

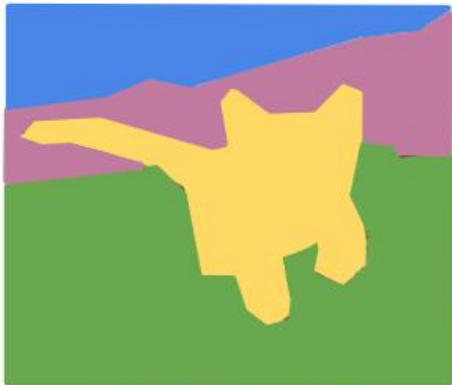
# Lecture 10: Semantic Segmentation

# Contents

1. Semantic Segmentation
2. Segmentation as clustering: k-means, mean-shift
3. Upsampling
4. FCN, U-Net, Tiramisu
5. Mask R-CNN

# Computer Vision Tasks

**Semantic Segmentation**



GRASS, CAT,  
TREE, SKY

No objects, just pixels

**Classification  
+ Localization**



CAT

Single Object

**Object  
Detection**



DOG, DOG, CAT

Multiple Object

**Instance  
Segmentation**



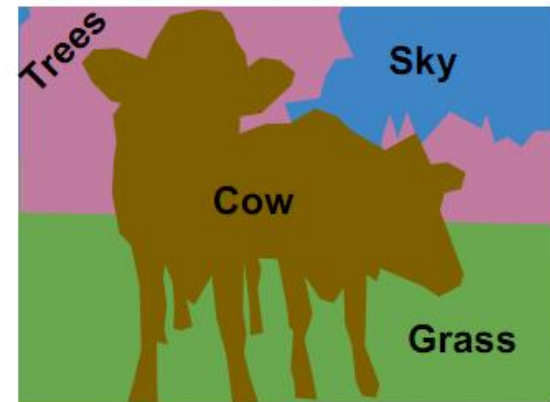
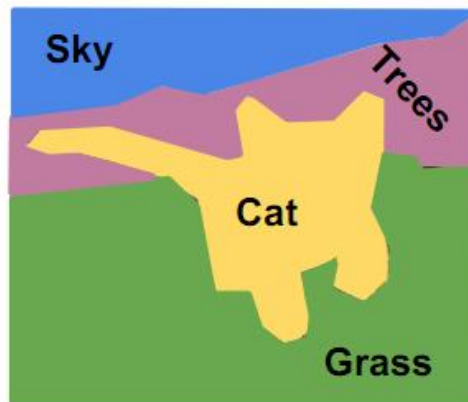
DOG, DOG, CAT

[This image is CC0 public domain](#)

# Semantic Segmentation

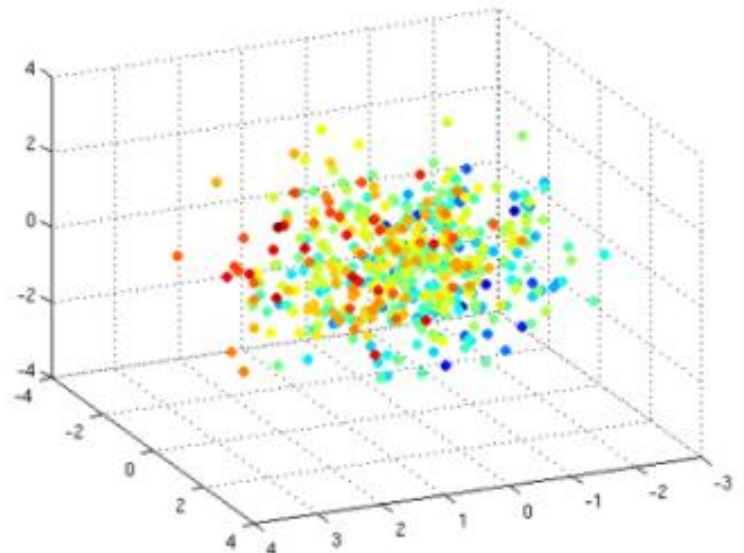
Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



# Segmentation as Clustering

- Pixels are points in a high-dimensional space
  - color: 3d
  - color + location: 5d
- Cluster pixels into segments



# Clustering: K-Means

## Algorithm:

1. Randomly initialize the cluster centers,  $c_1, \dots, c_K$
  2. Given cluster centers, determine points in each cluster
    - For each point  $p$ , find the closest  $c_i$ . Put  $p$  into cluster  $i$
  3. Given points in each cluster, solve for  $c_i$ 
    - Set  $c_i$  to be the mean of points in cluster  $i$
  4. If  $c_i$  have changed, repeat Step 2
- Properties
    - Will always converge to some solution
    - Can be a “local minimum”
      - Does not always find the global minimum of objective function:

