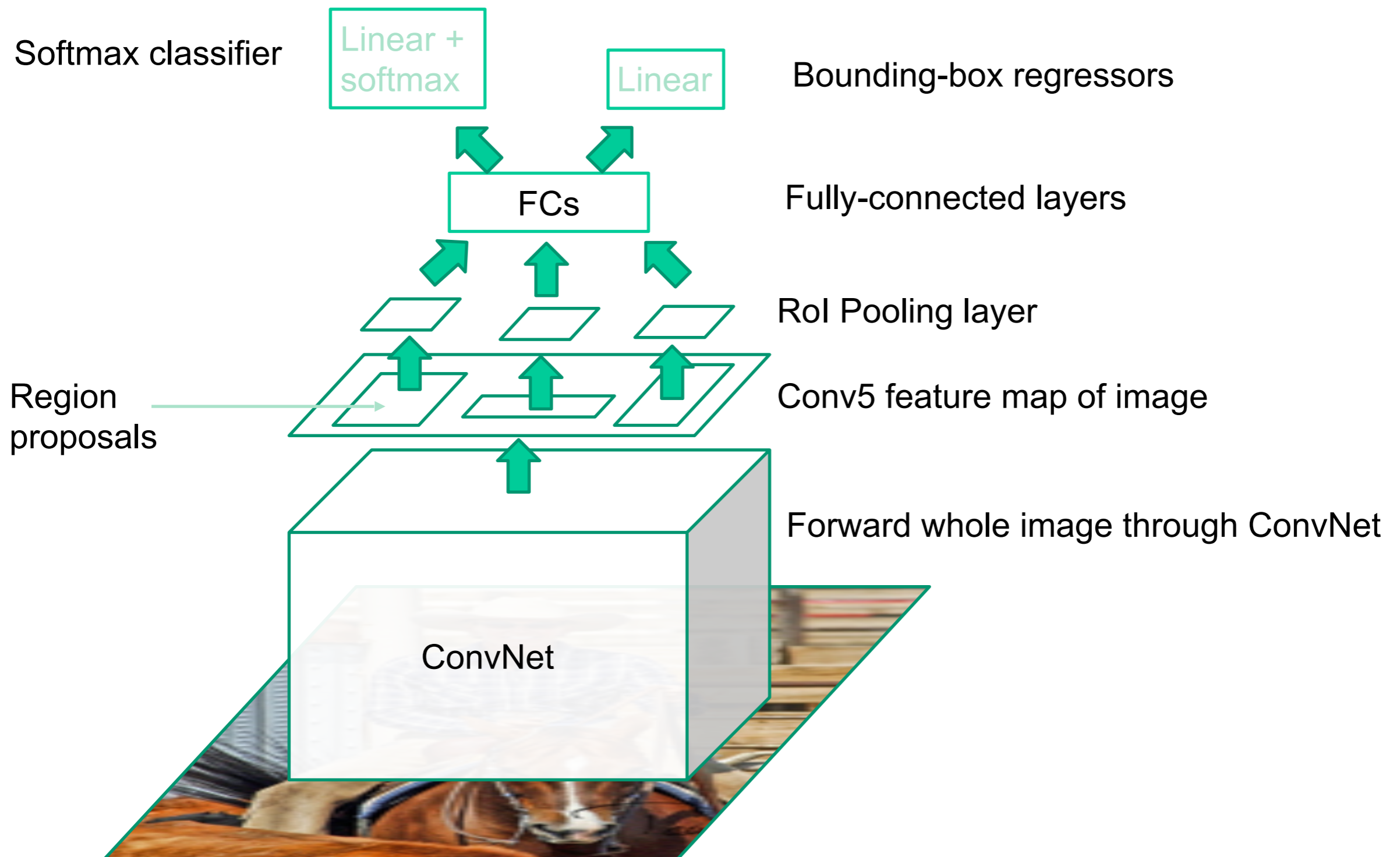


- **Regions:** ~2000 Selective Search proposals
- **Network:** AlexNet *pre-trained* on ImageNet (1000 classes), *fine-tuned* on PASCAL (21 classes)
- **Final detector:** warp proposal regions, extract fc7 network activations (4096 dimensions), classify with linear SVM
- **Bounding box regression** to refine box locations
- **Performance:** mAP of **53.7%** on PASCAL 2010 (vs. **35.1%** for Selective Search and **33.4%** for Deformable Part Models)

R-CNN pros and cons

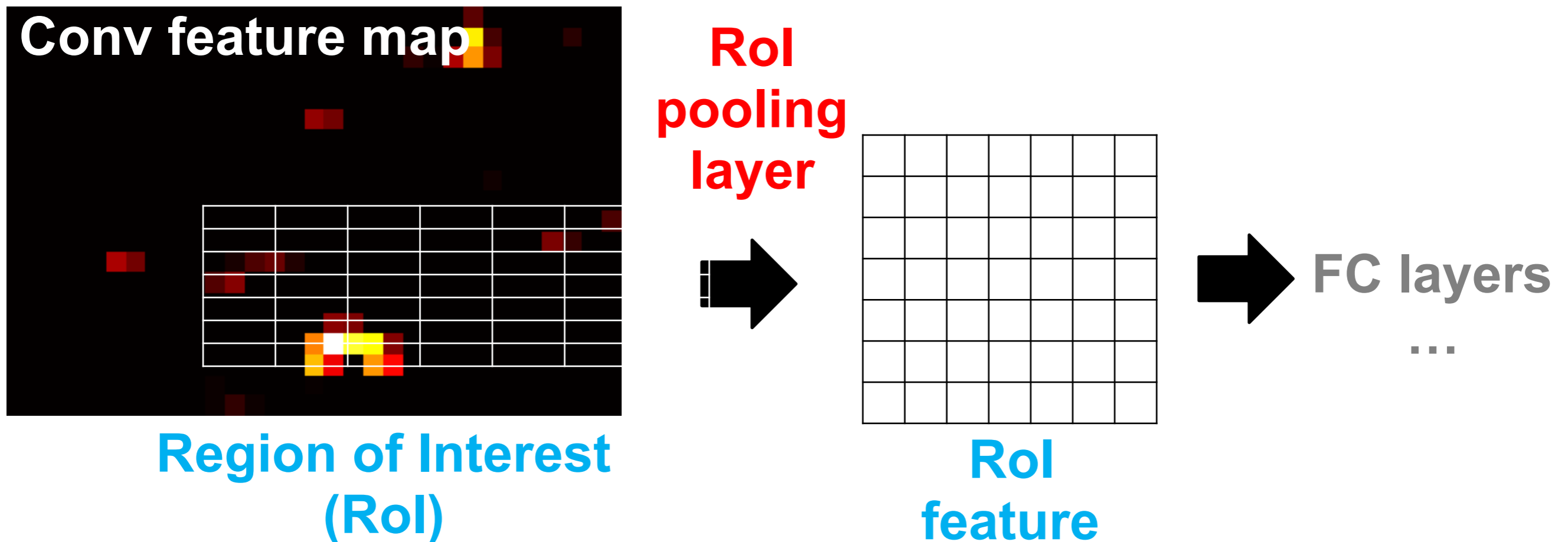
- Pros
 - Accurate!
 - Any deep architecture can immediately be “plugged in”
- Cons
 - Not a single end-to-end system
 - Fine-tune network with softmax classifier (log loss)
 - Train post-hoc linear SVMs (hinge loss)
 - Train post-hoc bounding-box regressions (least squares)
 - Training is slow (84h), takes a lot of disk space
 - 2000 CNN passes per image
 - Inference (detection) is slow (47s / image with VGG16)

