

Fully Convolutional Networks for Semantic Segmentation

Marek Leibl

Center for Machine Perception, Department of Cybernetics, FEL CTU

March 30, 2017

Presentation Outline

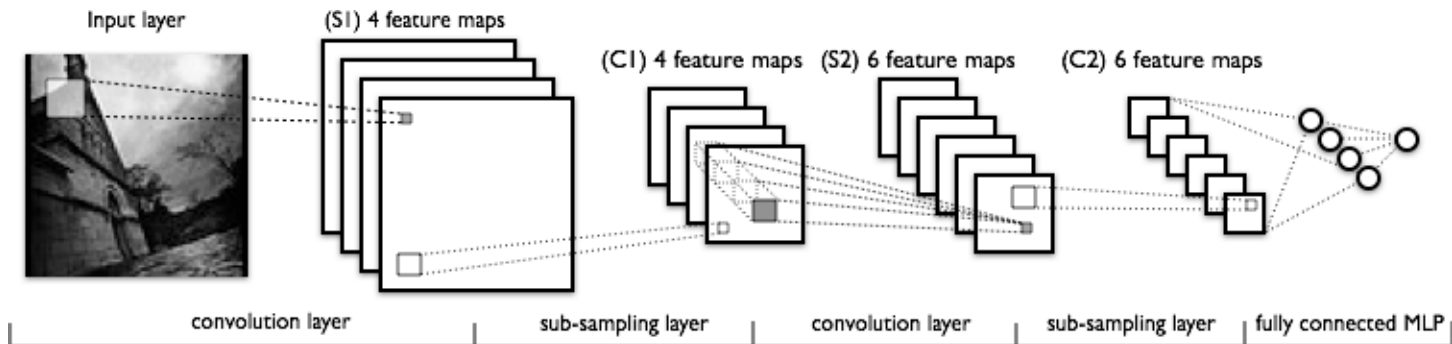
- Introduction ... convolutional neural networks
- Semantic Segmentation ... fully convolutional neural networks
- Conditional Random Fields ... and recurrent neural networks
- Applications ... drosophila eggs segmentation



Convolutional Neural Networks

Convolutional Neural Networks - Introduction

- Convolutional Neural Networks (CNNs)
 - layered artificial neural network
 - inspired by organization of the animal visual cortex
 - local connectivity
- Deep Convolutional Neural Networks
 - many convolutional layers often combined with max-pooling and relu