$$\sum_{i=1}^k \sum_{x \in C_i} \|x - c_i\|^2$$



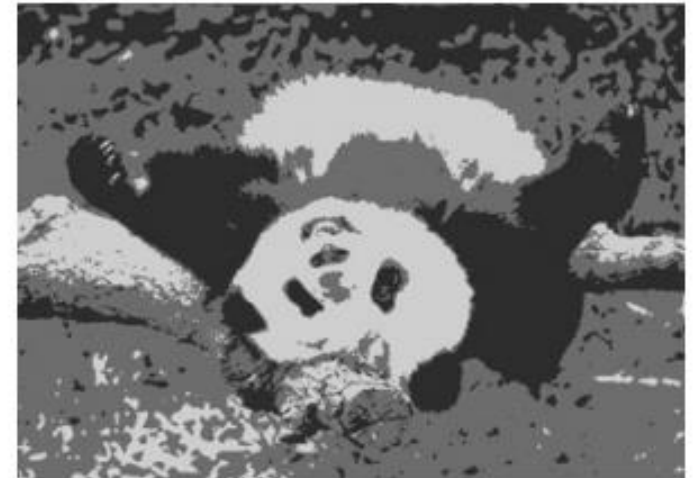
# Clustering: K-Means



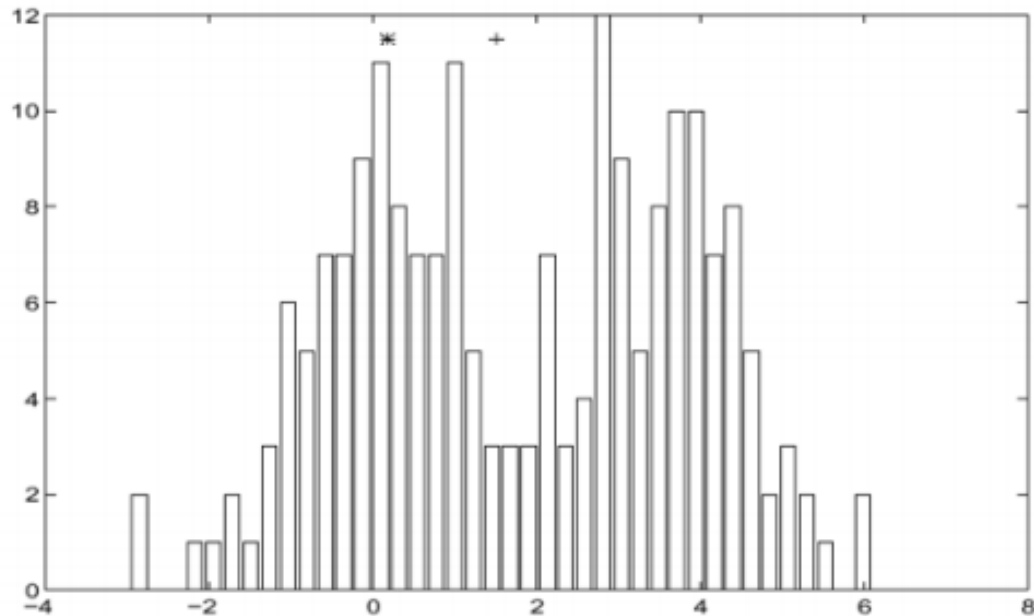
k=2



k=3



# Clustering: Mean-shift



1. Initialize random seed, and window  $W$

2. Calculate center of gravity (the “mean”) of  $W$ :  $\frac{1}{|W|} \sum_{x \in W} x$

- Can generalize to arbitrary windows/kernels

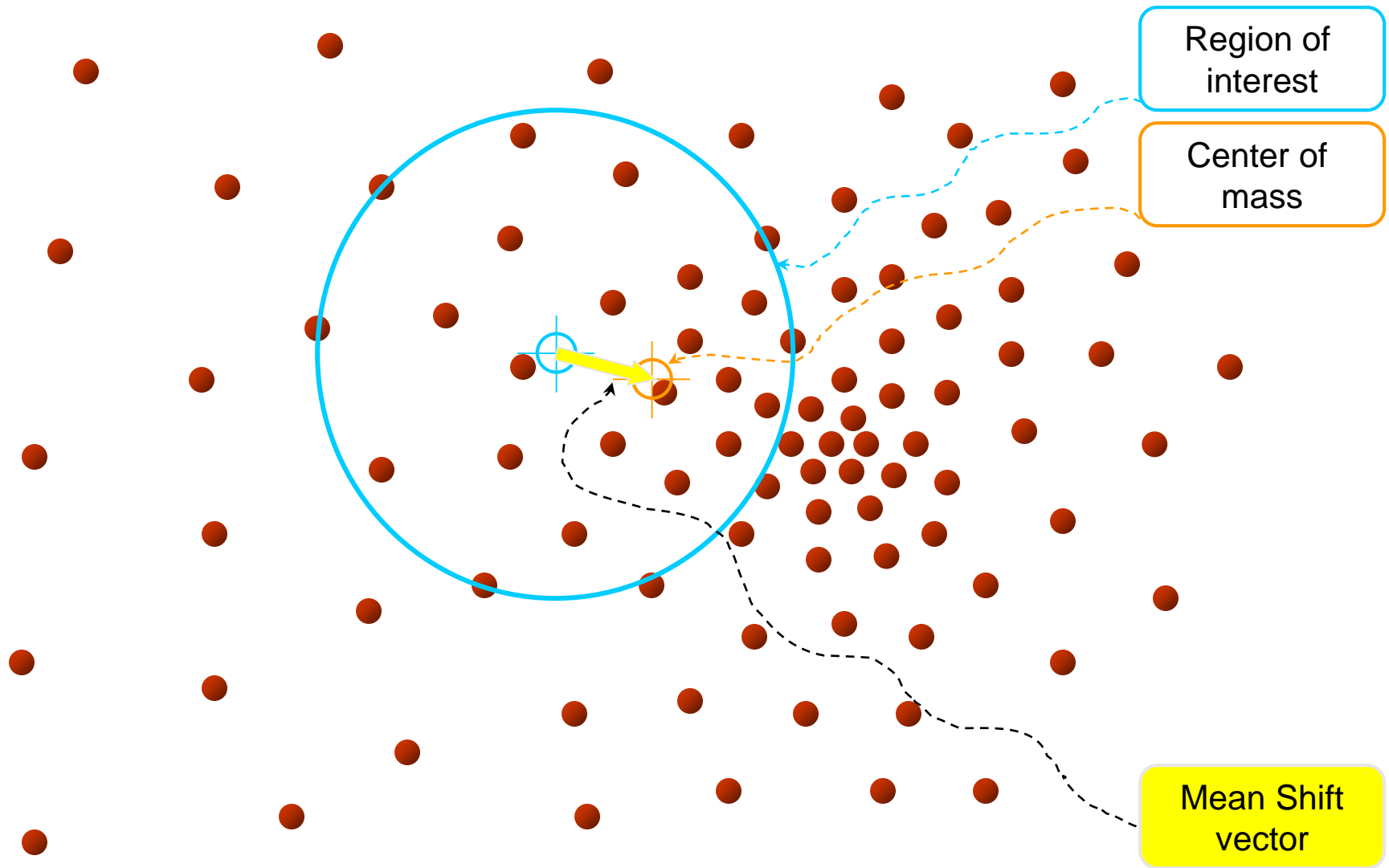
3. Shift the search window to the mean

4. Repeat Step 2 until convergence

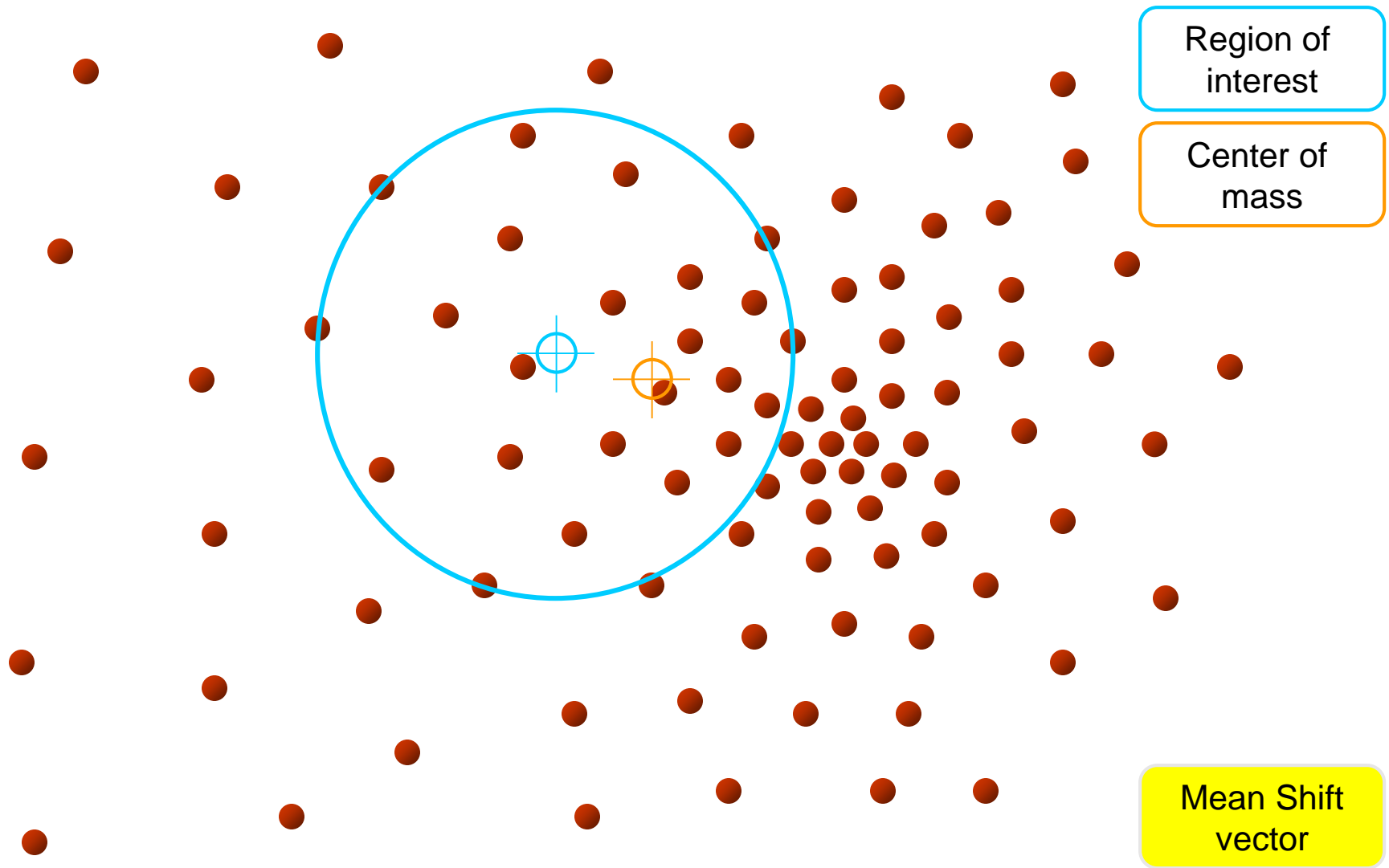