Fast R-CNN



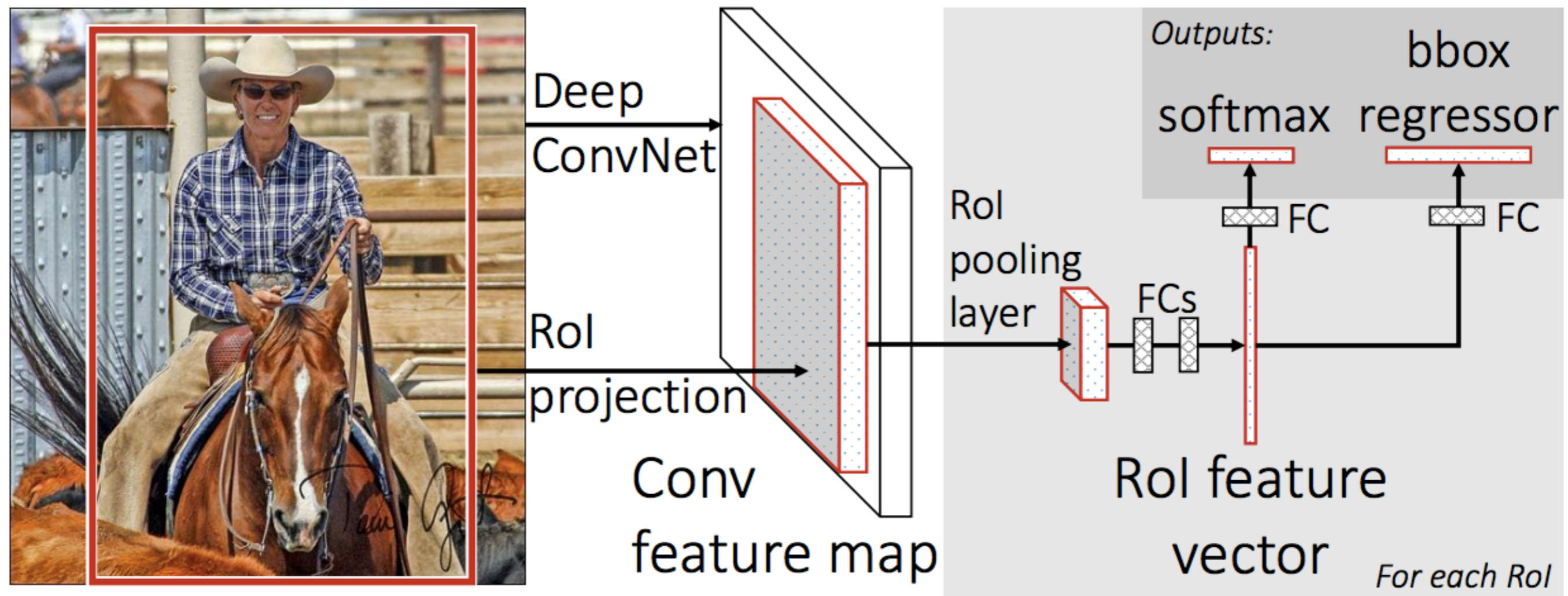
Roi pooling

- “Crop and resample” a fixed-size feature representing a region of interest out of the outputs of the last conv layer
 - Use nearest-neighbor interpolation of coordinates, max pooling

