

# Deep Tracking & Flow

Instructor - Simon Lucey

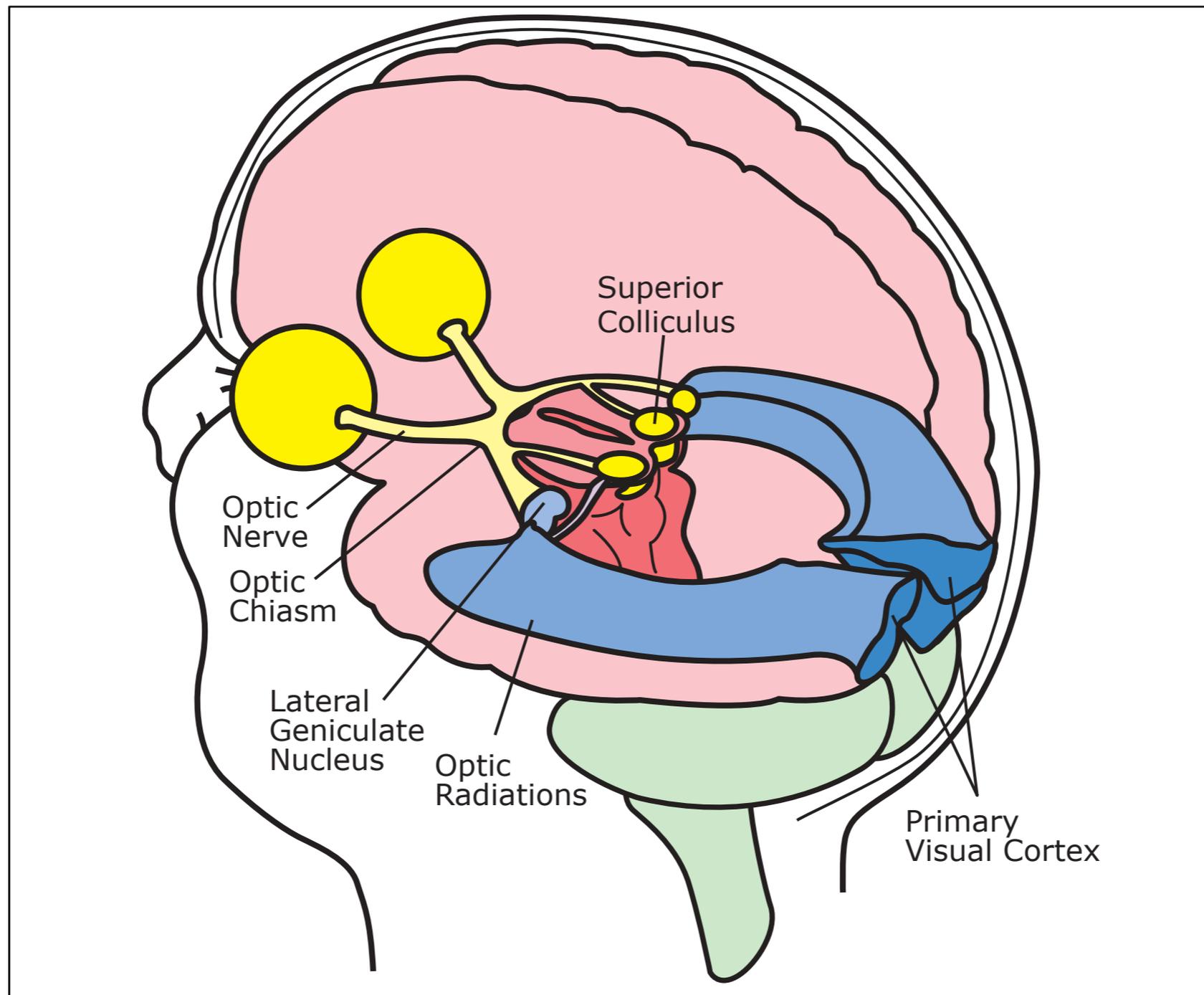
**16-623 - Designing Computer Vision Apps**