Only parameter: window size



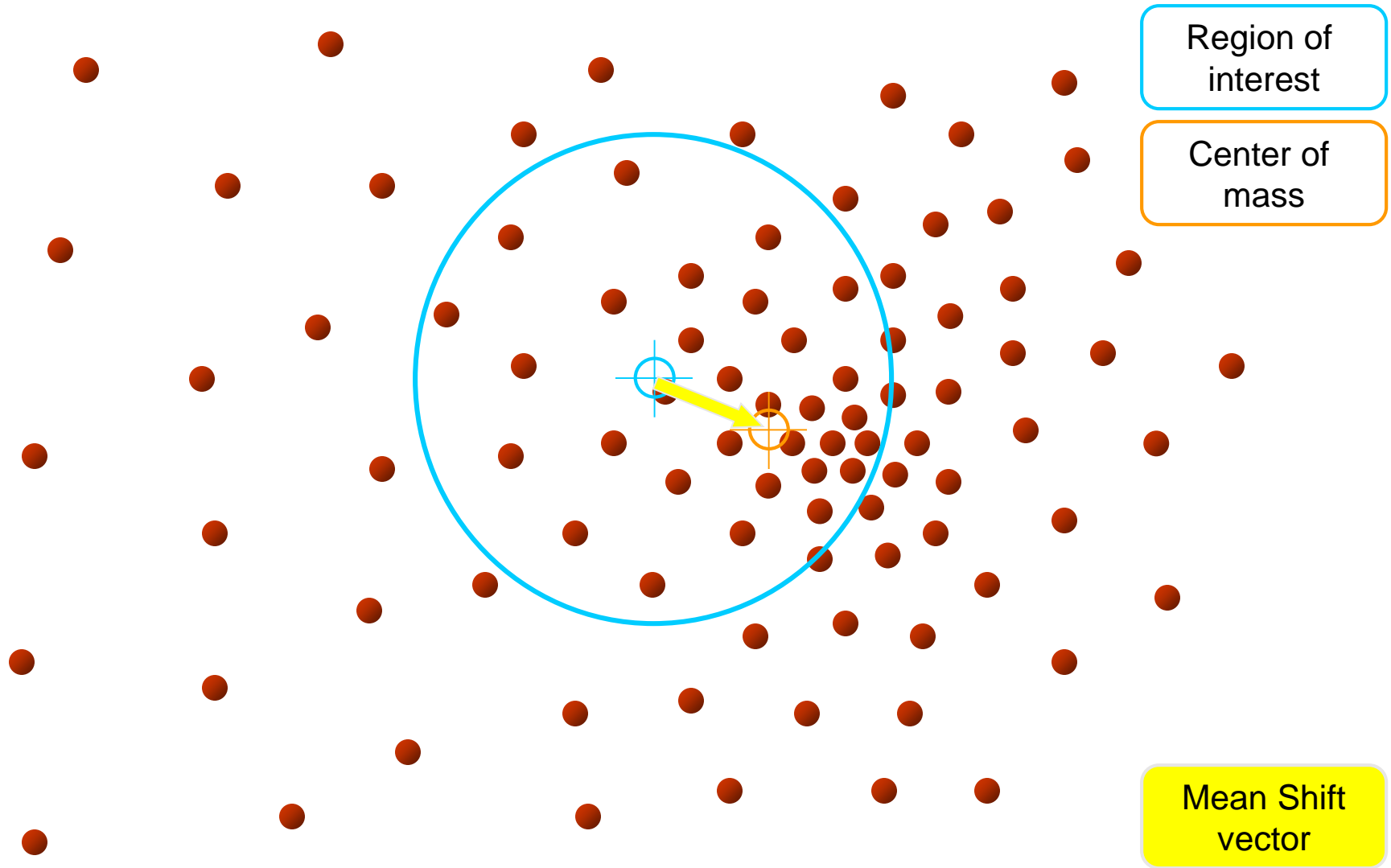
# Mean-Shift



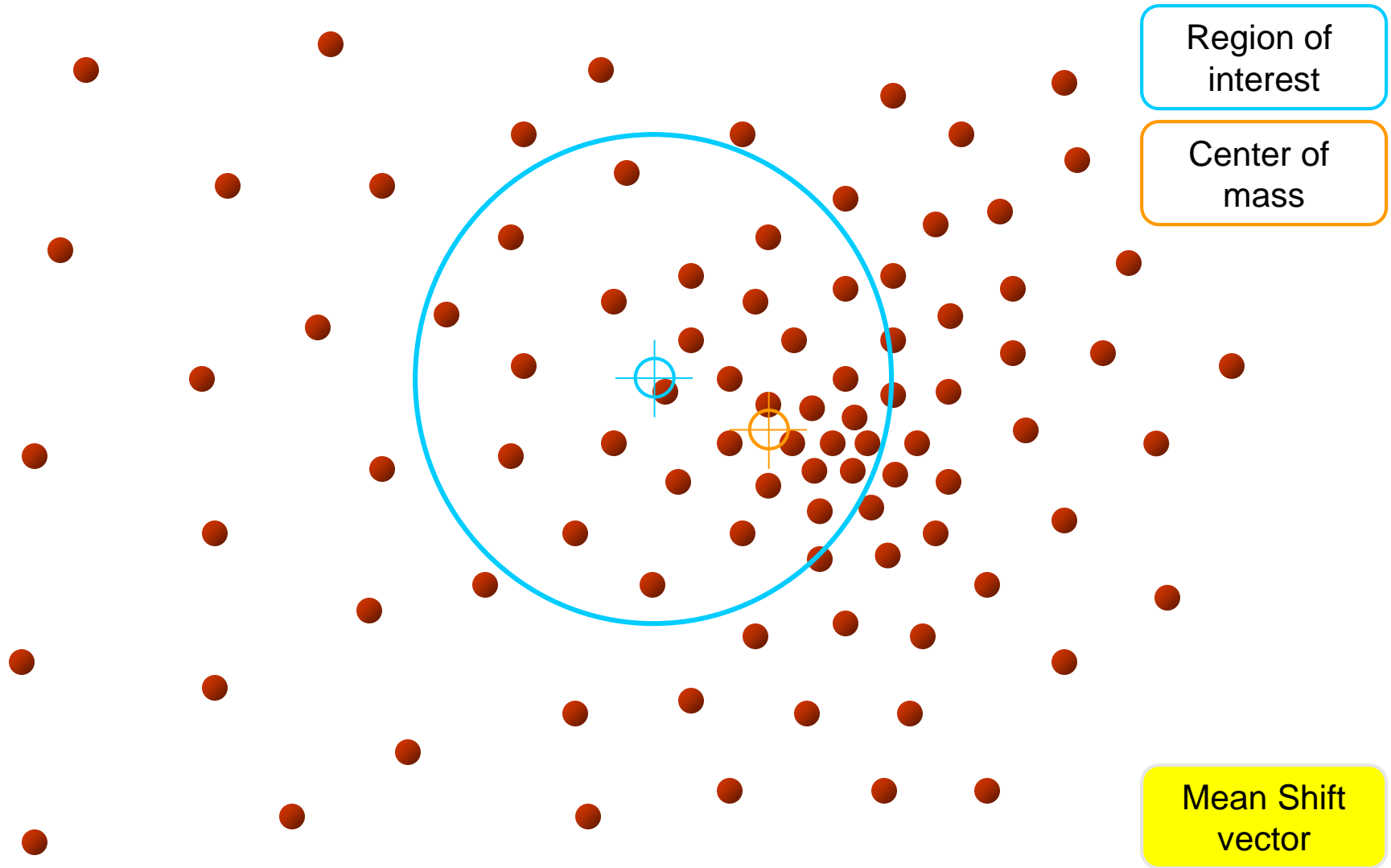
# Mean-Shift



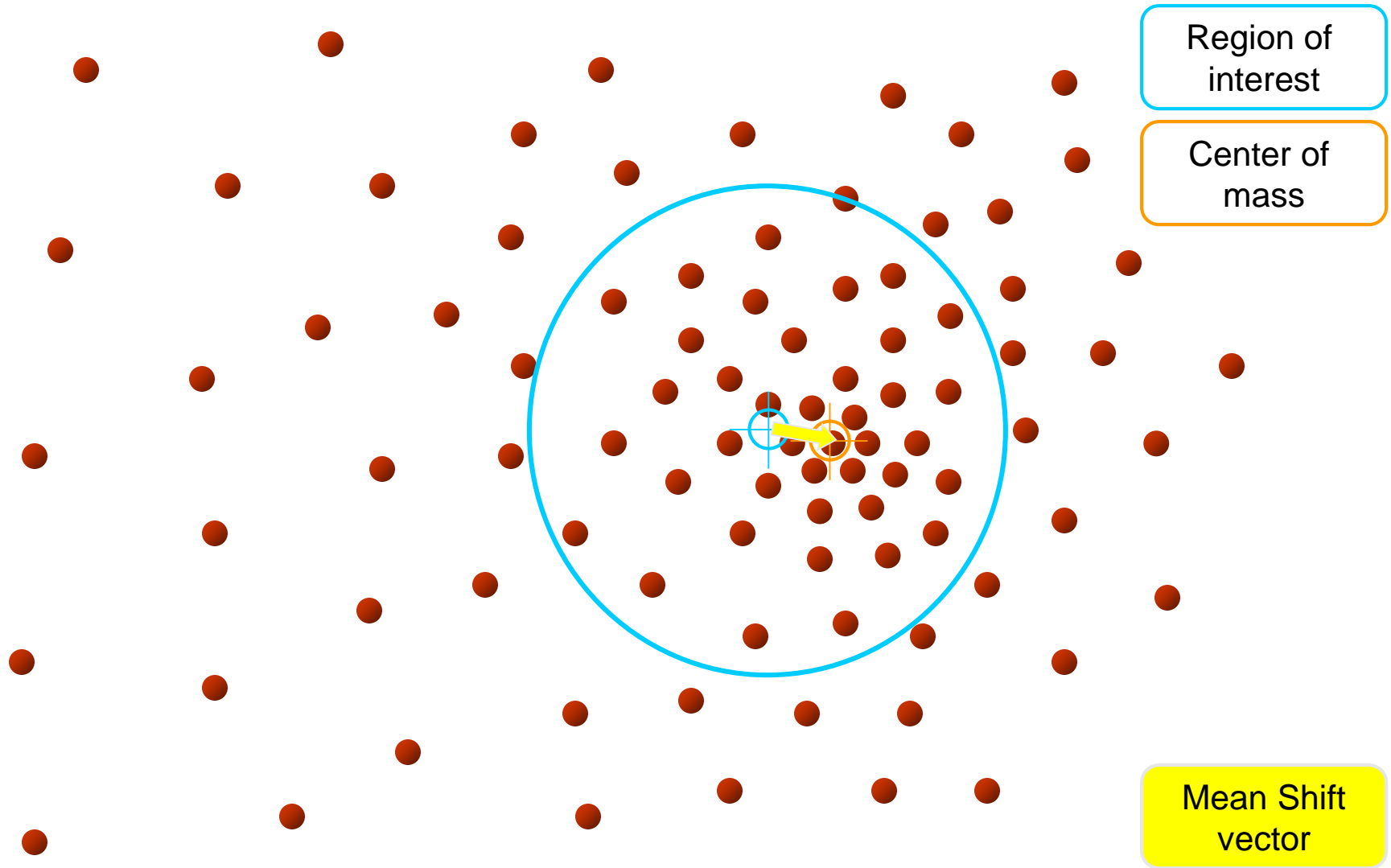
# Mean-Shift



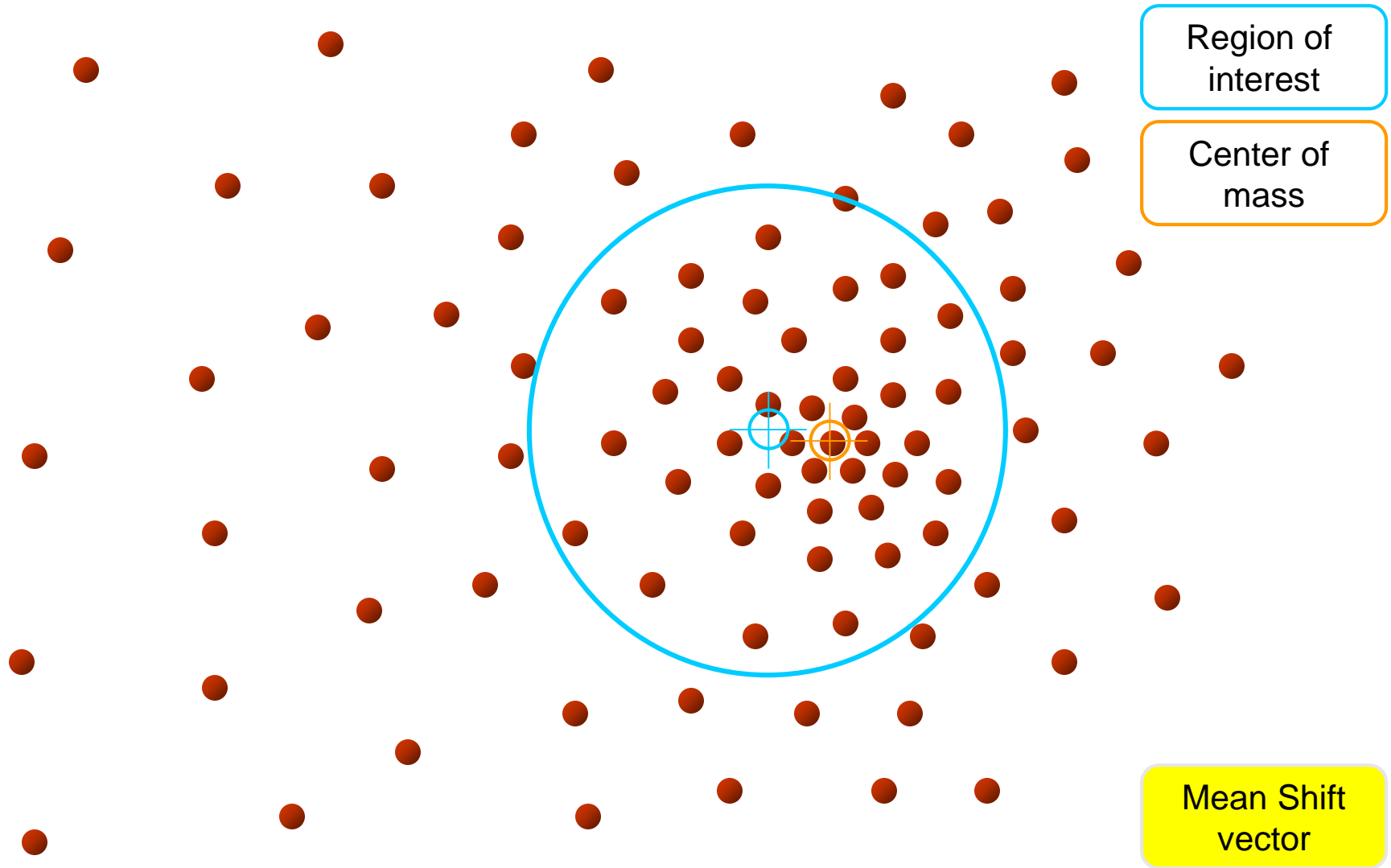
# Mean-Shift



# Mean-Shift

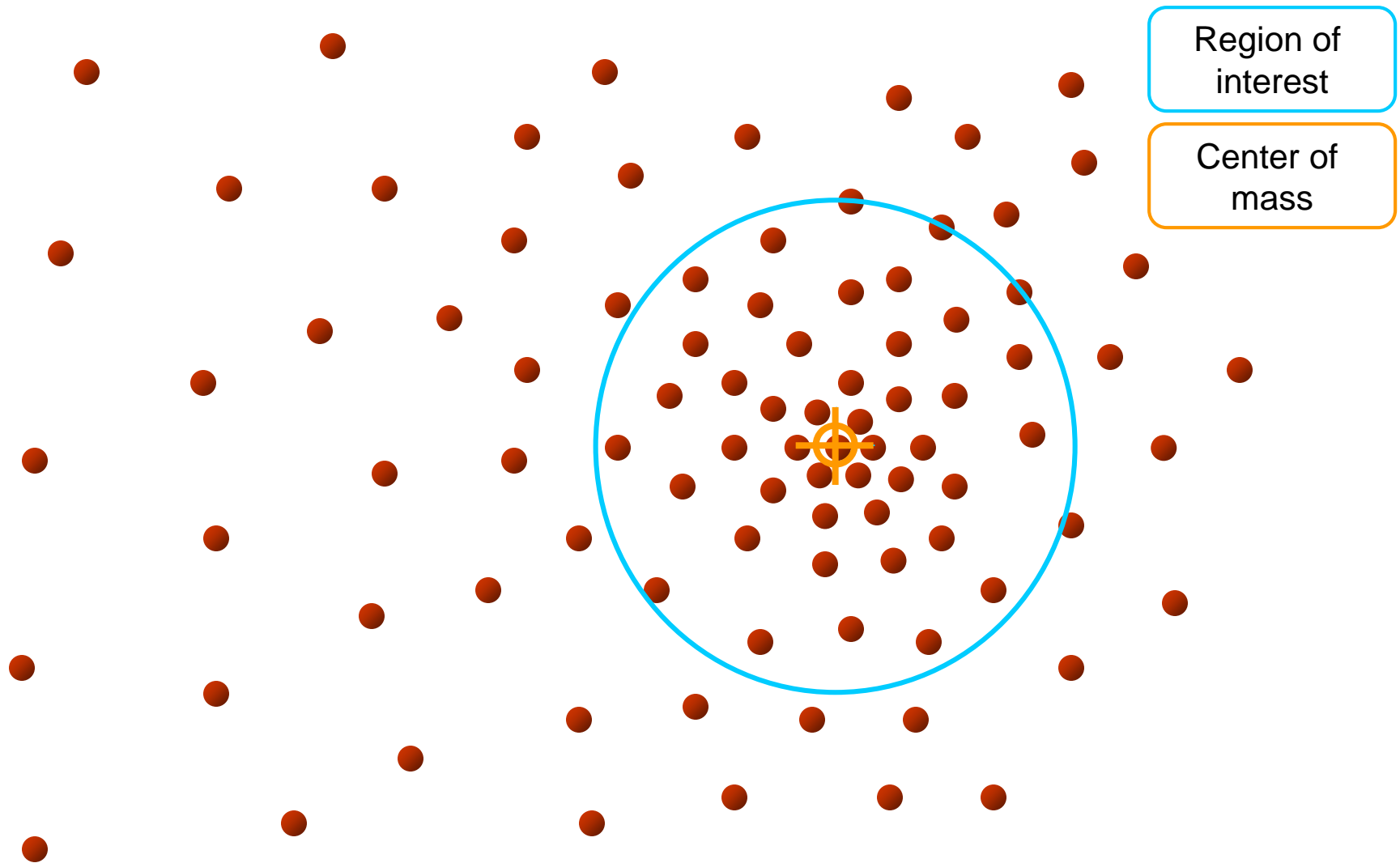


# Mean-Shift



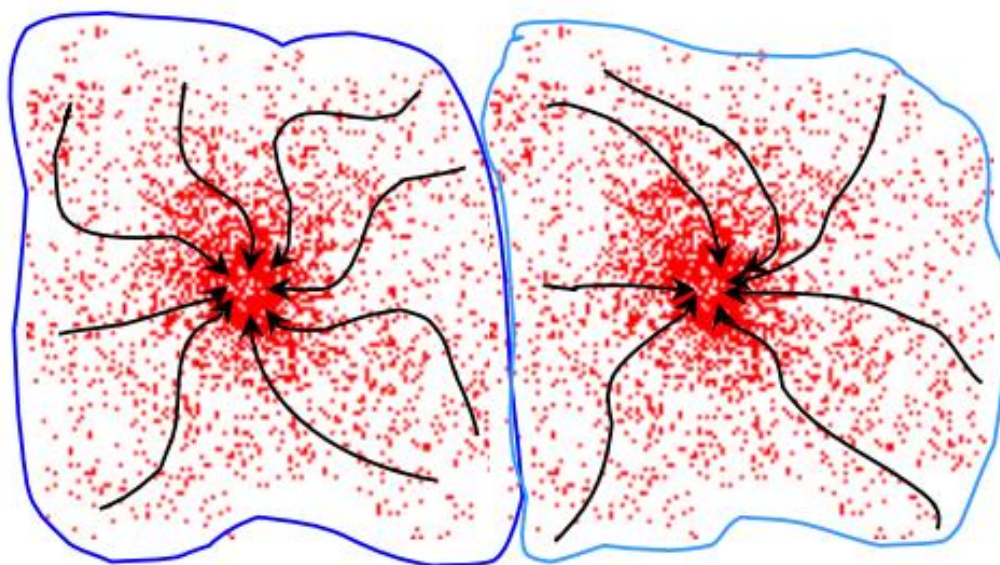


# Mean-Shift



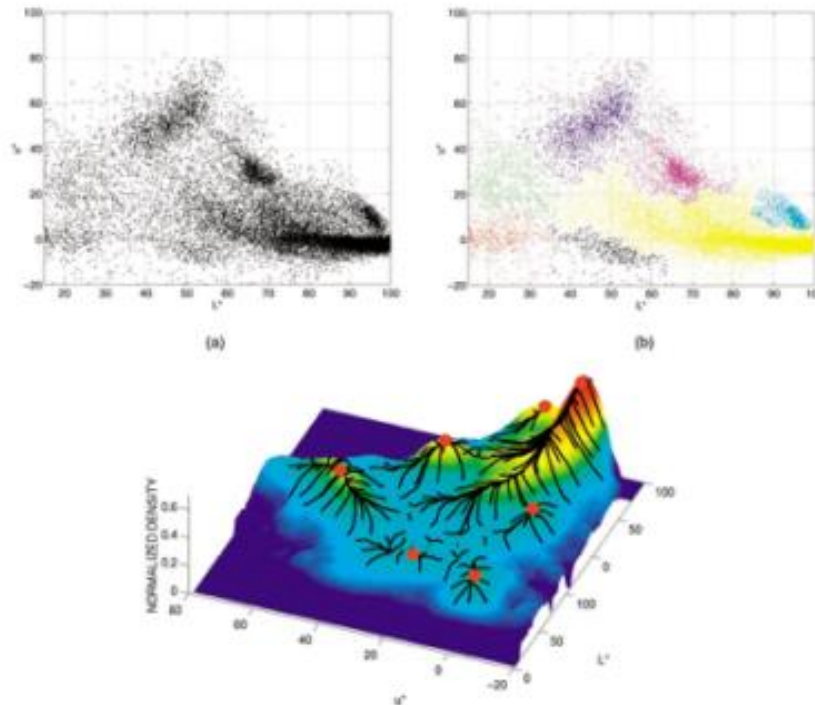
# Mean-Shift

- Cluster: all data points in the attraction basin of a mode
- Attraction basin: the region for which all trajectories lead to the same mode