CNNs - Convolution

- Convolution = building block of CNNs



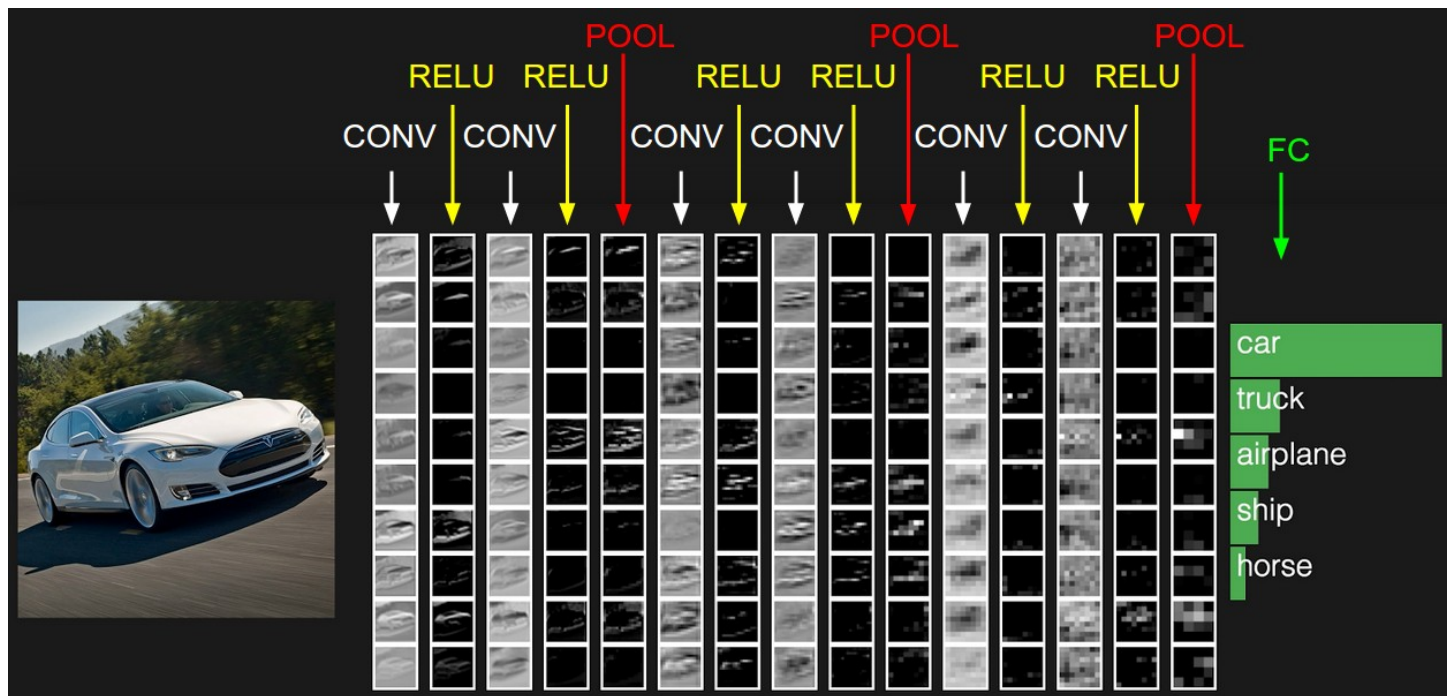
*

1	0	-1
2	0	-2
1	0	-1



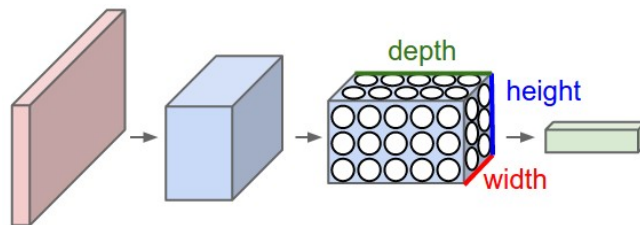
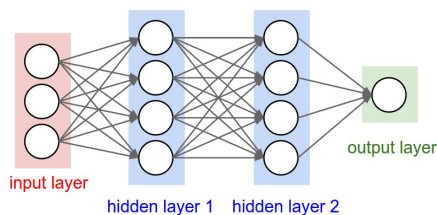
$$G[i, j] = \sum_u \sum_v H[u, v] \cdot F[i - u, j - v]$$

CNNs - Example



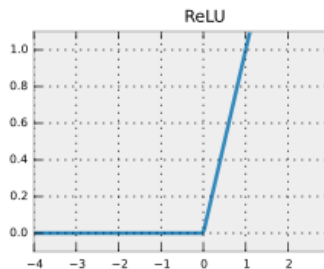
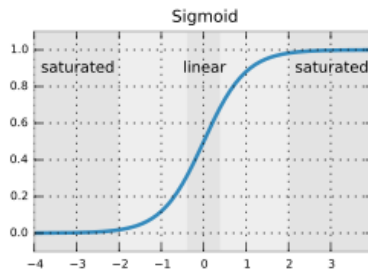
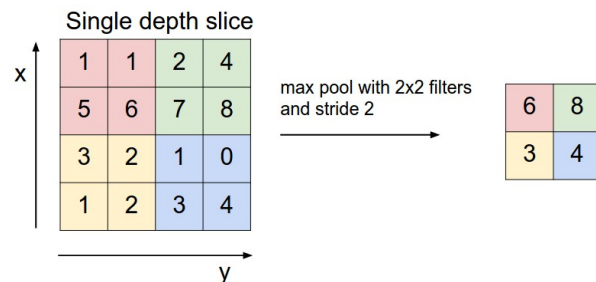
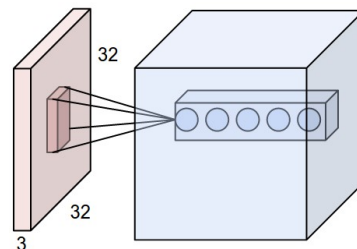
CNNs - Structure

- Neurons in each layer organized into 3D matrixes
 - 2 spatial dimensions
 - 1 dimension for channels (features)
- Local connectivity
 - each neuron relatively small number of input connections ... receptive field
- Shared weights
 - weights are invariant to translations in spatial dimensions
 - reduced search space and better robustness



CNNs - Building blocks

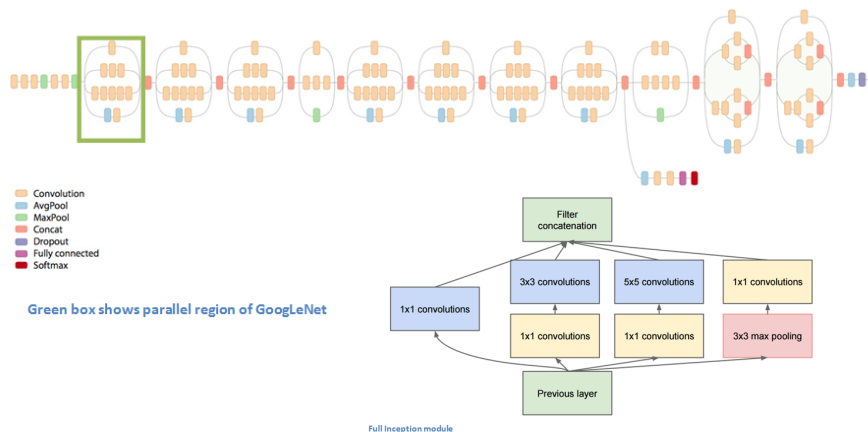
- Convolutional layer
 - trainable weights
 - typically followed by Relu
- Relu (rectified linear unit)
 - $f(x) = \max(0, x)$
 - introduces non-linearity to CNN
 - easy to compute, no saturation
- Max-pooling
 - reducing spatial size
- Softmax
 - normalized exponential function
 - $$f(x)_j = \frac{e^{x_j}}{\sum_i e^{x_i}}$$