# Today

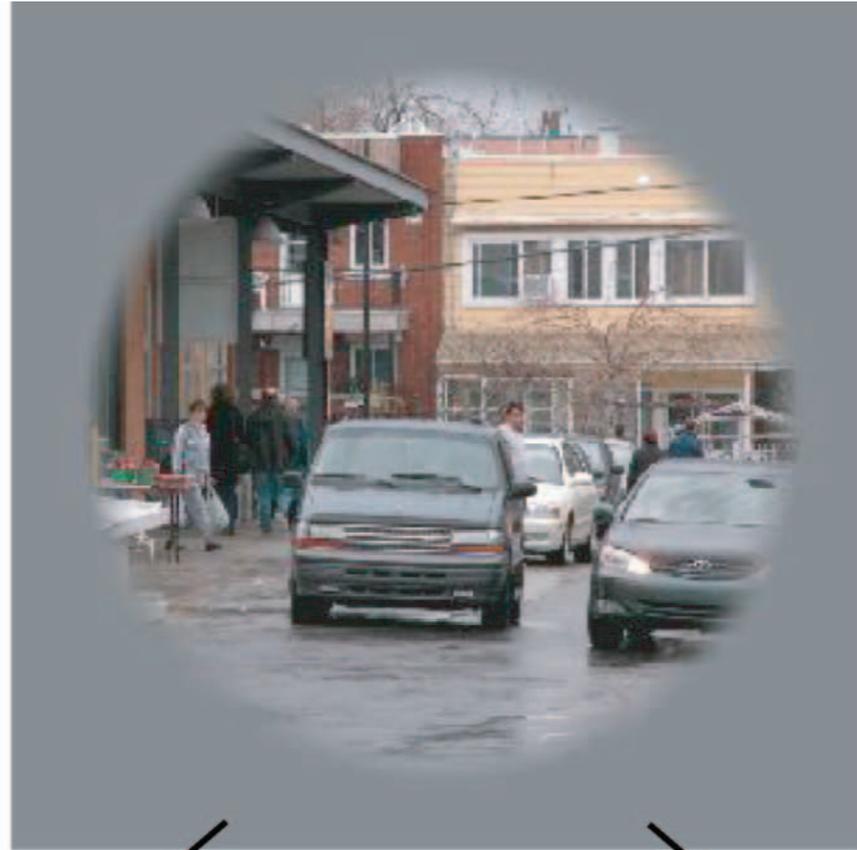
---

- Deep Features
- Deep Tracking
- Deep Flow