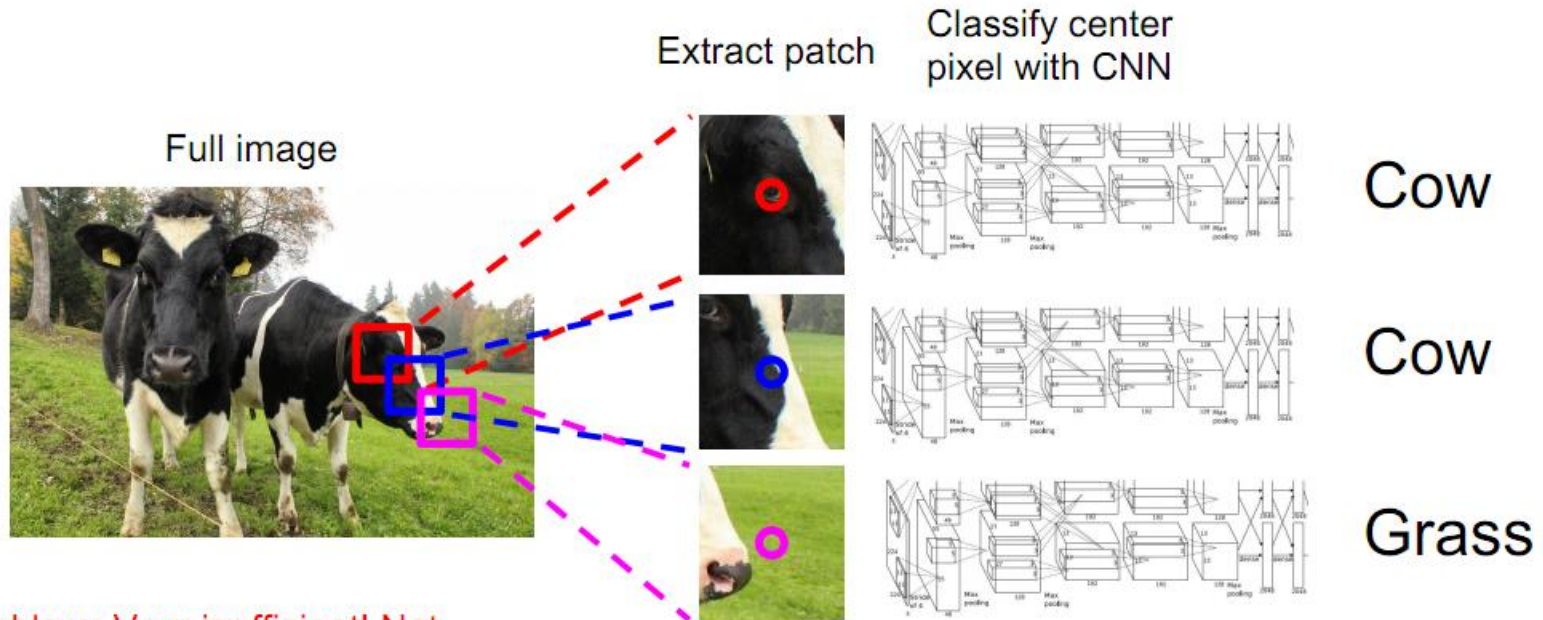


# Mean-Shift for segmentation

- Find features (color, gradients, texture, etc)
- Initialize windows at individual pixel locations
- Perform mean shift for each window until convergence
- Merge windows that end up near the same “peak” or mode



# Semantic Segmentation: Sliding Window



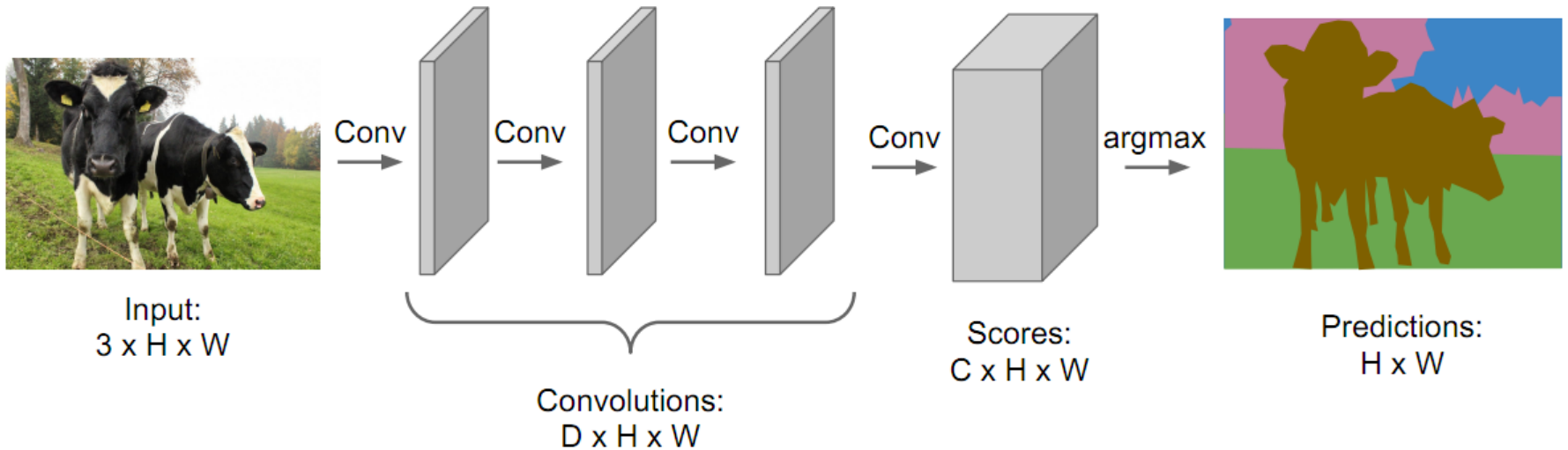
Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

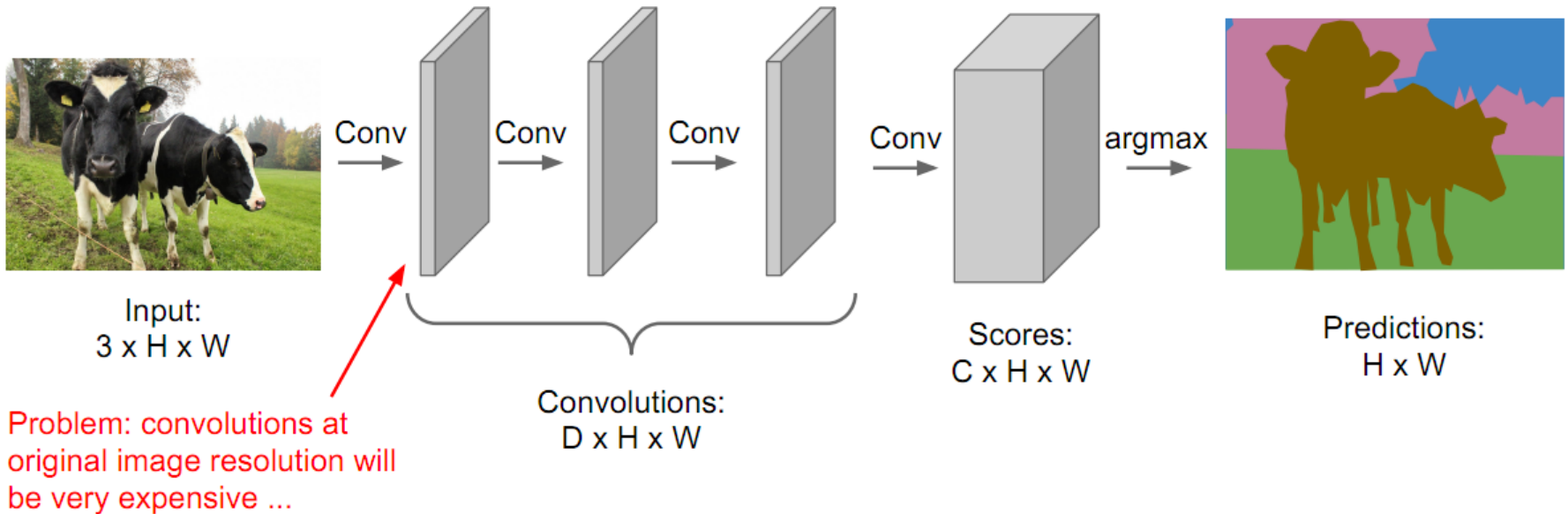
# Semantic Segmentation Idea: Fully Convolutional