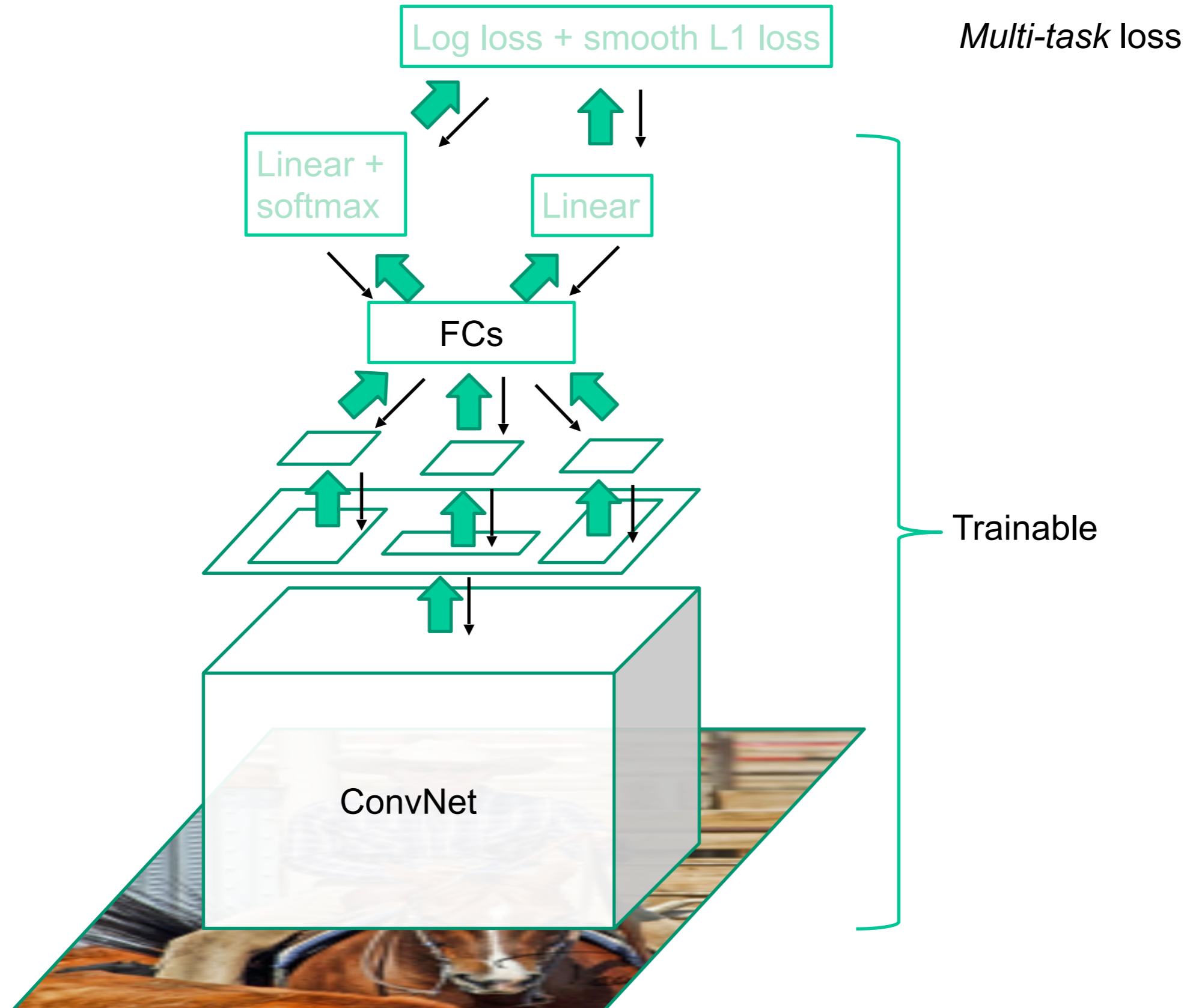


Prediction

- For each RoI, network predicts probabilities for $C+1$ classes (class 0 is background) and four bounding box offsets for C classes



Fast R-CNN training



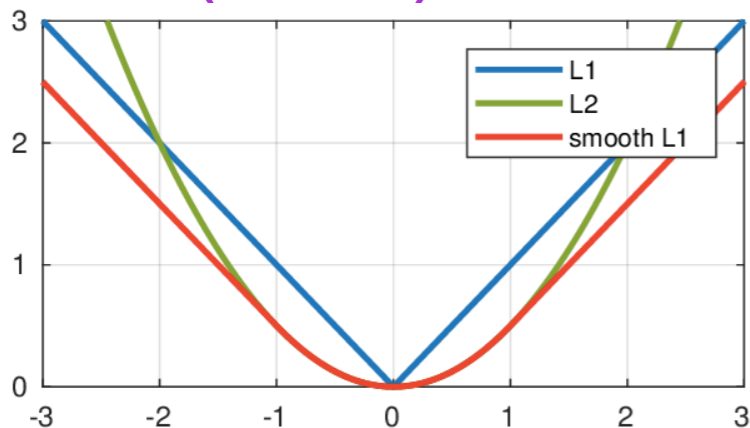
Multi-task loss

- Loss for ground truth class y , predicted class probabilities $P(y)$, ground truth box b , and predicted box \hat{b} :

$$L(y, P, b, \hat{b}) = \underbrace{-\log P(y)}_{\text{softmax loss}} + \lambda \mathbb{1}[y \geq 1] \underbrace{L_{\text{reg}}(b, \hat{b})}_{\text{regression loss}}$$

- Regression loss: *smooth L1 loss* on top of log space offsets relative to proposal

$$L_{\text{reg}}(b, \hat{b}) = \sum_i \text{smooth}_{L_1}(b_i - \hat{b}_i)$$



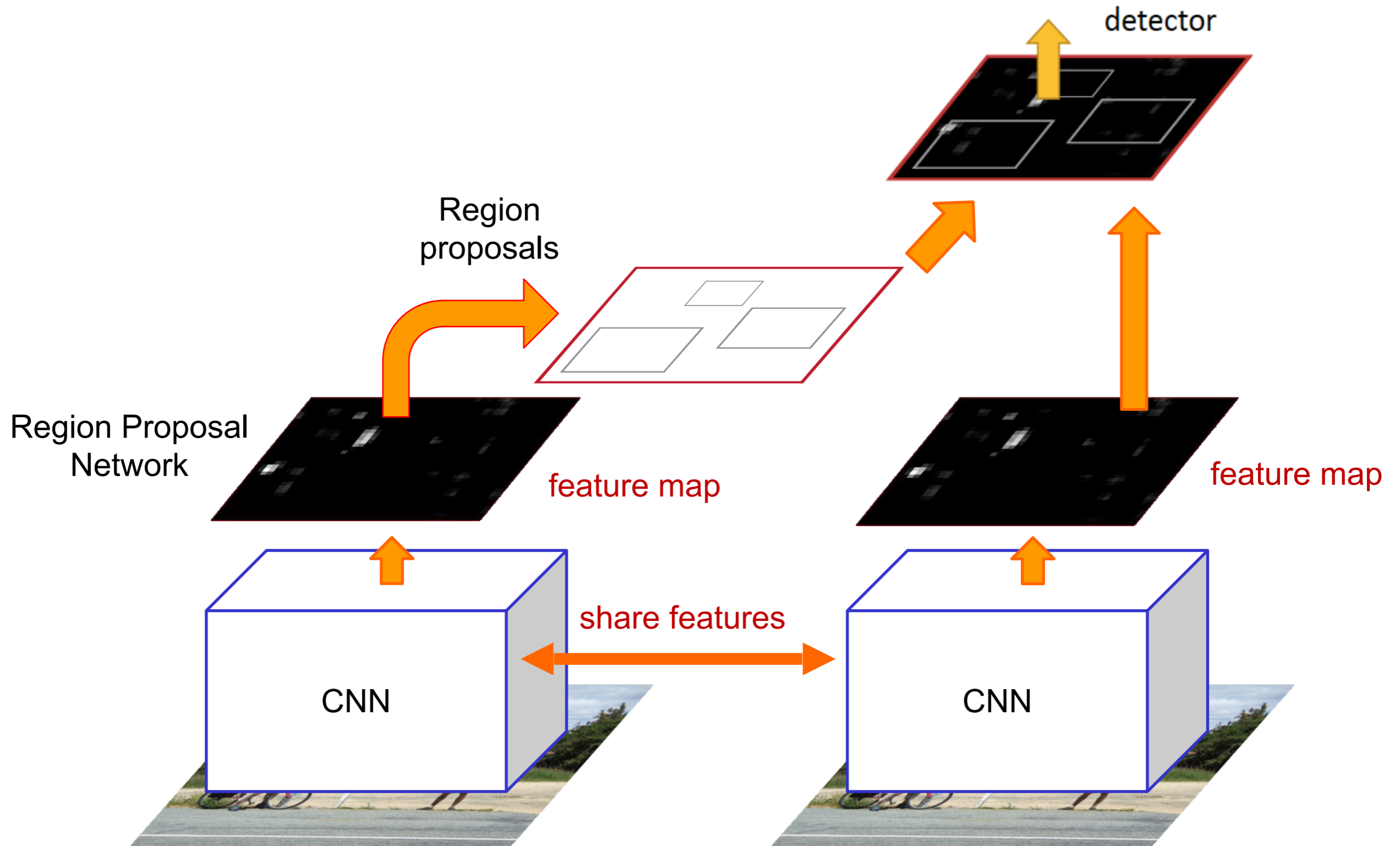
$$\text{smooth}_{L_1}(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases}$$