CNNs - Deep CNN architectures for image classification

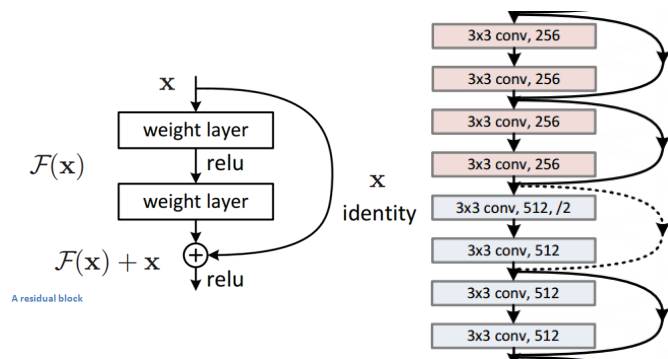
- GoogLeNet (2015)

- 9 Inception modules
- over 100 layers
- trained on a few GPUs within a week



- Microsoft ResNet (2015)

- 152 layers
- residual blocks
→ easier to optimize
- trained on an 8 GPU machine for two to three weeks





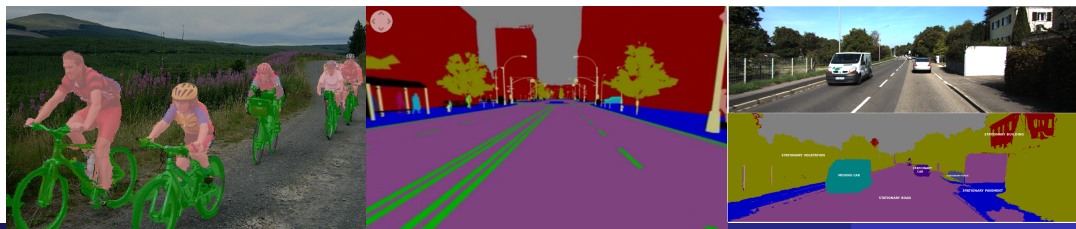
Semantic Segmentation

Semantic Segmentation vs Image Classification

- Image classification: determine class for a given image
 - super human performance on some tasks

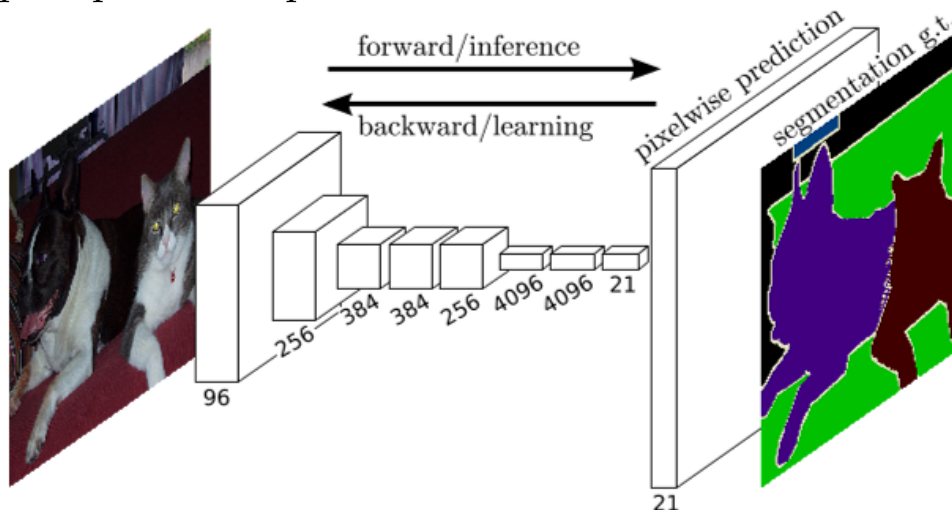


- Semantic segmentation: determine class for each pixel
 - harder problem



Fully Convolutional Neural Networks

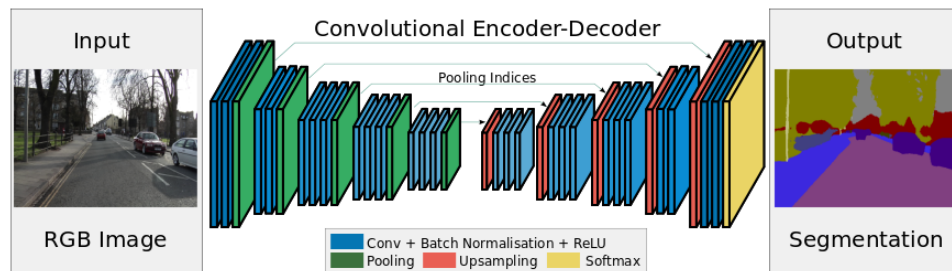
- Fully Convolutional Neural Network
 - no fully connected layer
 - output: pixel-wise prediction