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



# Semantic Segmentation Idea: Fully Convolutional

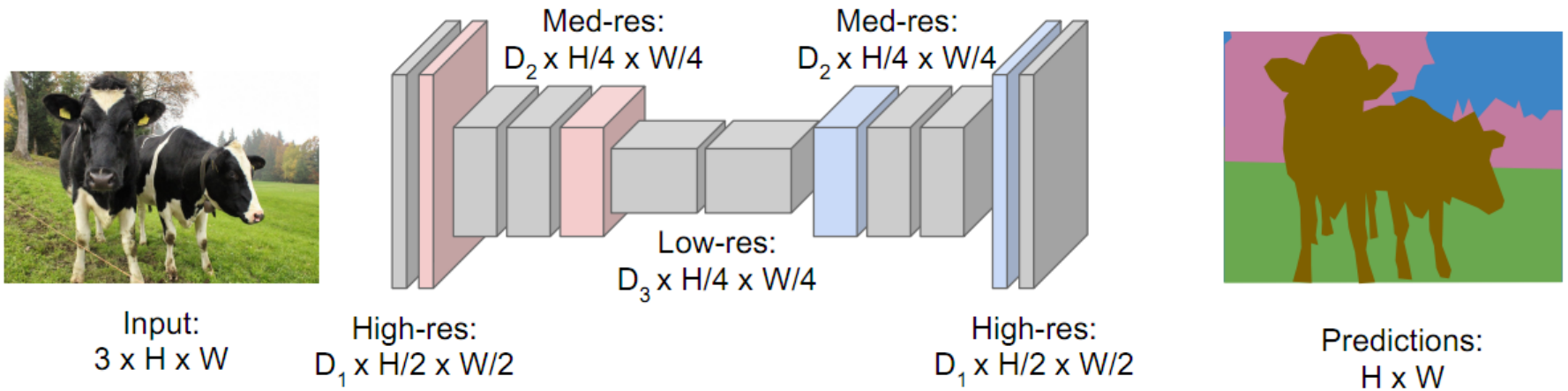
Design a network as a bunch of convolutional layers to make predictions for pixels all at once!





# Semantic Segmentation Idea: Fully Convolutional

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

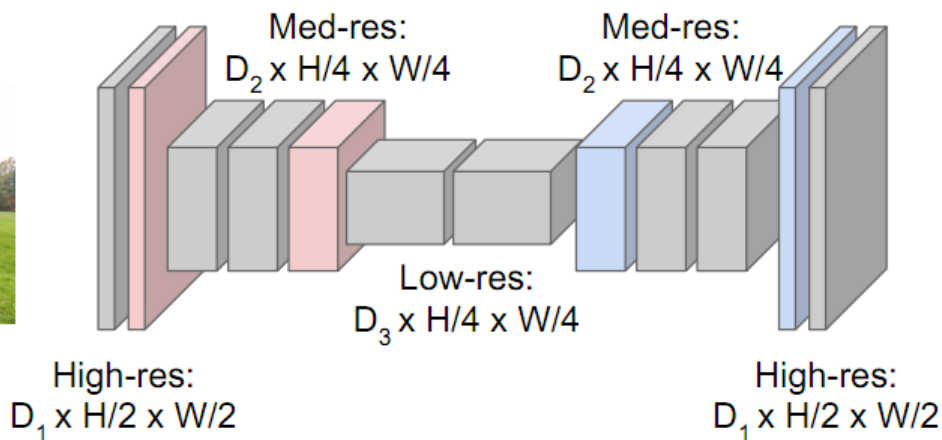
# Semantic Segmentation Idea: Fully Convolutional