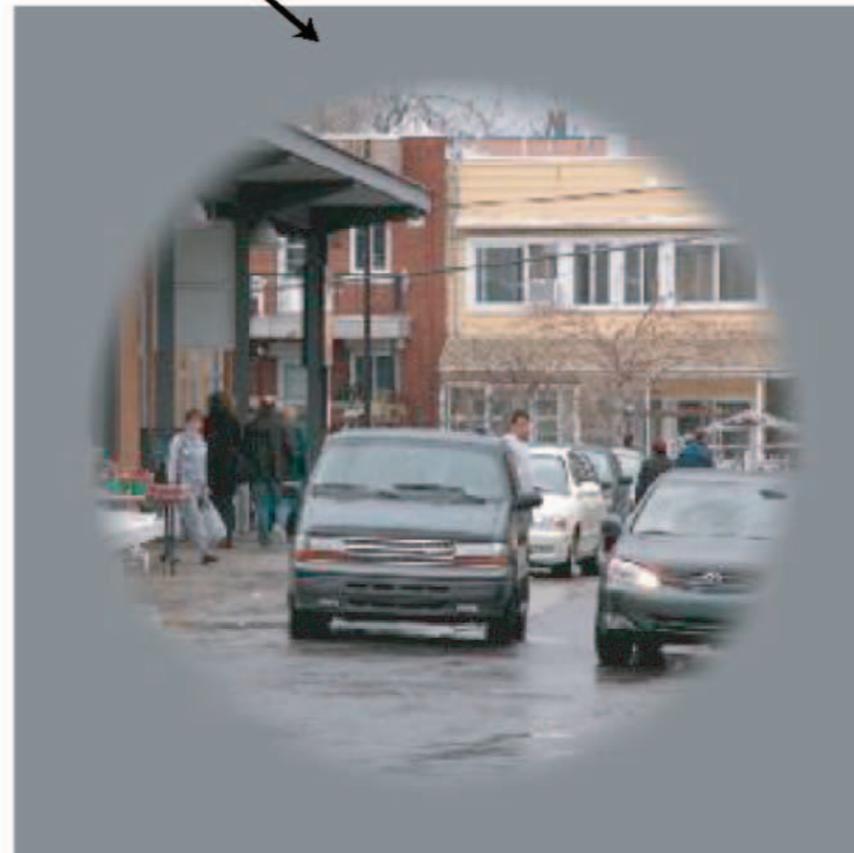
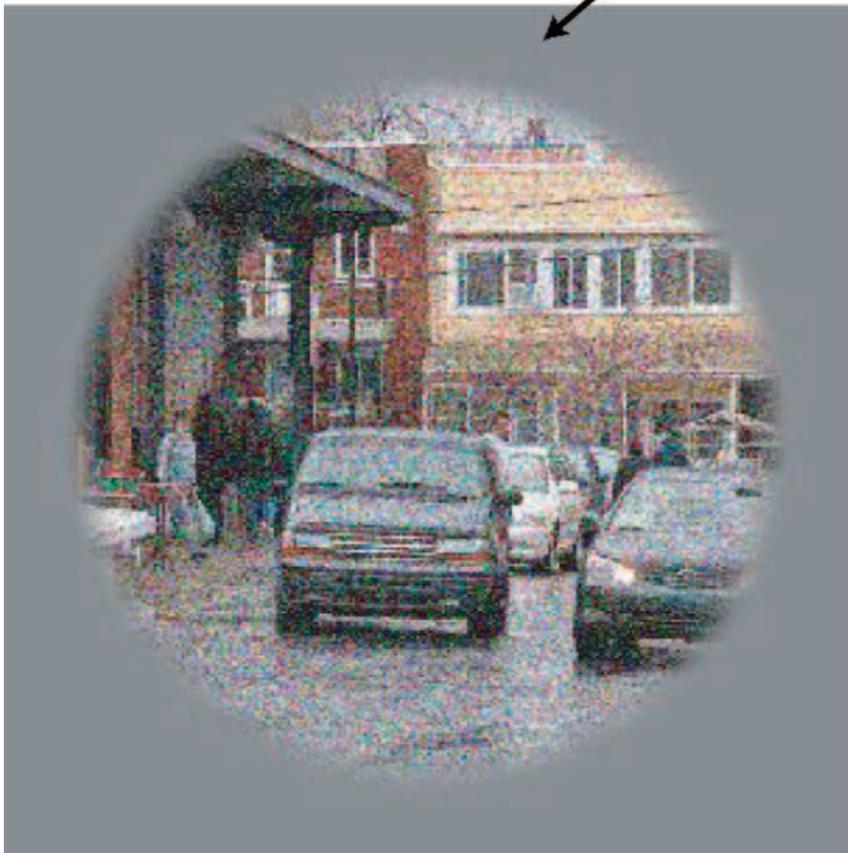
# Primary Visual Cortex