Fast R-CNN results

	Fast R-CNN	R-CNN
Train time (h)	9.5	84
- Speedup	8.8x	1x
Test time / image	0.32s	47.0s
Test speedup	146x	1x
mAP	66.9%	66.0%

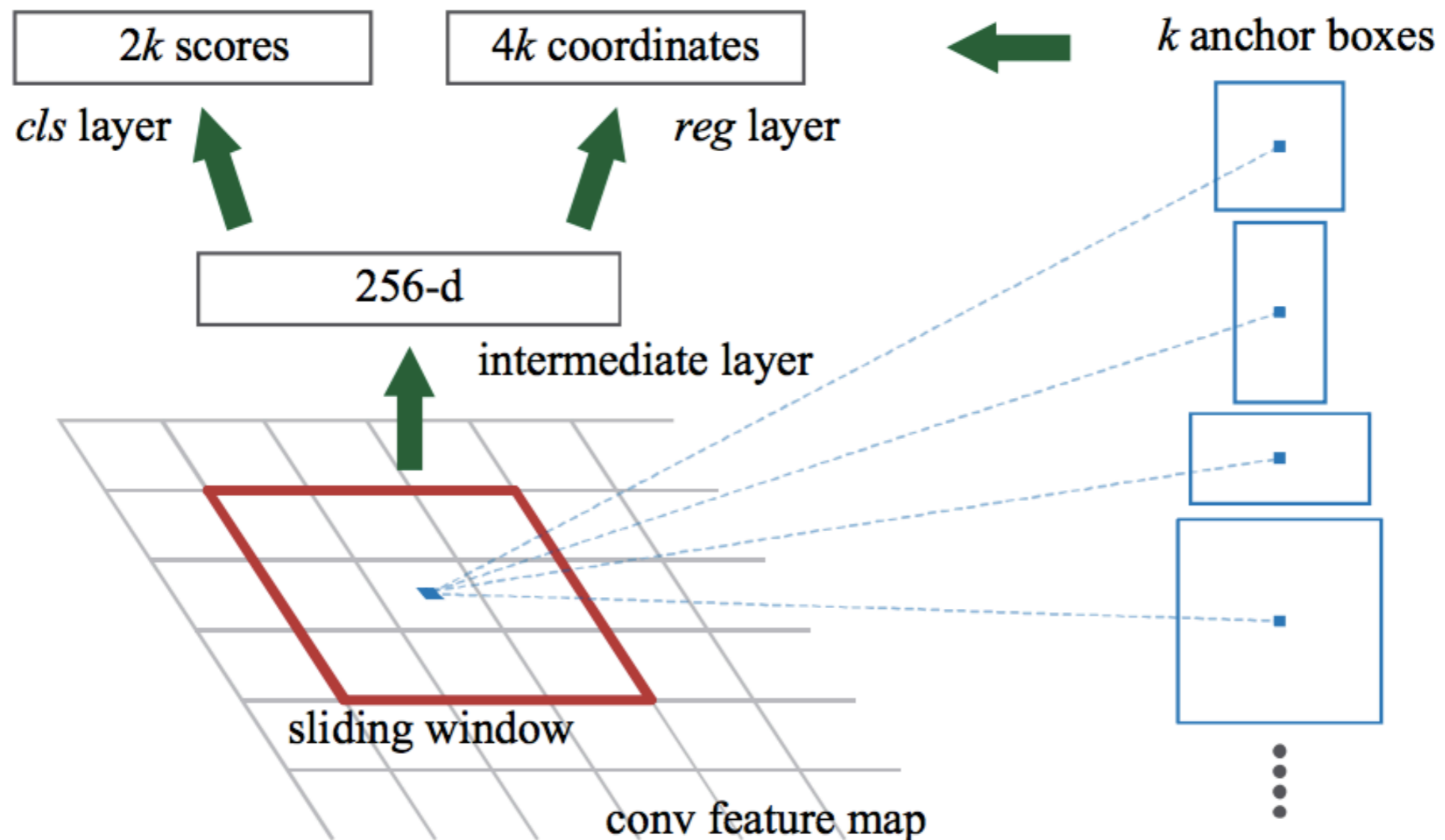
Timings exclude object proposal time, which is equal for all methods.
All methods use VGG16 from Simonyan and Zisserman.