**Downsampling:**  
Pooling, strided  
convolution



Input:  
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



**Upsampling:**  
???



Predictions:  
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015

Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

# In-Network upsampling: “Unpooling”

**Nearest Neighbor**

1	2
3	4



1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Input: 2 x 2

Output: 4 x 4

**“Bed of Nails”**

1	2
3	4

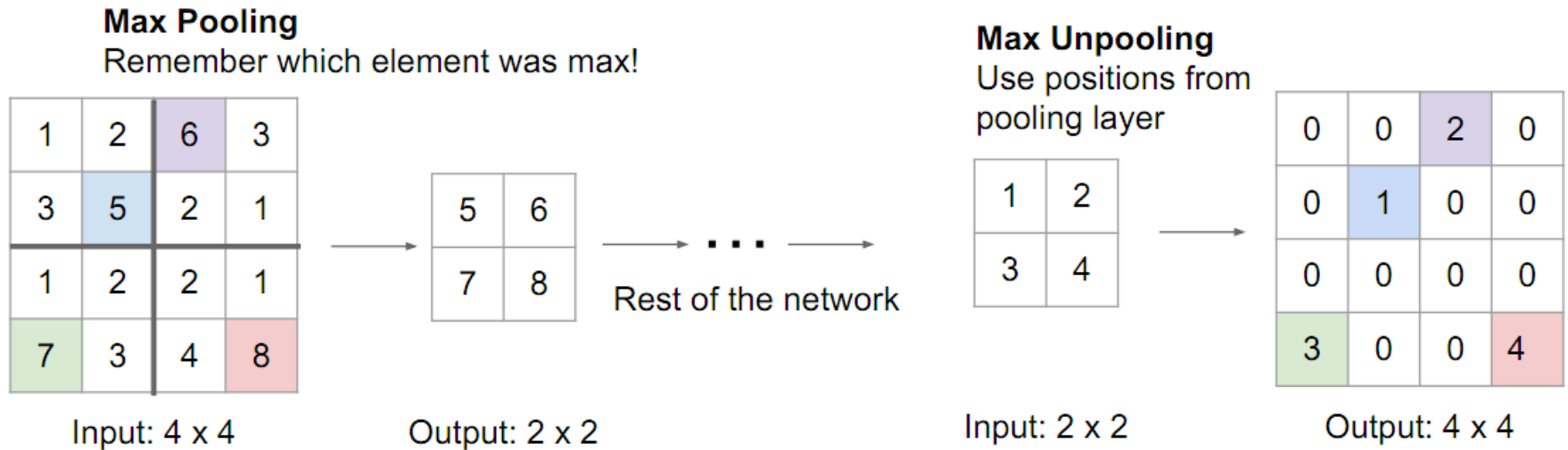


1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

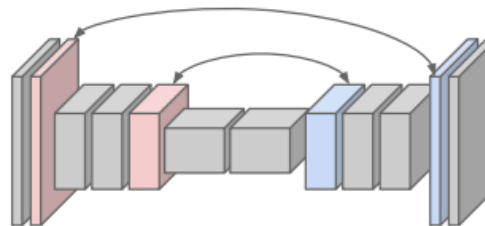
Input: 2 x 2

Output: 4 x 4

# In-Network upsampling: “Max Unpooling”

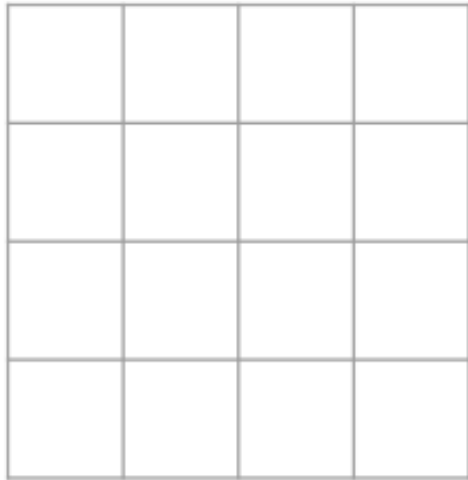


Corresponding pairs of  
downsampling and  
upsampling layers

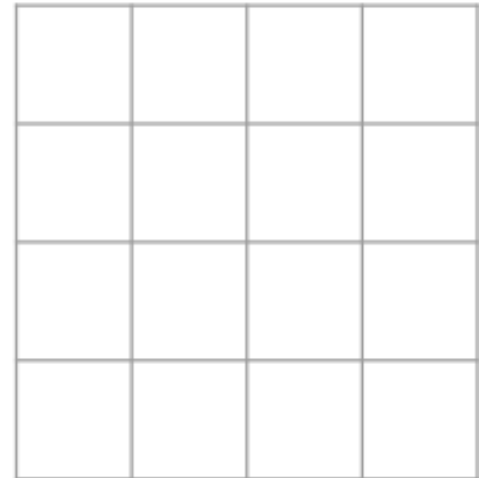


# Learnable Upsampling: Transpose Convolution

**Recall:** Typical 3 x 3 convolution, stride 1 pad 1



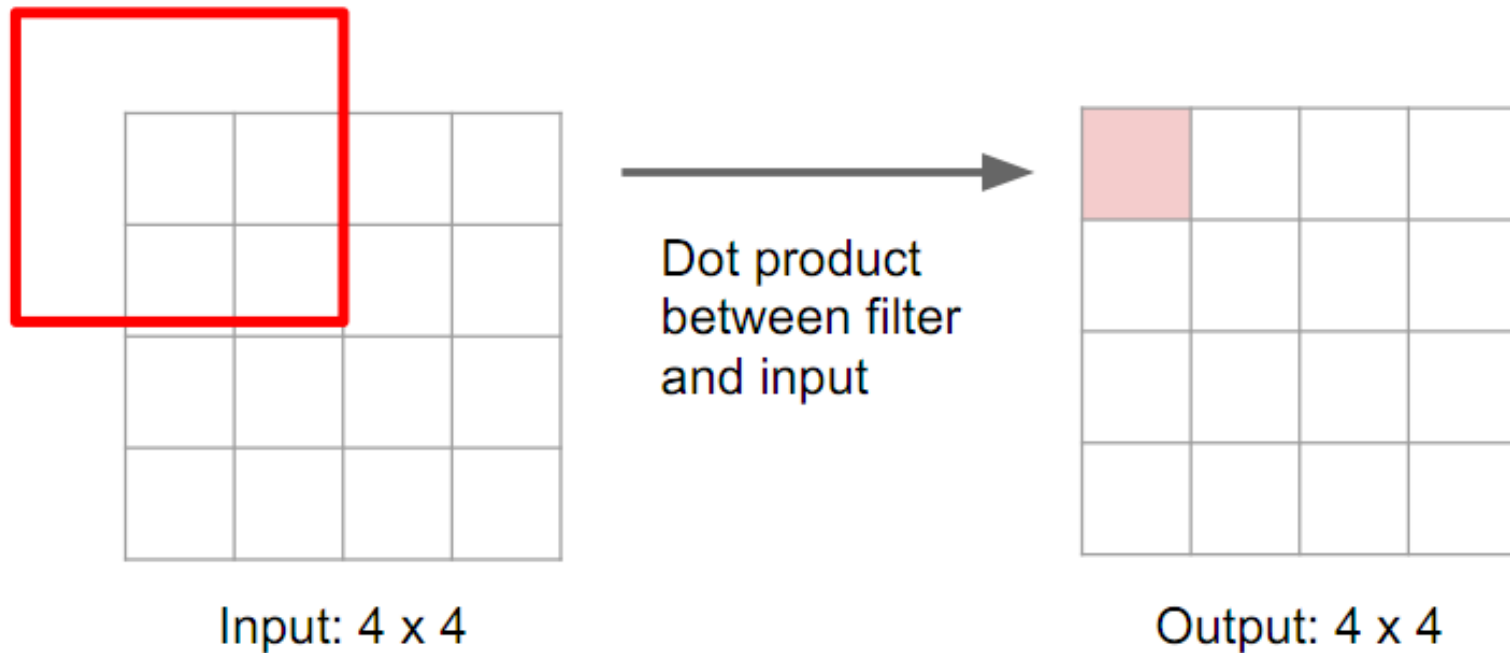
Input: 4 x 4



Output: 4 x 4

# Learnable Upsampling: Transpose Convolution

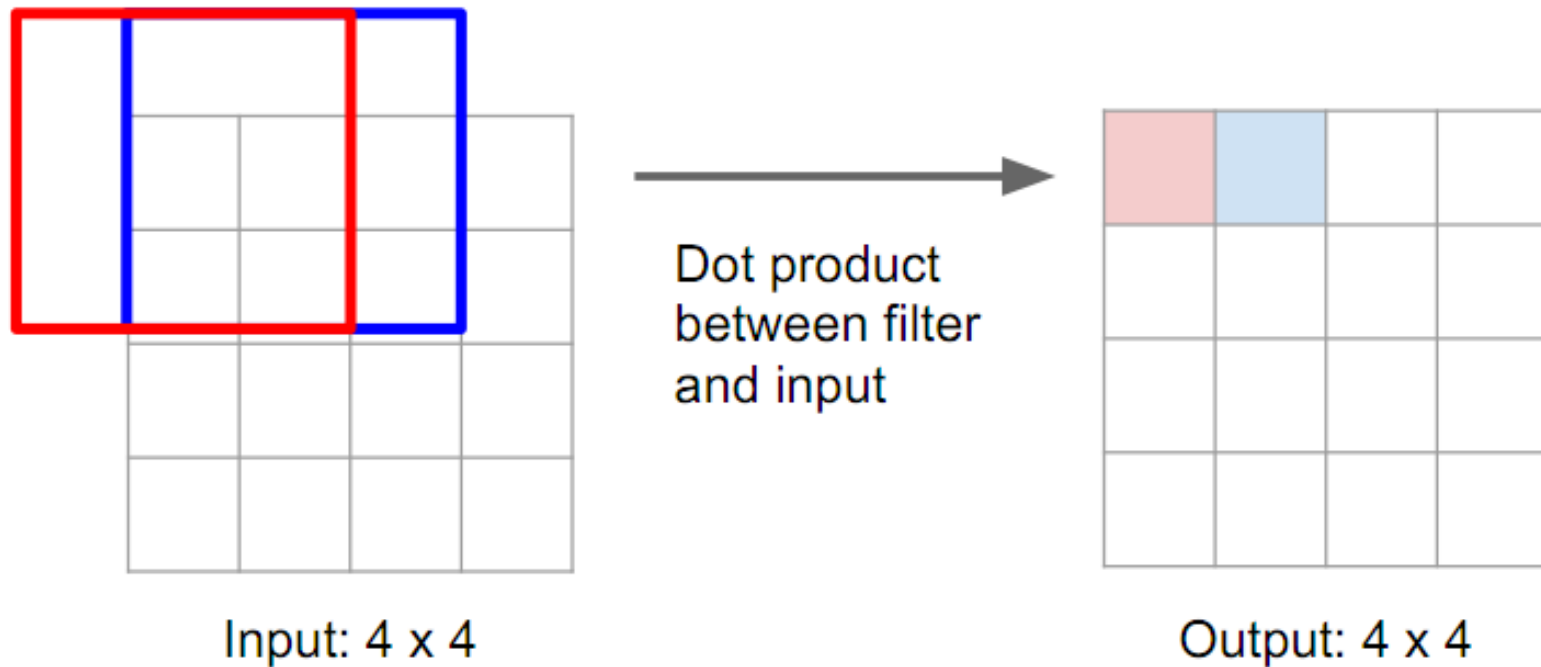
**Recall:** Normal 3 x 3 convolution, stride 1 pad 1





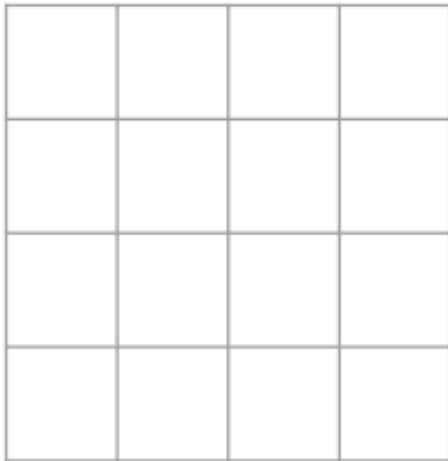
# Learnable Upsampling: Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 1 pad 1

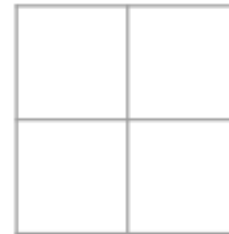


# Learnable Upsampling: Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 2 pad 1



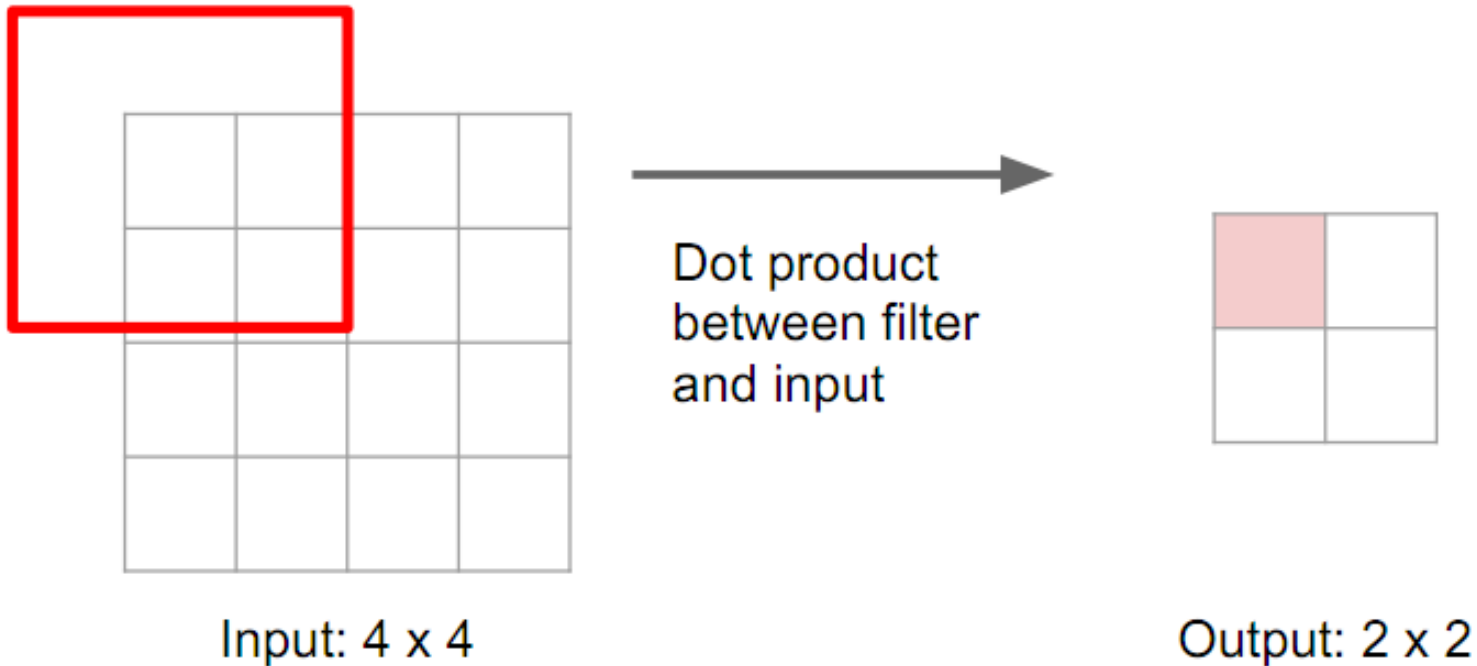
Input: 4 x 4



Output: 2 x 2

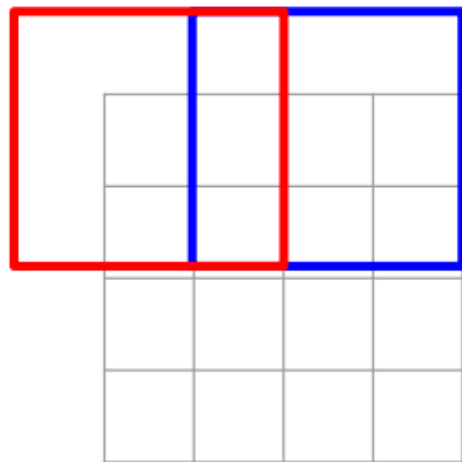
# Learnable Upsampling: Transpose Convolution

**Recall:** Normal 3 x 3 convolution, stride 2 pad 1



# Learnable Upsampling: Transpose Convolution

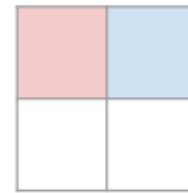
**Recall:** Normal 3 x 3 convolution, stride 2 pad 1



Input: 4 x 4



Dot product  
between filter  
and input



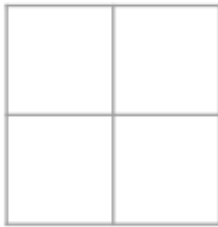
Output: 2 x 2

Filter moves 2 pixels in  
the input for every one  
pixel in the output

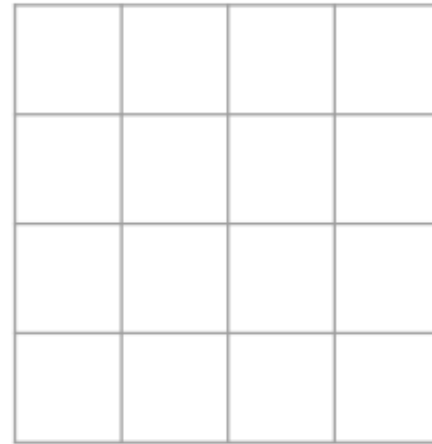
Stride gives ratio between  
movement in input and  
output

# Learnable Upsampling: Transpose Convolution

3 x 3 **transpose** convolution, stride 2 pad 1



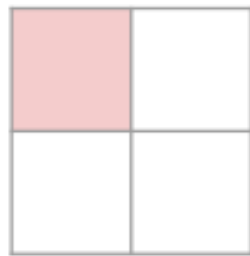
Input: 2 x 2



Output: 4 x 4

# Learnable Upsampling: Transpose Convolution

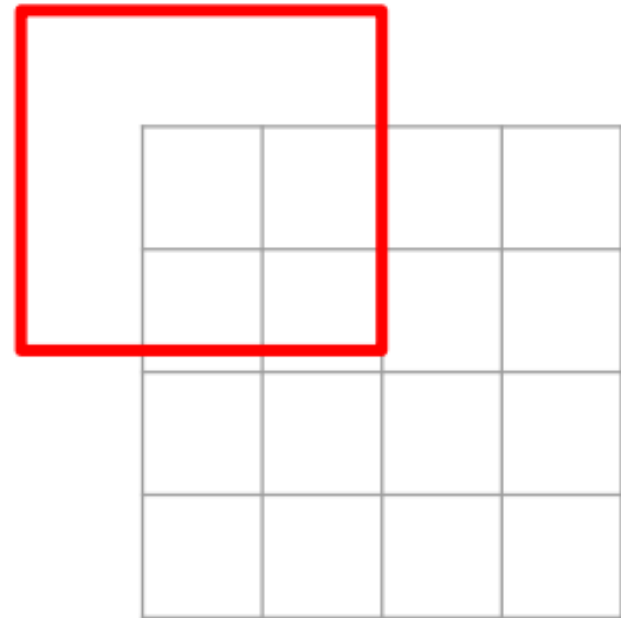
3 x 3 **transpose** convolution, stride 2 pad 1