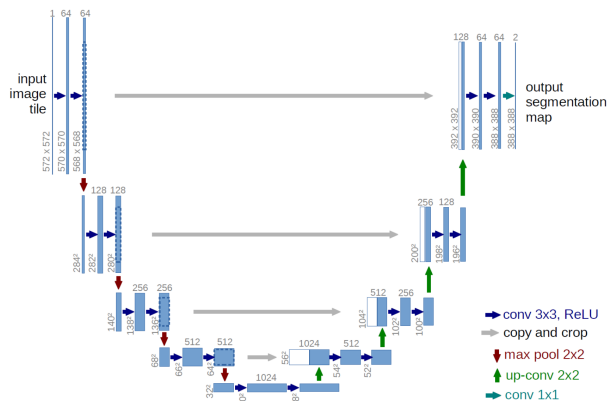
- Training Fully Convolutional NNs
 - needed label for each pixel (ground truth)
 - need to be manually segmented by human
→ training set is usually expensive

Fully Convolutional Neural Networks: Architectures

- SegNet (2015)
 - deep encoder-decoder architecture
 - University of Cambridge

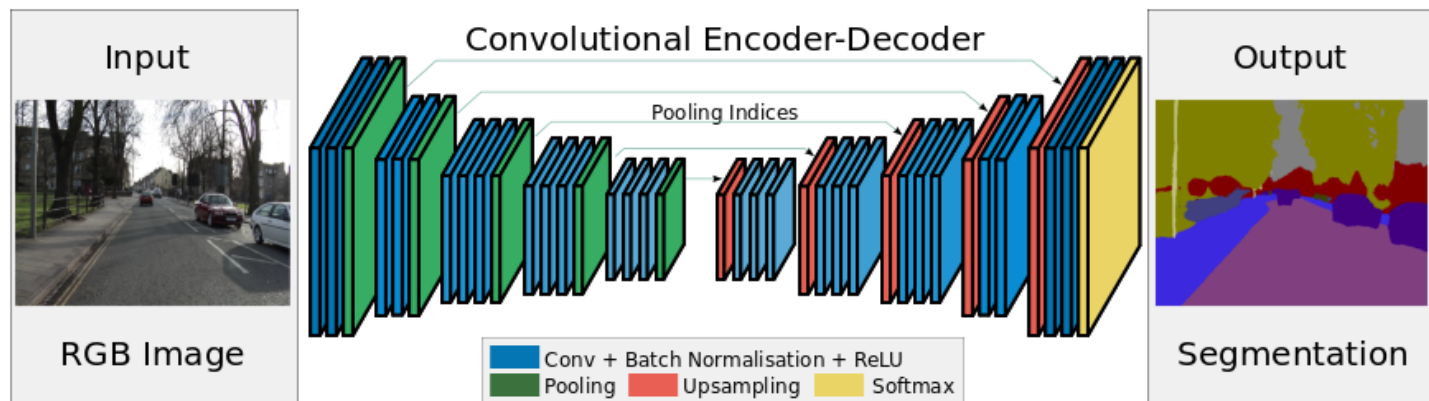


- U-Net (2015)
 - U-shaped architecture
 - biomedical image segmentation



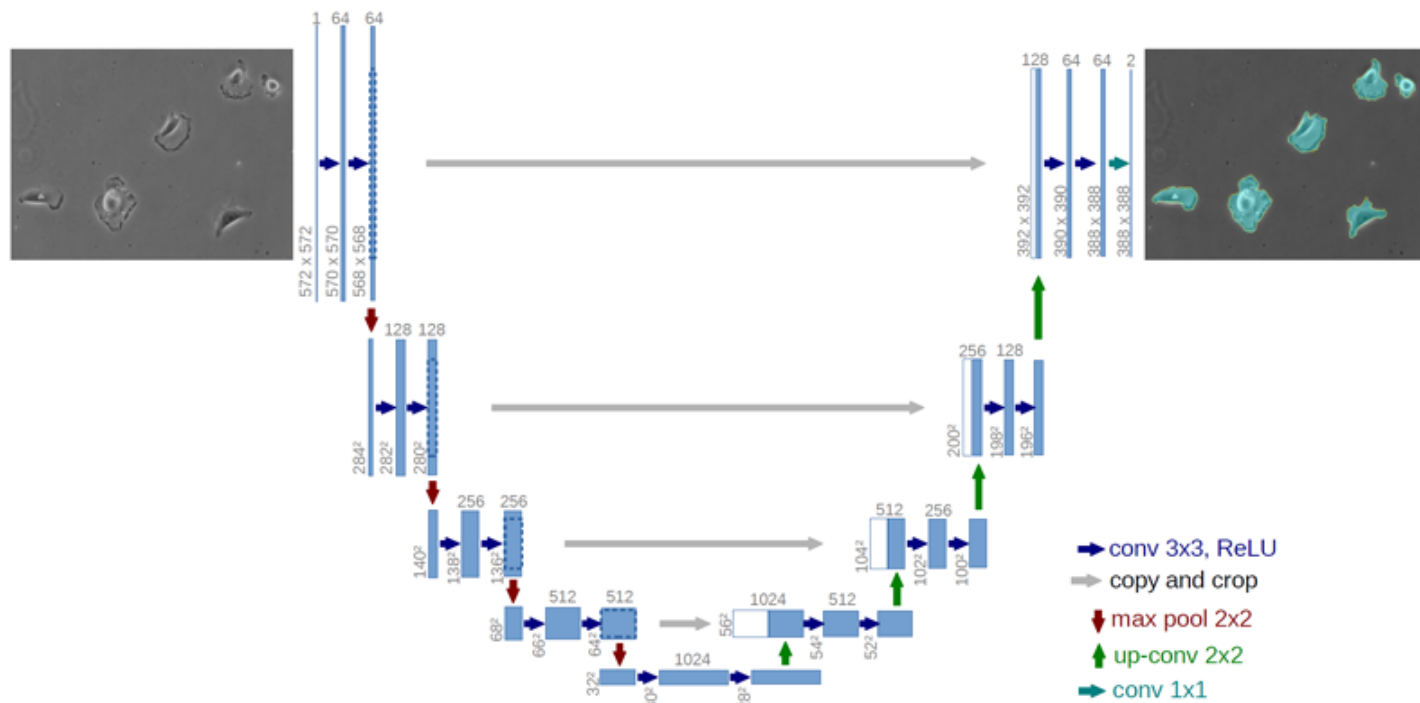
Fully Convolutional Neural Networks: SegNet