Faster R-CNN

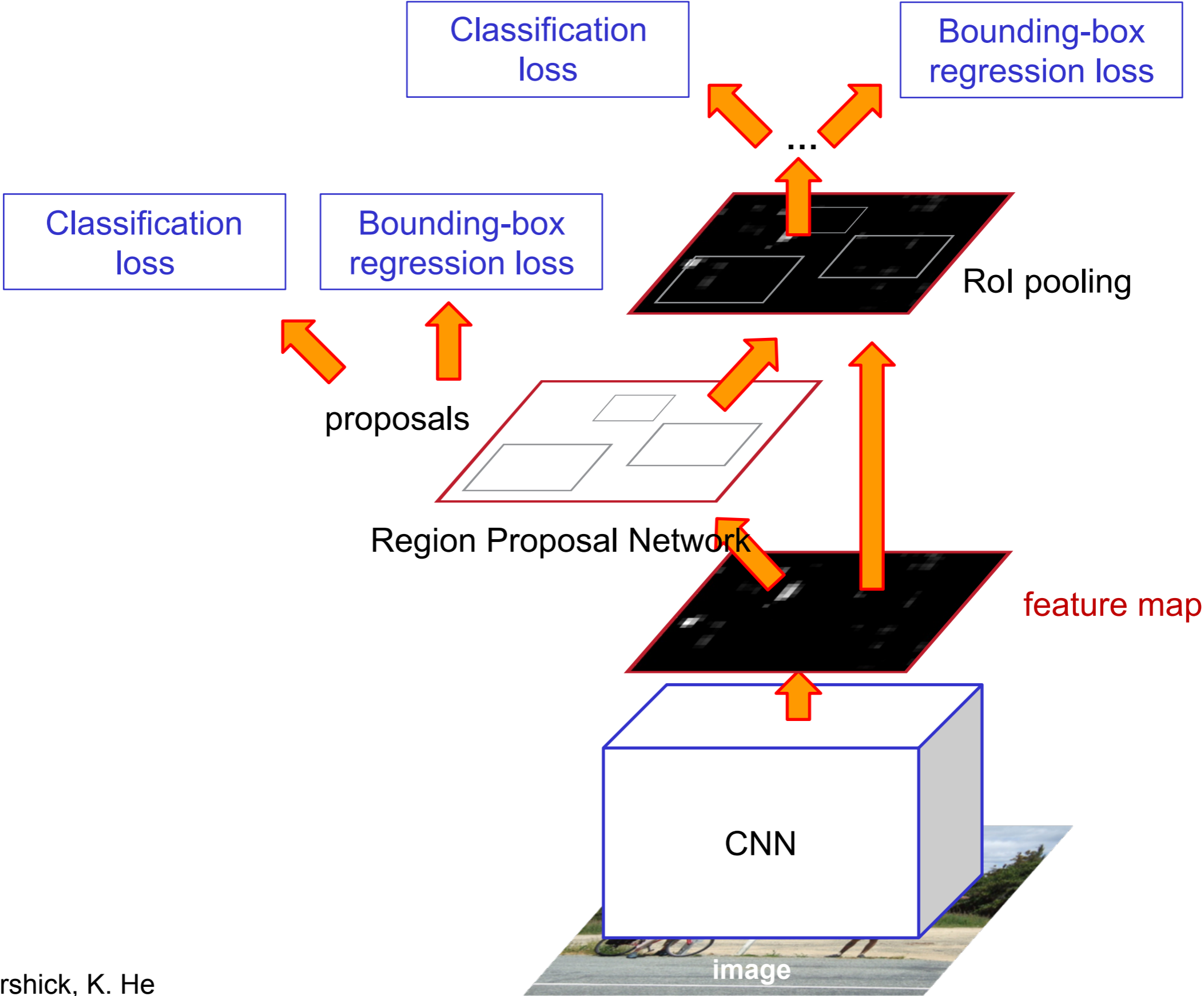


Region proposal network (RPN)

- Slide a small window (3x3) over the conv5 layer
 - Predict object/no object
 - Regress bounding box coordinates with reference to *anchors* (3 scales x 3 aspect ratios)



One network, four losses



Source: R. Girshick, K. He

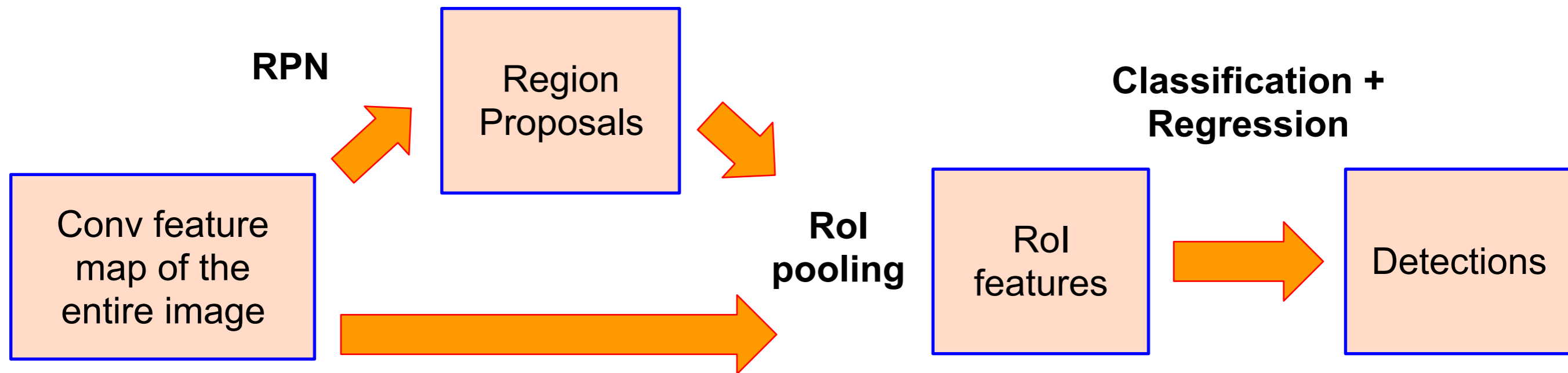
Faster R-CNN results

system	time	07 data	07+12 data
R-CNN	~50s	66.0	-
Fast R-CNN	~2s	66.9	70.0
Faster R-CNN	198ms	69.9	73.2

detection mAP on PASCAL VOC 2007, with VGG-16 pre-trained on ImageNet

Streamlined detection architectures

- The Faster R-CNN pipeline separates proposal generation and region classification:

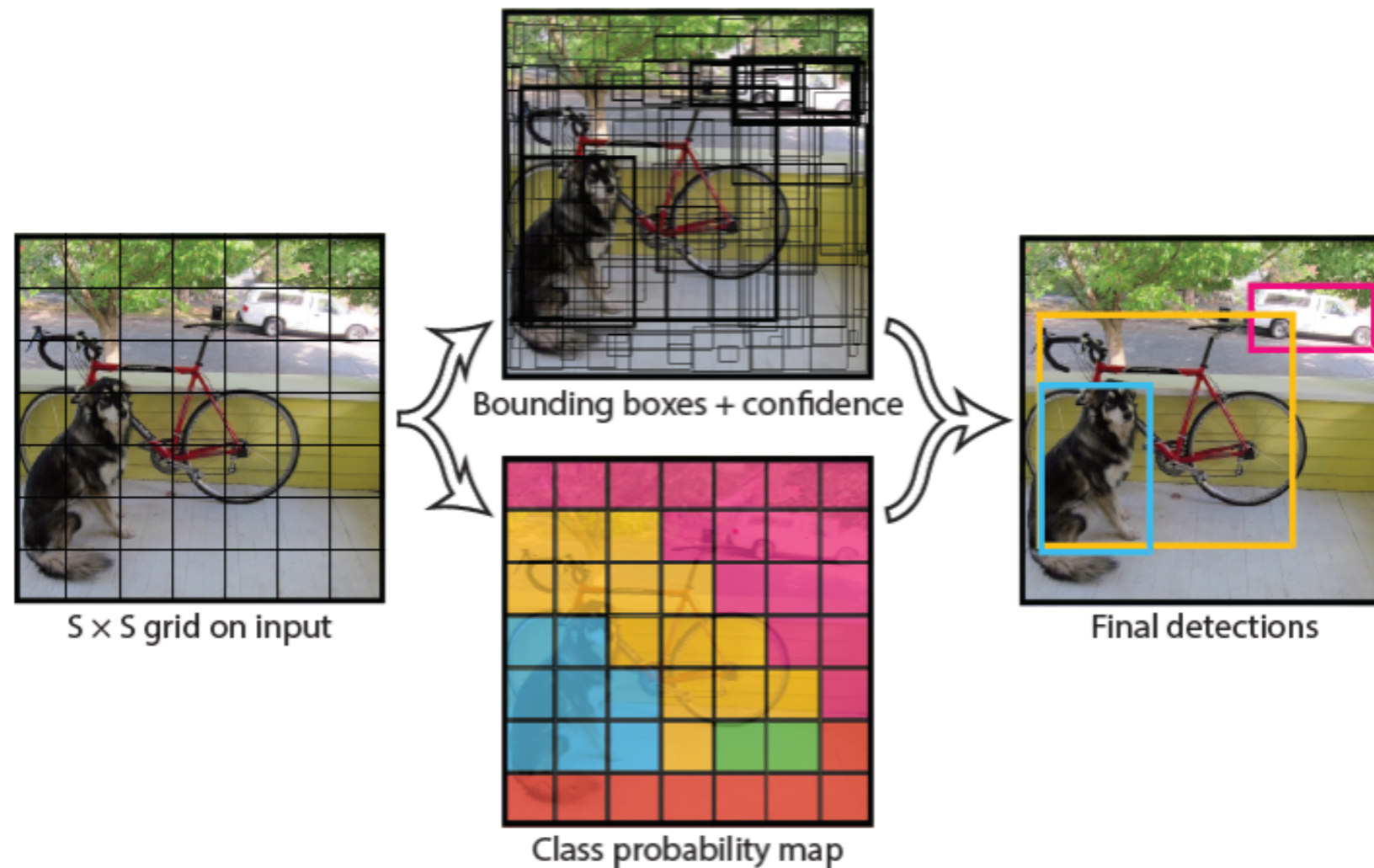


- Is it possible to do detection in one shot?



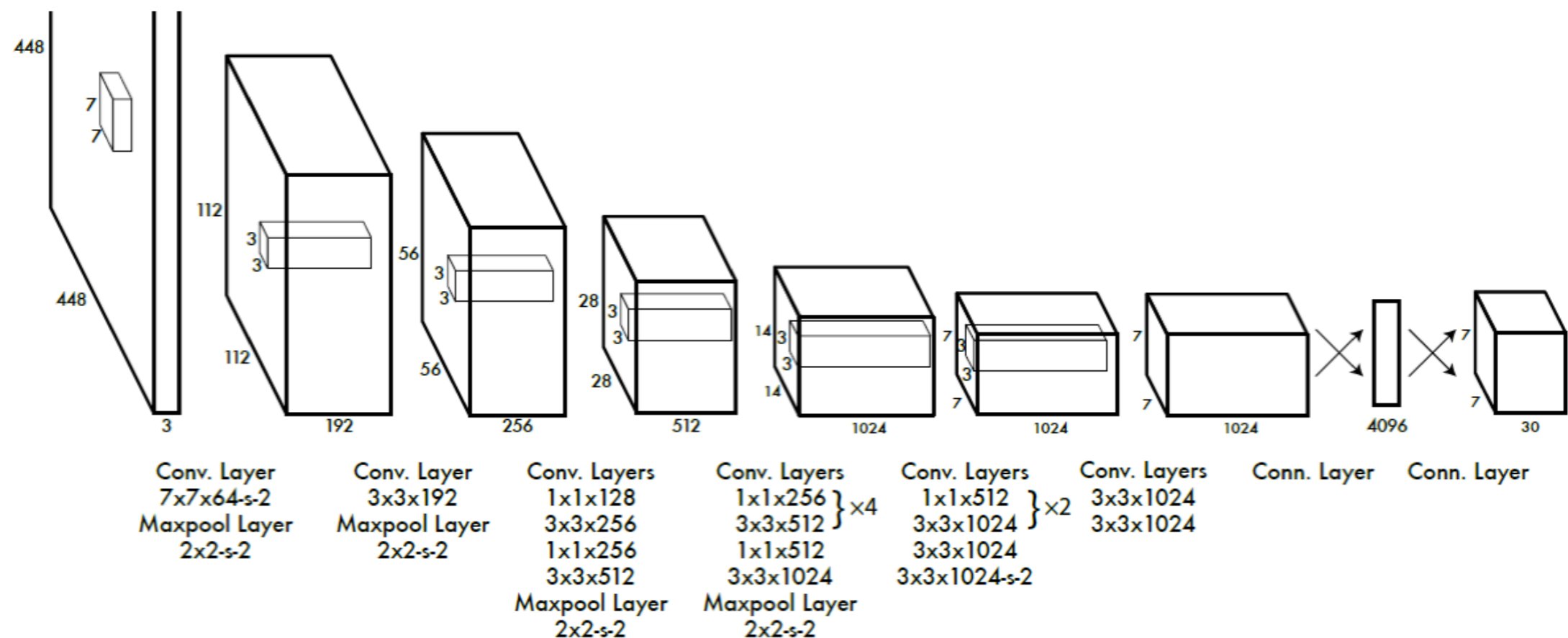
YOLO

- Divide the image into a coarse grid and directly predict class label and a few candidate boxes for each grid cell



YOLO

1. Take conv feature maps at 7x7 resolution
2. Add two FC layers to predict, at each location, a score for each class and 2 bboxes w/ confidences
 - For PASCAL, output is 7x7x30 (30 = 20 + 2*(4+1))



J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, [You Only Look Once: Unified, Real-Time Object Detection](#), CVPR 2016