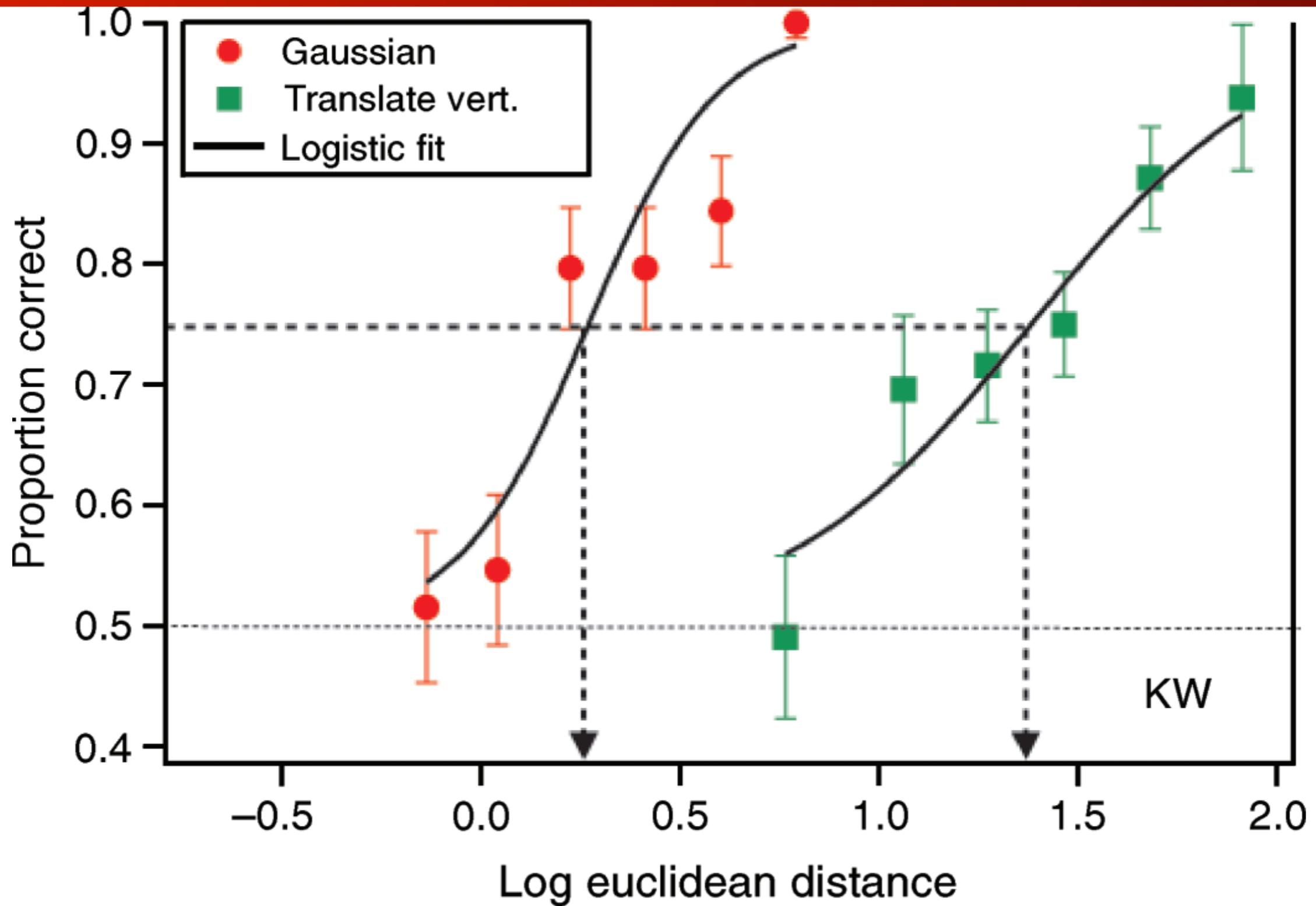
# Spatial Sensitivity



Kingdom, Field, Olmos, 2007

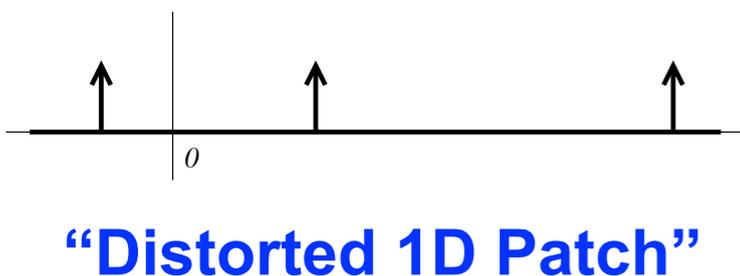
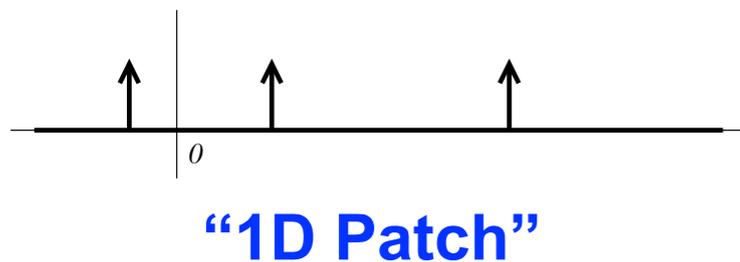