- SegNet (2015)
 - deep encoder-decoder architecture
 - University of Cambridge
 - <http://mi.eng.cam.ac.uk/projects/segnet/>



Fully Convolutional Neural Networks: U-Net

- U-Net (2015)
 - U-shaped architecture
 - biomedical image segmentation

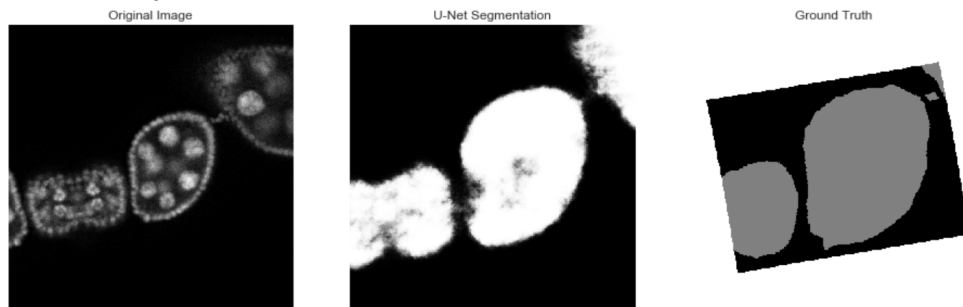




Conditional Random Fields & Convolutional Neural Networks

Drawbacks of fully connected CNNs

- CNNs outperforms “classical” methods on majority of semantic segmentation tasks
- they still fail in some cases

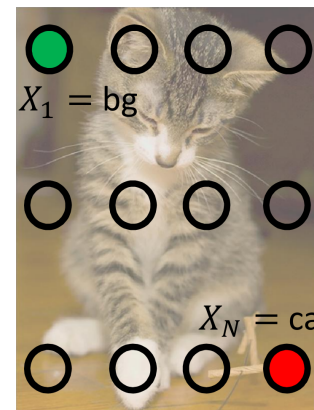
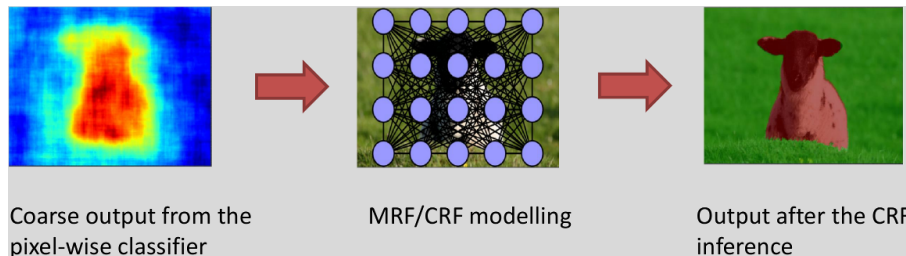
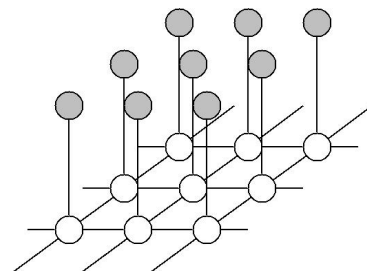


- Cross entropy loss $E = - \sum_x \sum_c p_{xc} \log(y_{xc})$
 - does not consider global consistency
 - can only use information within its receptive field
 - → “holes” in the segmentation map
 - → inconsistent segmentation