YOLO

- Objective function:

$$\begin{aligned} & \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[(x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2 \right] \\ & + \lambda_{\text{coord}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} \left[\left(\sqrt{w_i} - \sqrt{\hat{w}_i} \right)^2 + \left(\sqrt{h_i} - \sqrt{\hat{h}_i} \right)^2 \right] \\ & + \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{obj}} (C_i - \hat{C}_i)^2 \\ & + \lambda_{\text{noobj}} \sum_{i=0}^{S^2} \sum_{j=0}^B \mathbb{1}_{ij}^{\text{noobj}} (C_i - \hat{C}_i)^2 \\ & + \sum_{i=0}^{S^2} \mathbb{1}_i^{\text{obj}} \sum_{c \in \text{classes}} (p_i(c) - \hat{p}_i(c))^2 \end{aligned}$$

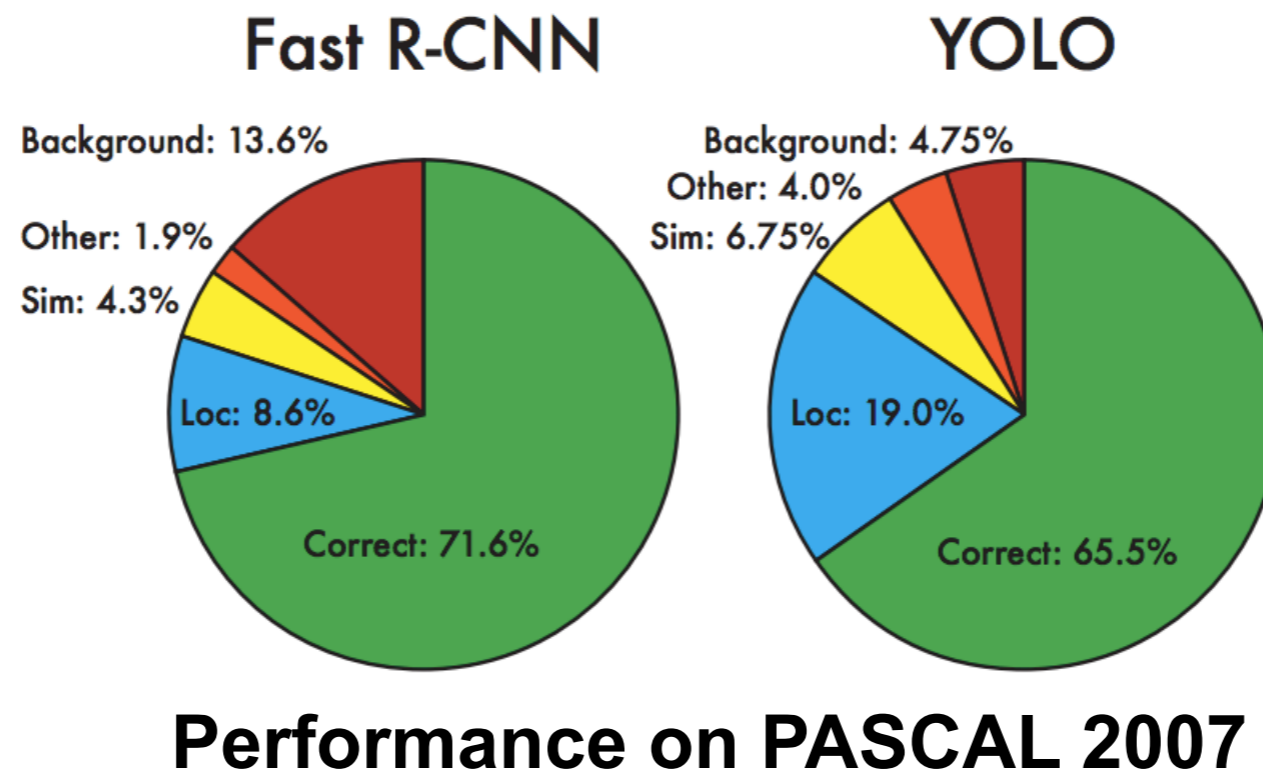
Regression

Object/no object confidence

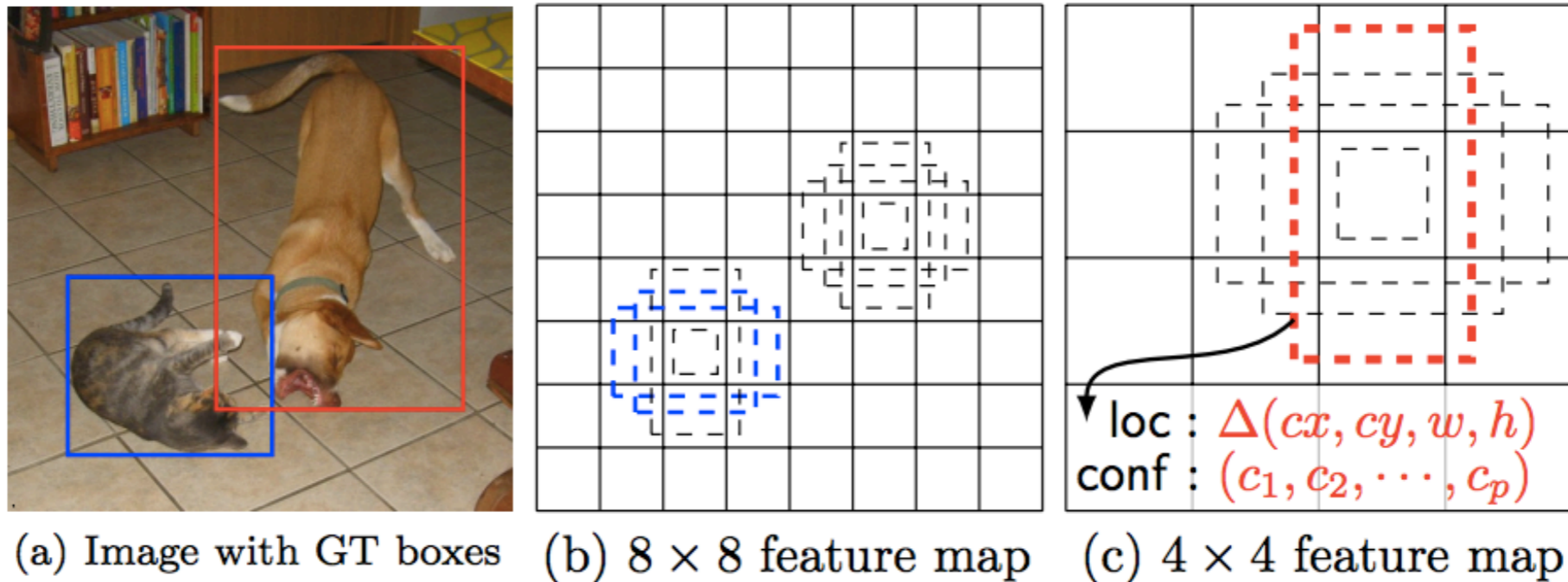
Class prediction

YOLO: Results

- Each grid cell predicts only two boxes and can only have one class – this limits the number of nearby objects that can be predicted
- Localization accuracy suffers compared to Fast(er) R-CNN due to coarser features, errors on small boxes
- 7x speedup over Faster R-CNN (45-155 FPS vs. 7-18 FPS)