# Spatial Sensitivity



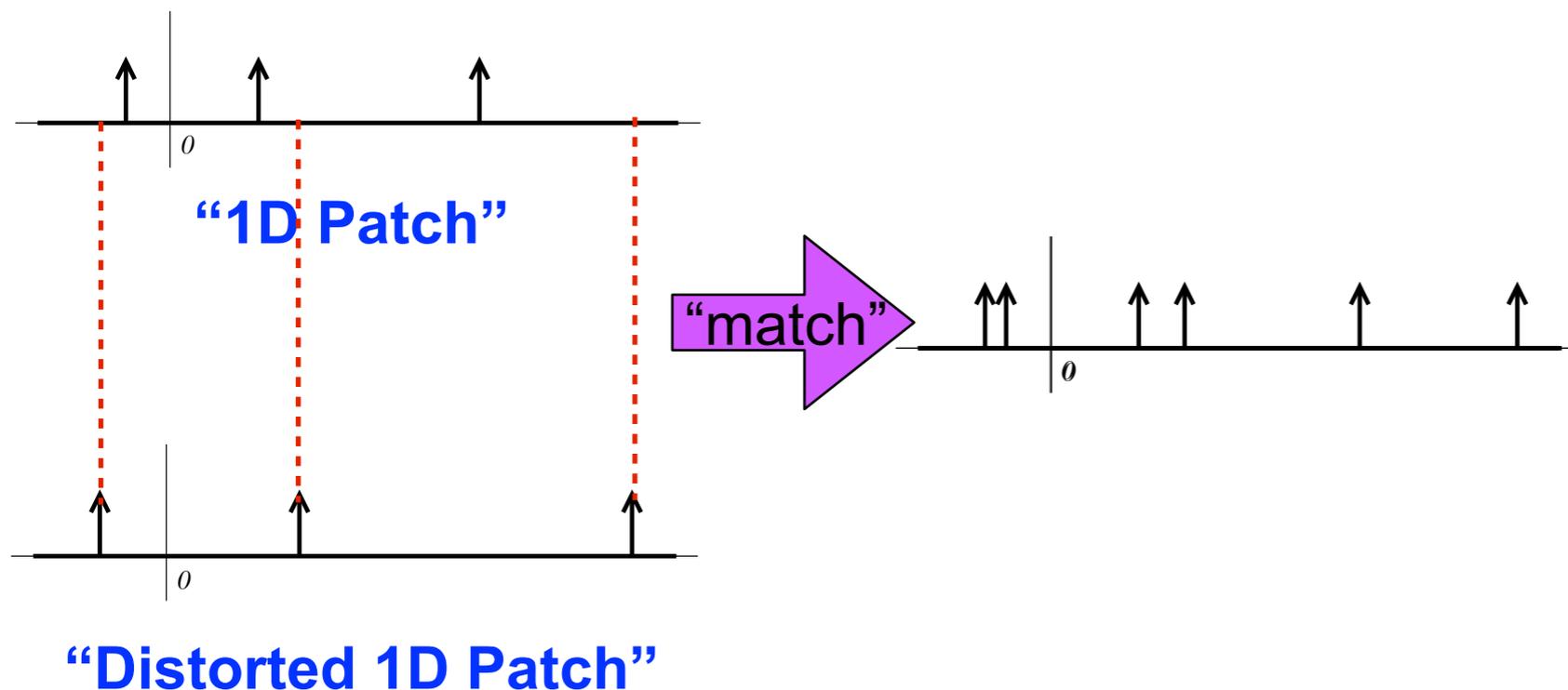
# Handling Geometric Distortion

- As pointed out in seminal work by Berg and Malik (CVPR'01) the effectiveness of SSD will degrade with significant viewpoint change.
- Two options to match local image patches:-



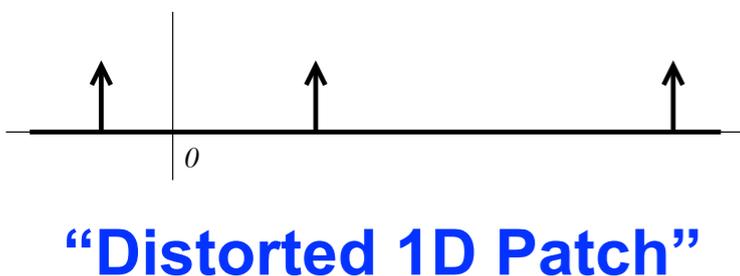
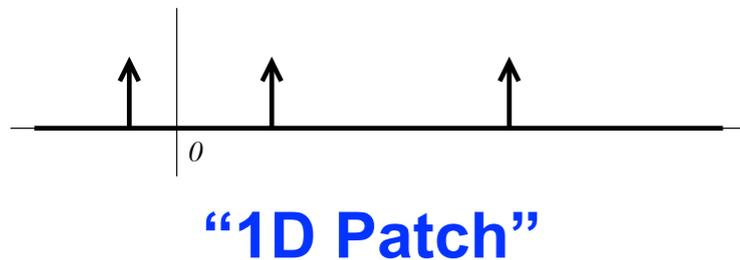
# Handling Geometric Distortion