Input: 2 x 2



Input gives  
weight for  
filter

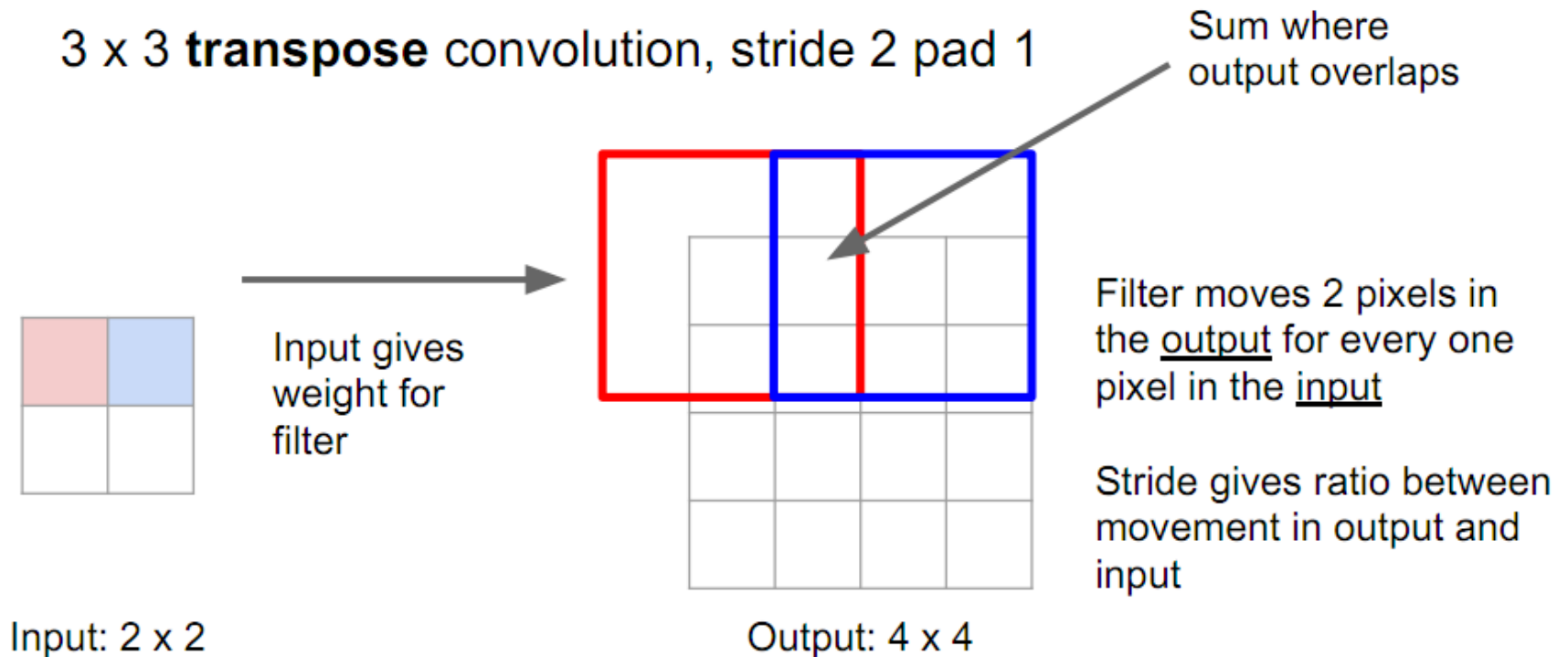


Output: 4 x 4



# Learnable Upsampling: Transpose Convolution

3 x 3 **transpose** convolution, stride 2 pad 1

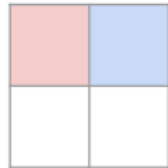


# Learnable Upsampling: Transpose Convolution

## Other names:

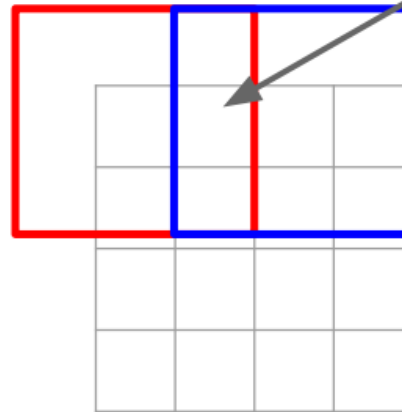
- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution

3 x 3 **transpose** convolution, stride 2 pad 1



Input: 2 x 2

Input gives weight for filter



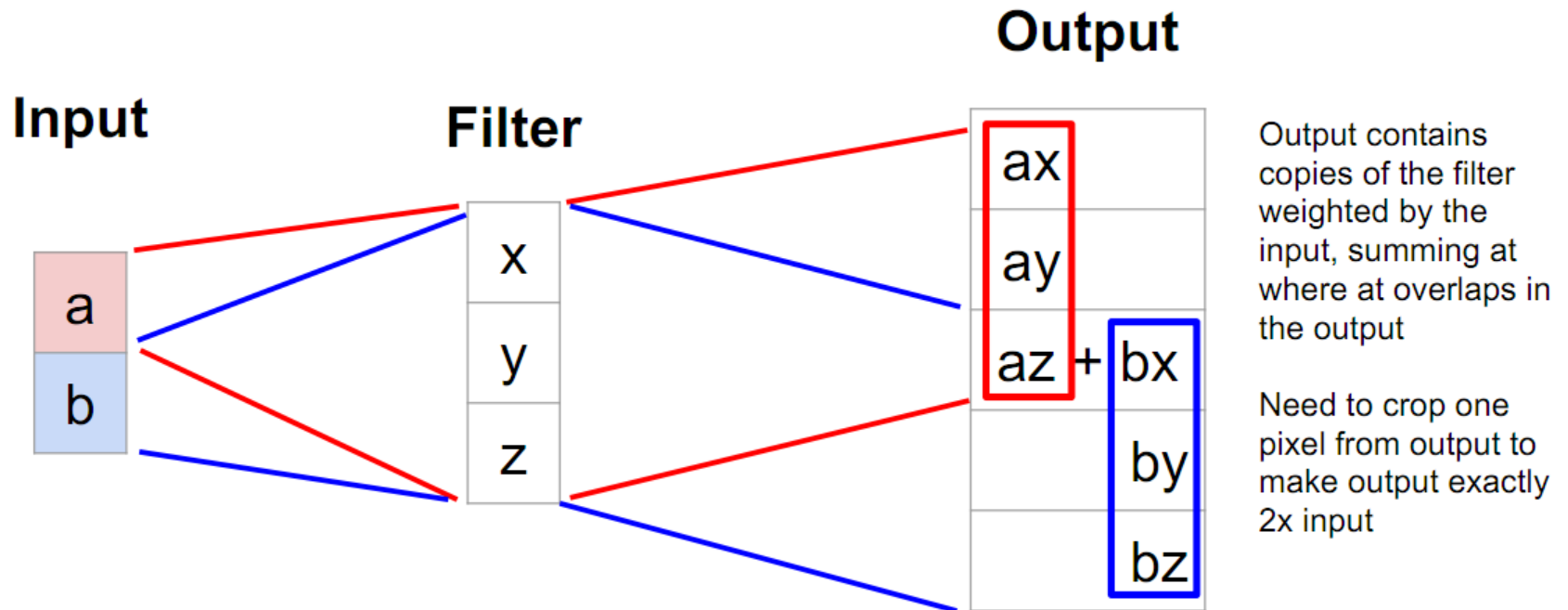
Output: 4 x 4

Sum where output overlaps

Filter moves 2 pixels in the output for every one pixel in the input

Stride gives ratio between movement in output and input

# Learnable Upsampling: 1D Example



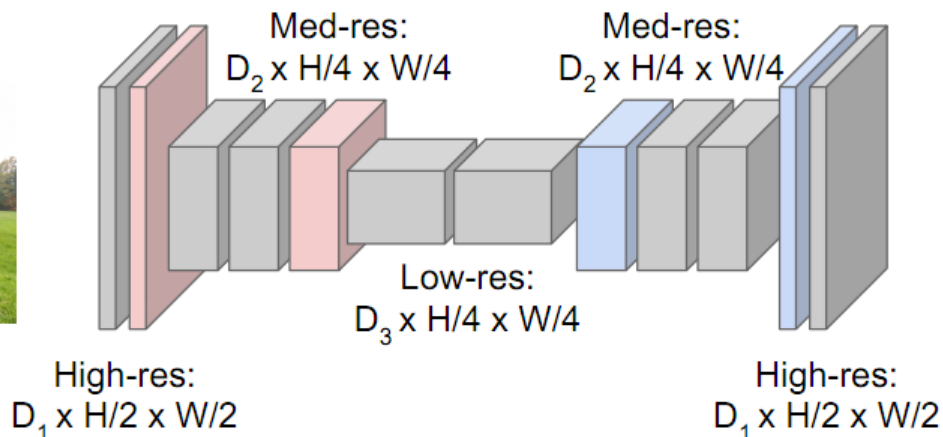
# Semantic Segmentation Idea: Fully Convolutional

**Downsampling:**  
Pooling, strided  
convolution



Input:  
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!