Conditional Random Fields: Introduction

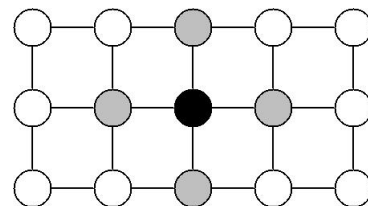
- Conditional Random Field (CRF)
 - probabilistic graphical model
 - two types of variables:
 - observations (e.g. original image)
 - random variables (e.g. pixel-wise segmentation)
 - used for structured prediction
e.g. semantic segmentation
 - these days outperformed by CNNs
- CRF as a post-processing method to improve CNN output



Conditional Random Fields: Formal Definition

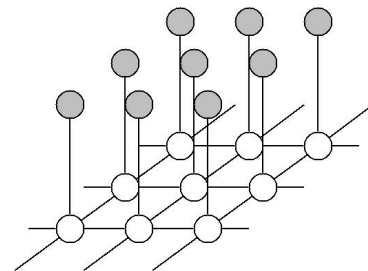
- Markov Random Field (MRF)

- $G = (V, E)$ undirected graph
- $Y = \{Y_v\}_{v \in V}$
- $Y_v | Y_{N(v)} \perp\!\!\!\perp Y_{V \setminus N[v]} | Y_{N(v)}$



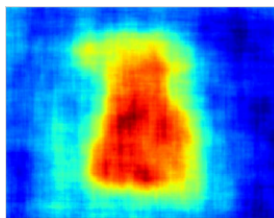
- Conditional Random Field (CRF)

- $G = (V, E)$ undirected graph
 - $V = \{1, \dots, N\}$... pixels
 - E ... define neighbors
- $I = \{I_i\}_1^N$... observation for each pixel
- $X = \{X_i\}_1^N$... random variable for each pixel
- $X|I$ is a Markov Random Field (MRF)
 - $\rightarrow P(X = x|I) = \frac{1}{Z(I)} \exp \{-E(x|I)\}$

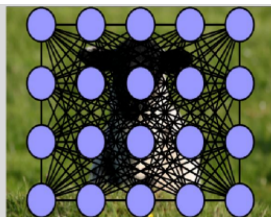


Conditional Random Fields Postprocessing

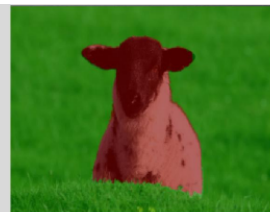
- Goal: maximize $P(X = x|I) = \frac{1}{Z(I)} \exp \{-E_I(x)\}$
 - I ... CNN segmentation
 - X ... final segmentation
- \rightarrow minimize energy function:
 - $E_I(x) = \sum_i \psi_u(x_i) + \sum_{i < j} \psi_p(x_i, x_j)$
 - NP-hard for general graphs
- \rightarrow Mean Field Approximation



Coarse output from the
pixel-wise classifier



MRF/CRF modelling



Output after the CRF
inference

CRF: Mean Field Iteration as RNN

$$Q_i(x_i = l) = \frac{1}{Z_i} \exp \left\{ -\psi_u(x_i) - \sum_{l' \in \mathcal{L}} \mu(l, l') \sum_{m=1}^K w^{(m)} \sum_{j \neq i} k^{(m)}(\mathbf{f}_i, \mathbf{f}_j) Q_j(l') \right\}.$$

Algorithm 1 Mean field in fully connected CRFs

Initialize Q

while not converged **do**

$\tilde{Q}_i^{(m)}(l) \leftarrow \sum_{j \neq i} k^{(m)}(\mathbf{f}_i, \mathbf{f}_j) Q_j(l)$ for all m

$\hat{Q}_i(x_i) \leftarrow \sum_{l \in \mathcal{L}} \mu^{(m)}(x_i, l) \sum_m w^{(m)} \tilde{Q}_i^{(m)}(l)$

$Q_i(x_i) \leftarrow \exp\{-\psi_u(x_i) - \hat{Q}_i(x_i)\}$

normalize $Q_i(x_i)$

end while

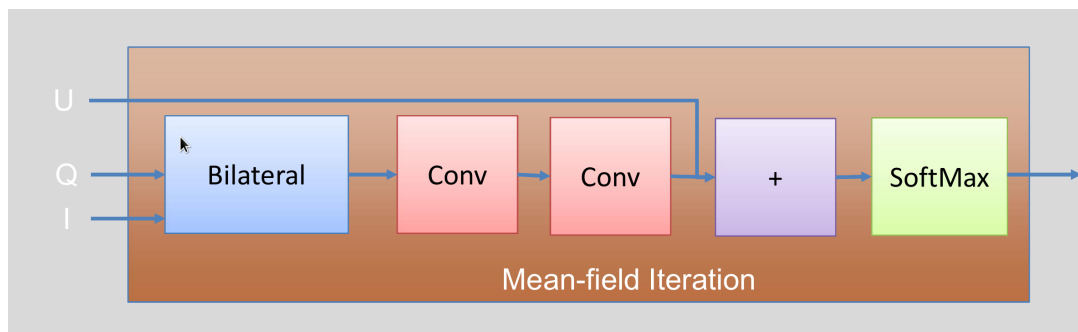
$\triangleright Q_i(x_i) \leftarrow \frac{1}{Z_i} \exp\{-\phi_u(x_i)\}$