- As pointed out in seminal work by Berg and Malik (CVPR'01) the effectiveness of SSD will degrade with significant viewpoint change.
- Two options to match local image patches:-



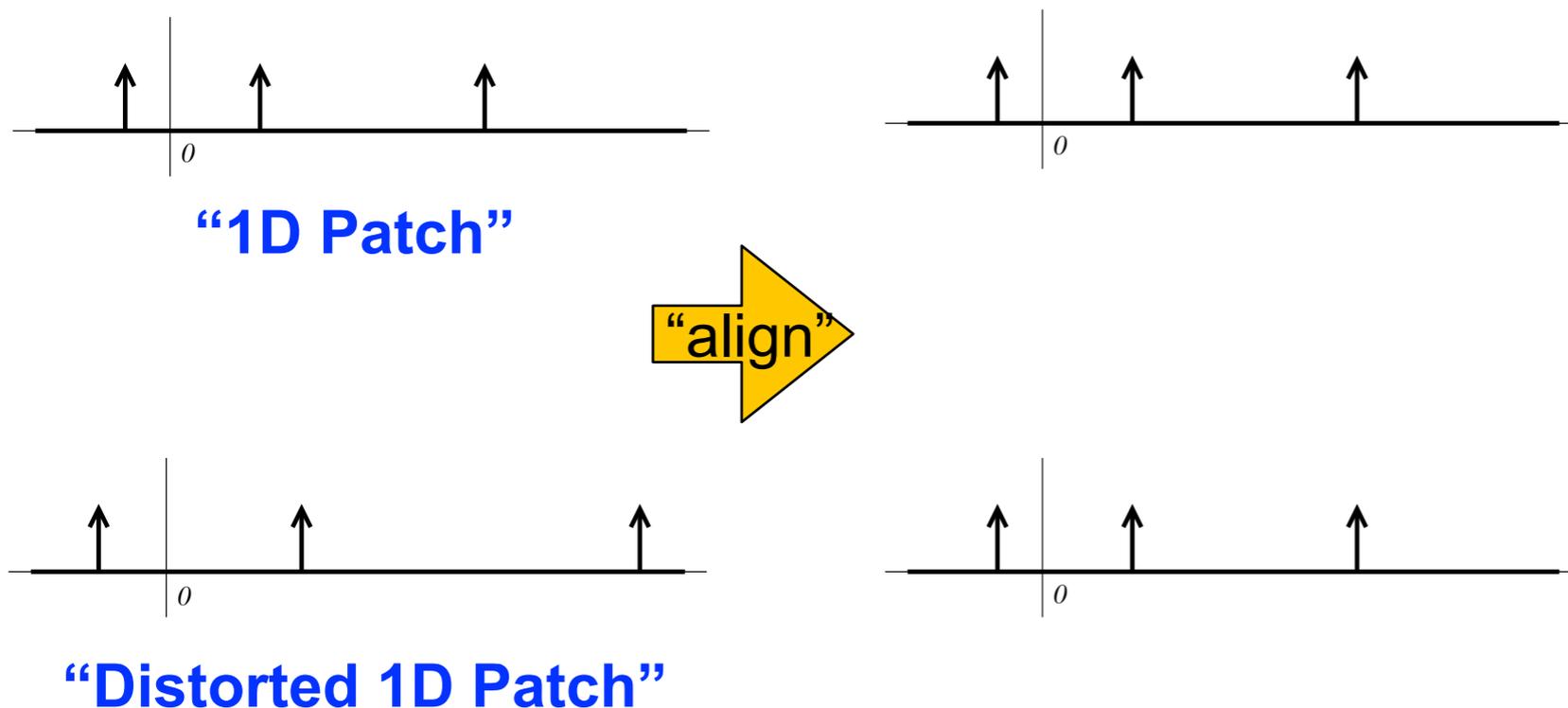
# Handling Geometric Distortion

- As pointed out in seminal work by Berg and Malik (CVPR'01) the effectiveness of SSD will degrade with significant viewpoint change.
- Two options to match local image patches:-
  1. simultaneously estimate the distortion and position of matching patch