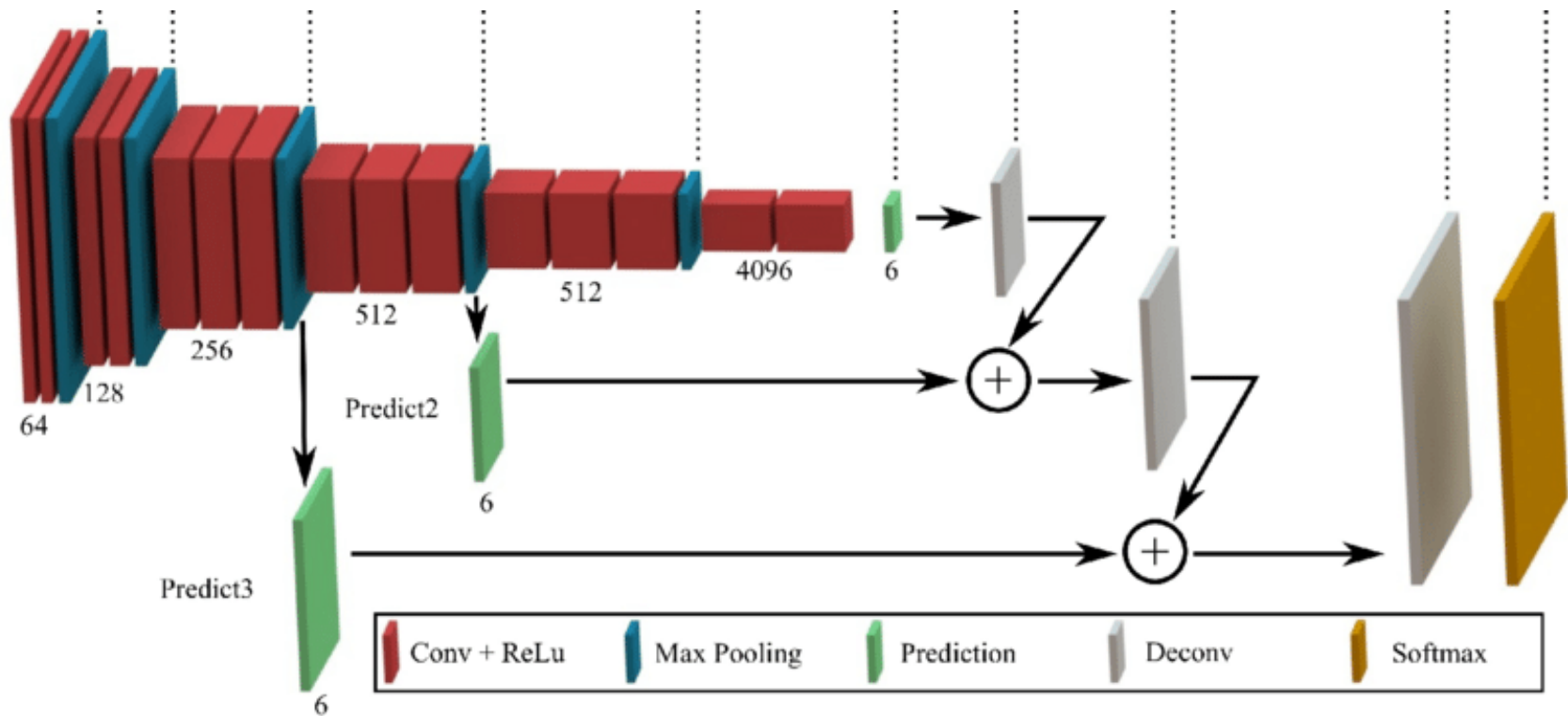


**Upsampling:**  
Unpooling or strided  
transpose convolution

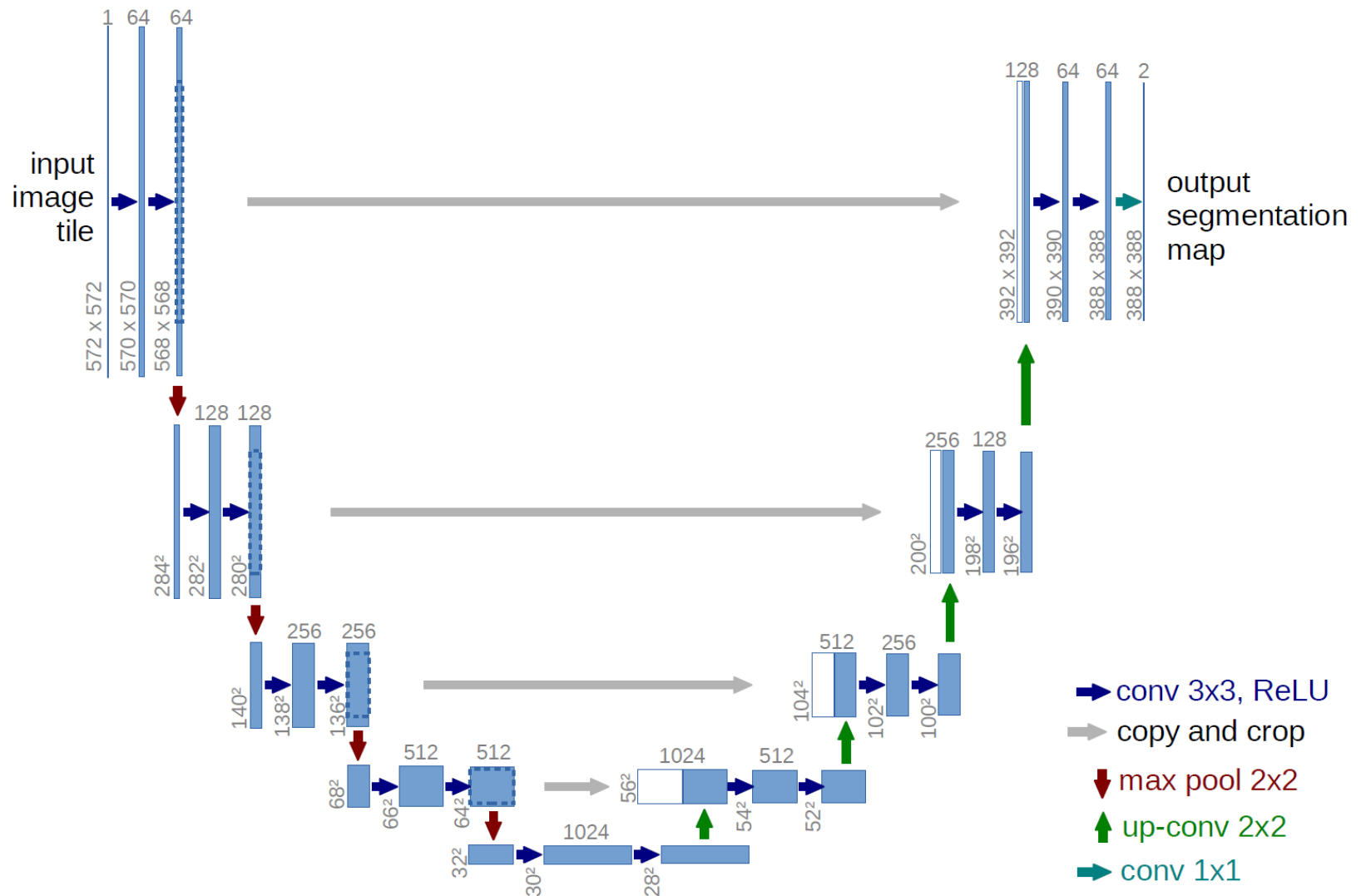


Predictions:  
 $H \times W$

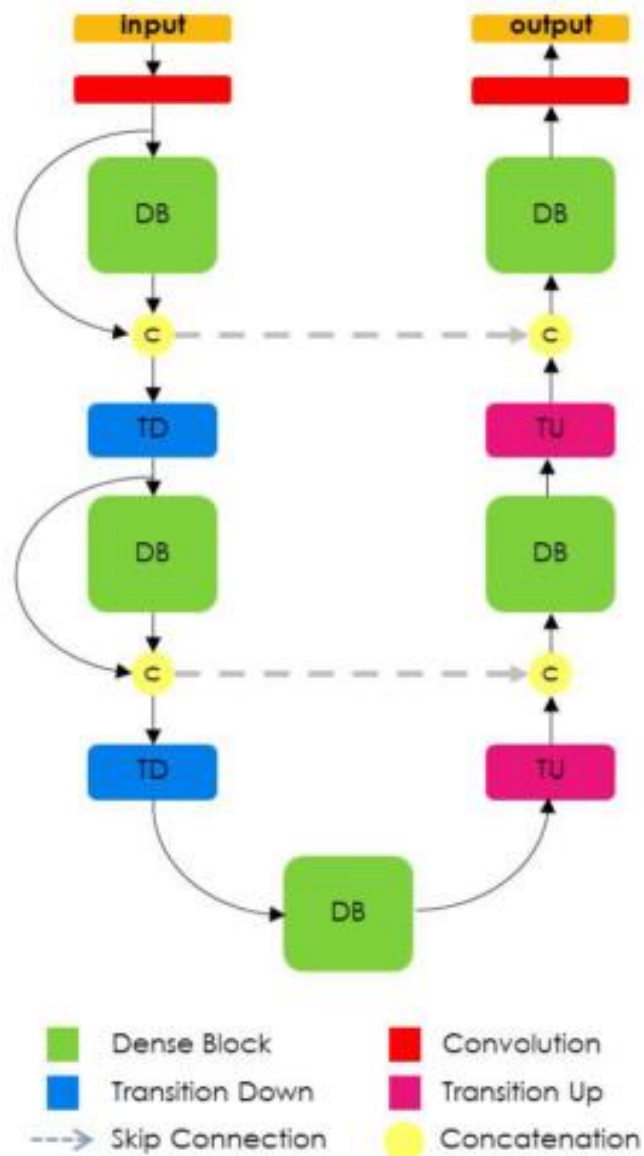
# FCN with 2 skip connections



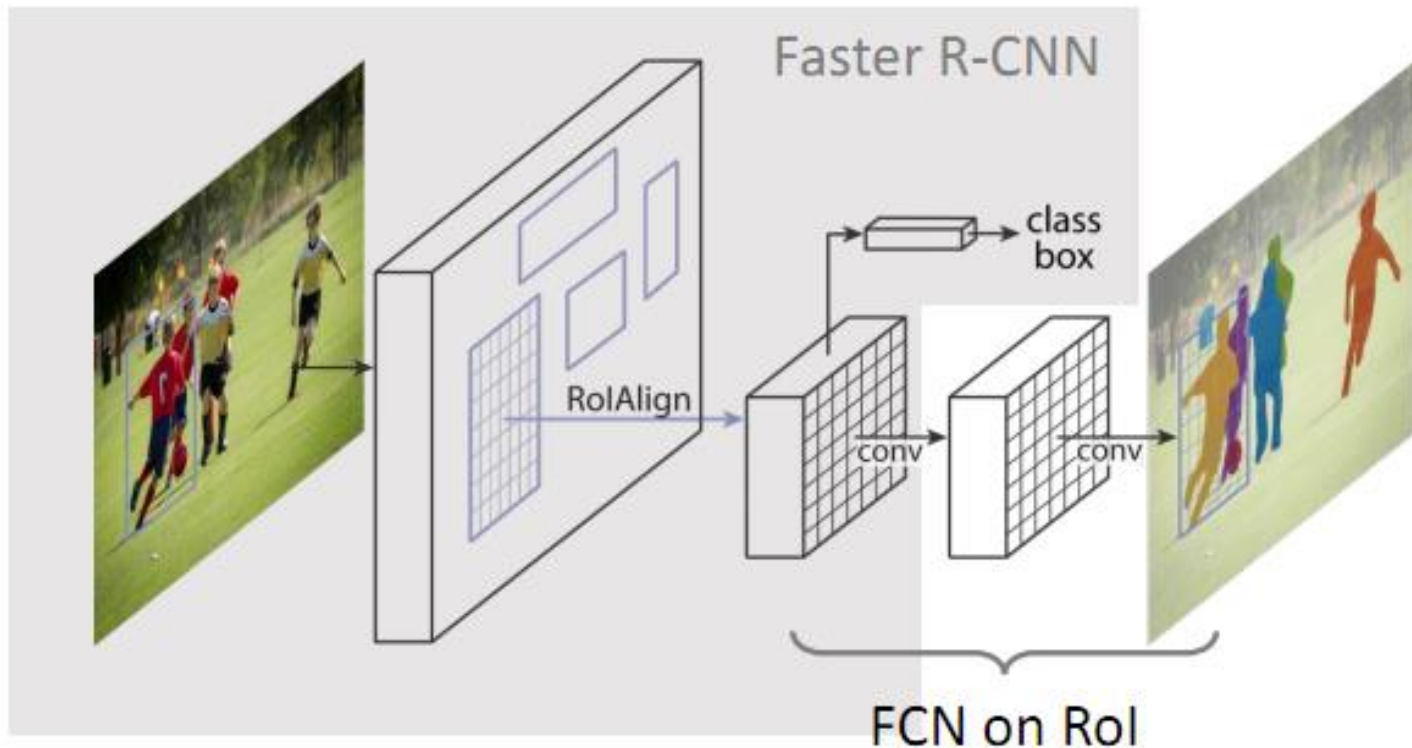
# U-Net



# Tiramisu: Fully Convolutional DenseNets



# Mask R-CNN



Mask RCNN =

- 1. Object detector using Faster RCNN +
- 2. fully convolutional network (FCN) on region of interest (RoI)