\triangleright See Section 6 for convergence analysis

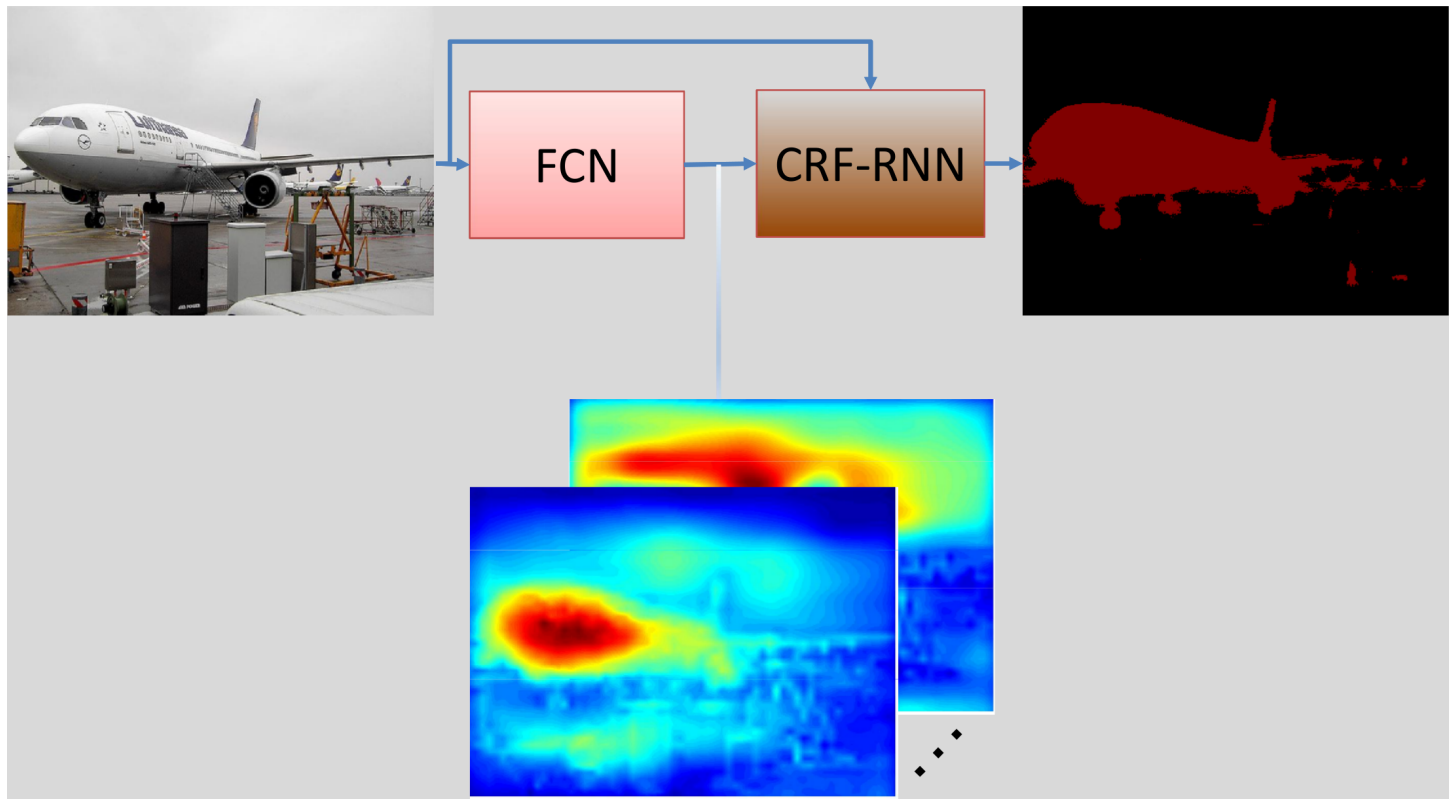
\triangleright **Message passing** from all X_j to all X_i

\triangleright **Compatibility transform**

\triangleright **Local update**



CRF as RNN

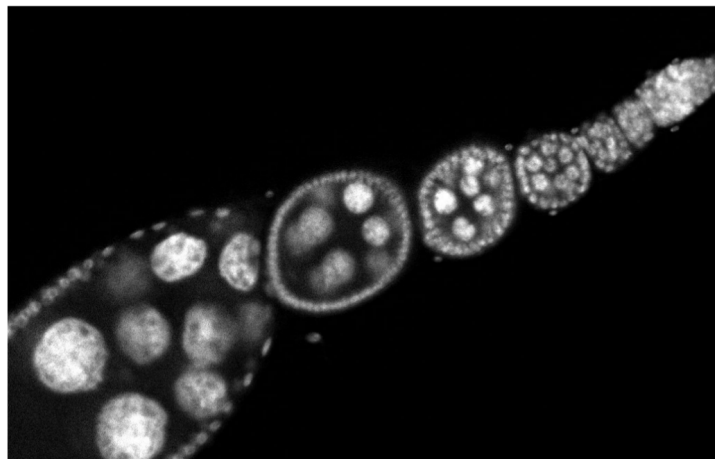


A grayscale microscopic image of a row of Drosophila eggs. The eggs are oval-shaped and arranged in a diagonal line from the bottom left towards the top right. Each egg shows internal cellular structures, including a prominent nucleus and various organelles. The background is a uniform gray.