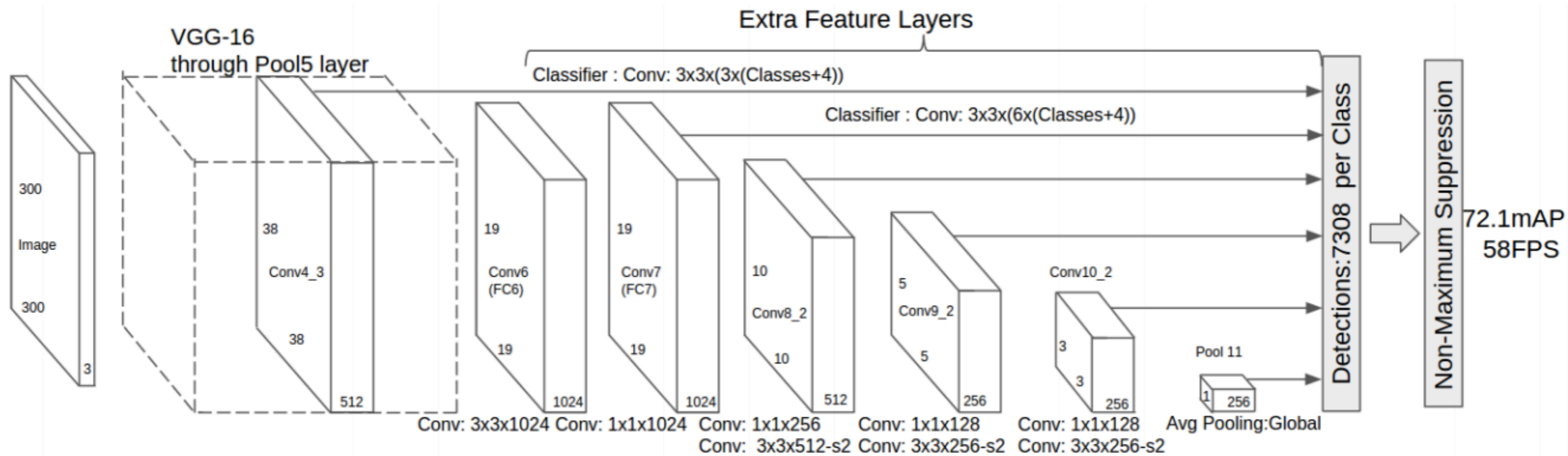
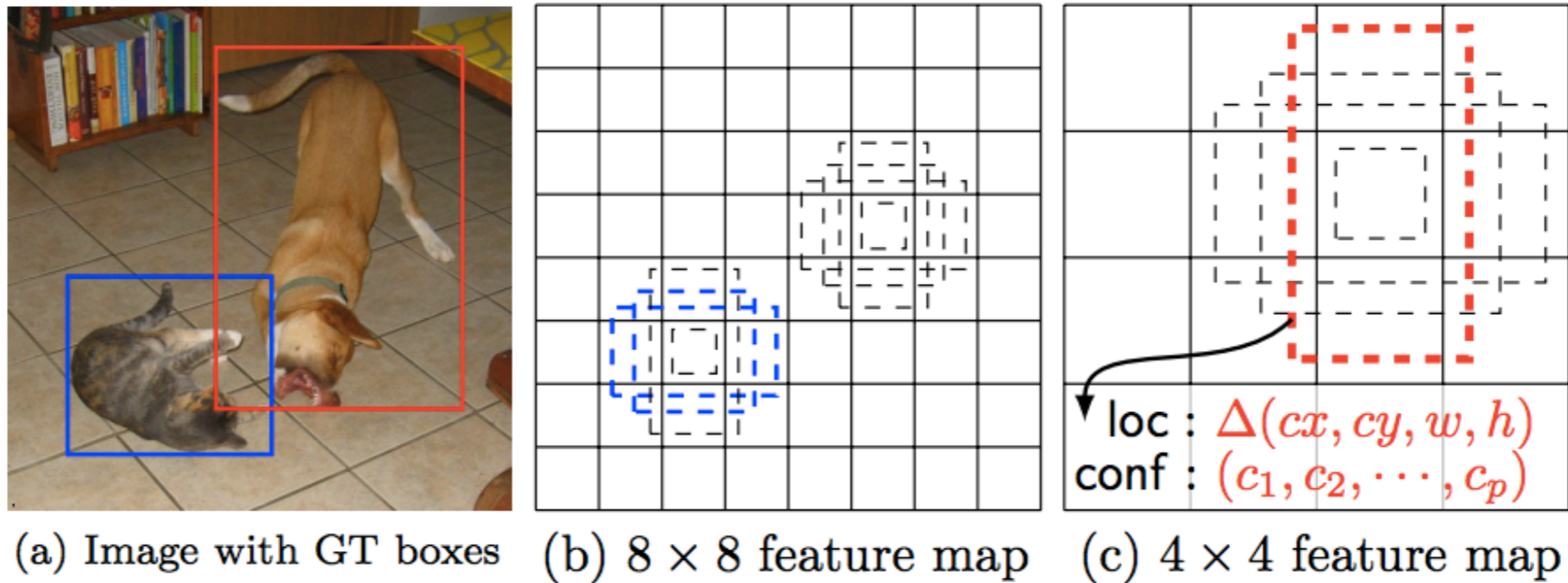


SSD



- Similarly to YOLO, predict bounding boxes directly from conv maps
- Unlike YOLO, do not use FC layers and predict different size boxes from conv maps at different resolutions
- Similarly to RPN, use anchors

SSD



W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. Berg, [SSD: Single Shot MultiBox Detector](#), ECCV 2016.

Object detection overview

- Different choices for “base network”:
 - VGG, ResNet...
- Object detection architecture:
 - Region-based: R-CNN, Fast R-CNN, Faster R-CNN
 - Single shot: YOLO, SSD
- Faster R-CNN is slower but more accurate.
- YOLO/SSD are faster but less accurate.