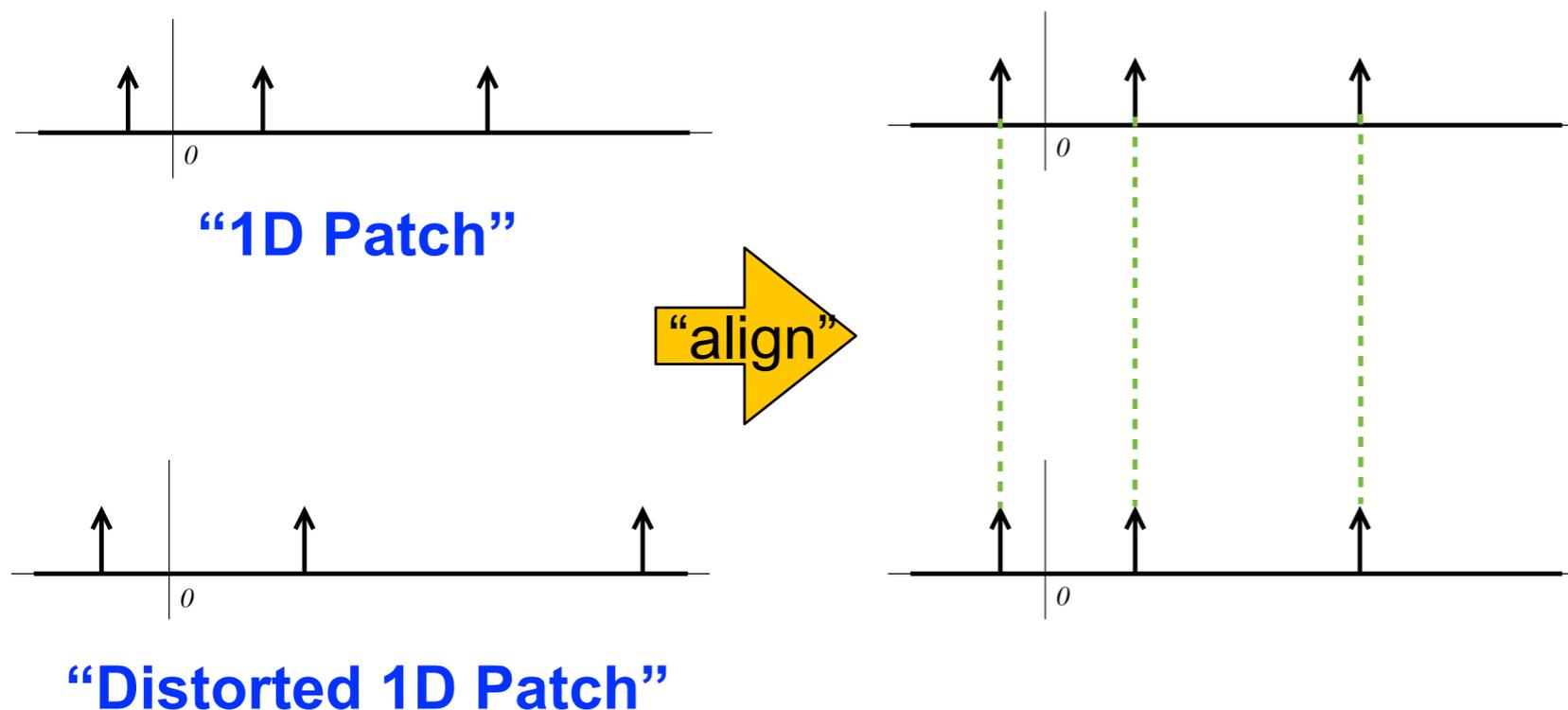
# Handling Geometric Distortion

- As pointed out in seminal work by Berg and Malik (CVPR'01) the effectiveness of SSD will degrade with significant viewpoint change.
- Two options to match local image patches:-
  1. simultaneously estimate the distortion and position of matching patch