Drosophila Eggs Segmentation

Problem Definition

- Binary segmentation of microscopy scans
 - classes: the eggs and the background



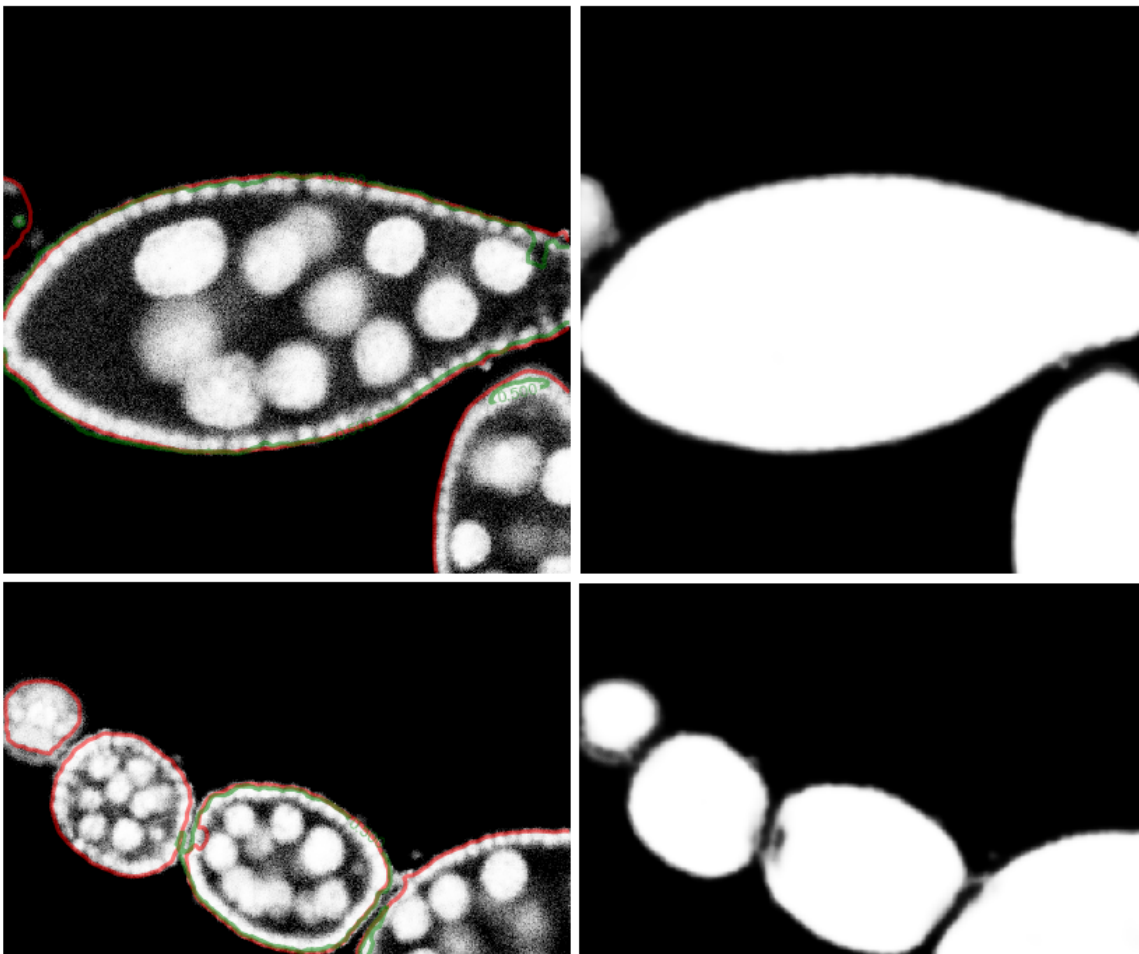
(a) Input slice



(b) Desired output segmentation

- 75 manually annotated images of drosophila eggs
- 4000 unlabeled images

Results



A photograph of three cyclists riding through a field of purple flowers. The cyclist on the left is wearing a red and black striped jersey. The cyclist in the center is wearing a blue and yellow jersey and a yellow helmet. The cyclist on the right is wearing a green and white jersey and a red helmet. They are all riding towards the camera. The background shows a line of trees and a hill under a cloudy sky.

Thank You!

- <http://cs231n.stanford.edu/>
- <https://adeshpande3.github.io/adeshpande3.github.io/The-9-Deep-Learning-Papers-You-Need-To-Know-About.html>
- <http://www.robots.ox.ac.uk/~szheng/CRFasRNN.html>
- <https://www.tensorflow.org>