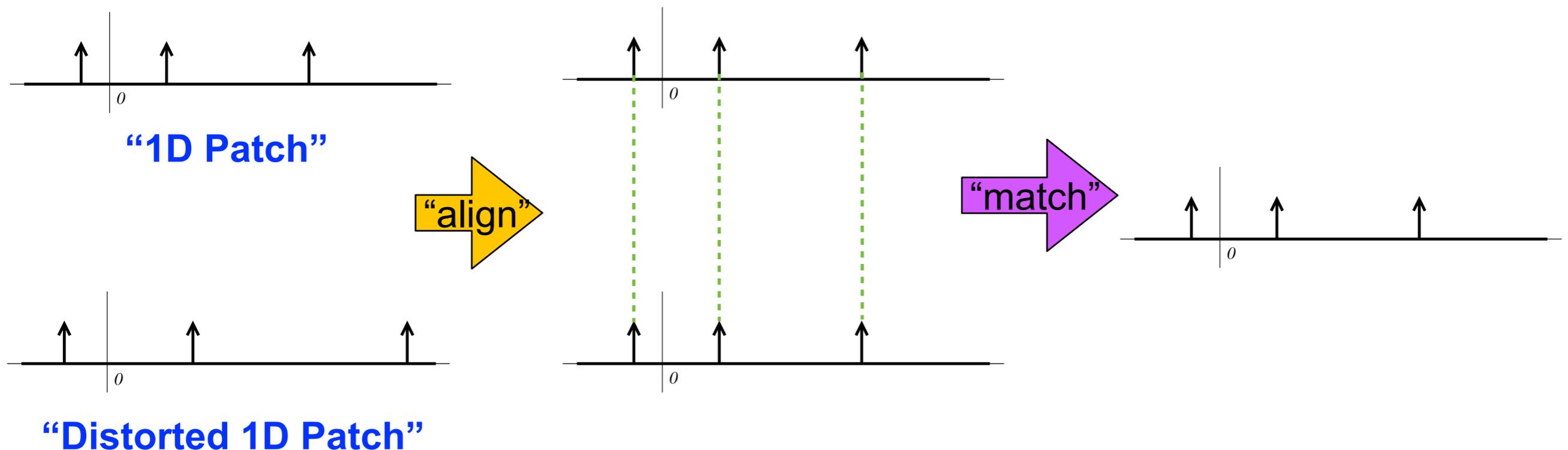
# Handling Geometric Distortion

- As pointed out in seminal work by Berg and Malik (CVPR'01) the effectiveness of SSD will degrade with significant viewpoint change.
- Two options to match local image patches:-
  1. simultaneously estimate the distortion and position of matching patch



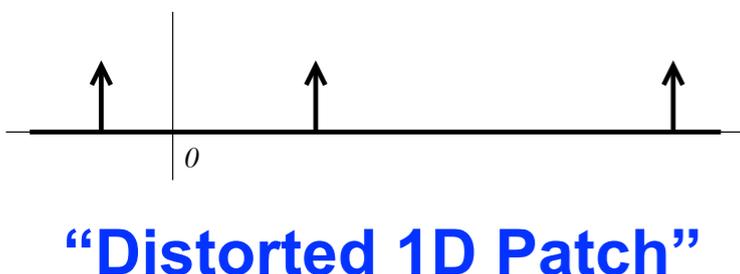
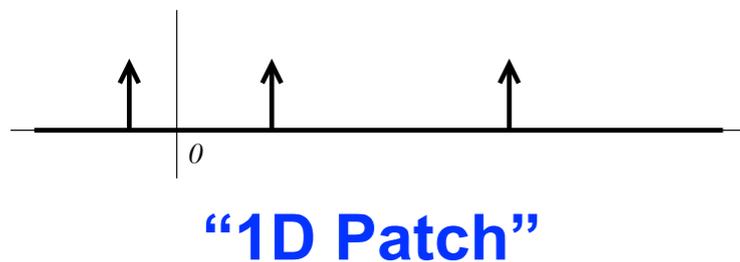
# Handling Geometric Distortion

- As pointed out in seminal work by Berg and Malik (CVPR'01) the effectiveness of SSD will degrade with significant viewpoint change.
- Two options to match local image patches:-
  1. simultaneously estimate the distortion and position of matching patch



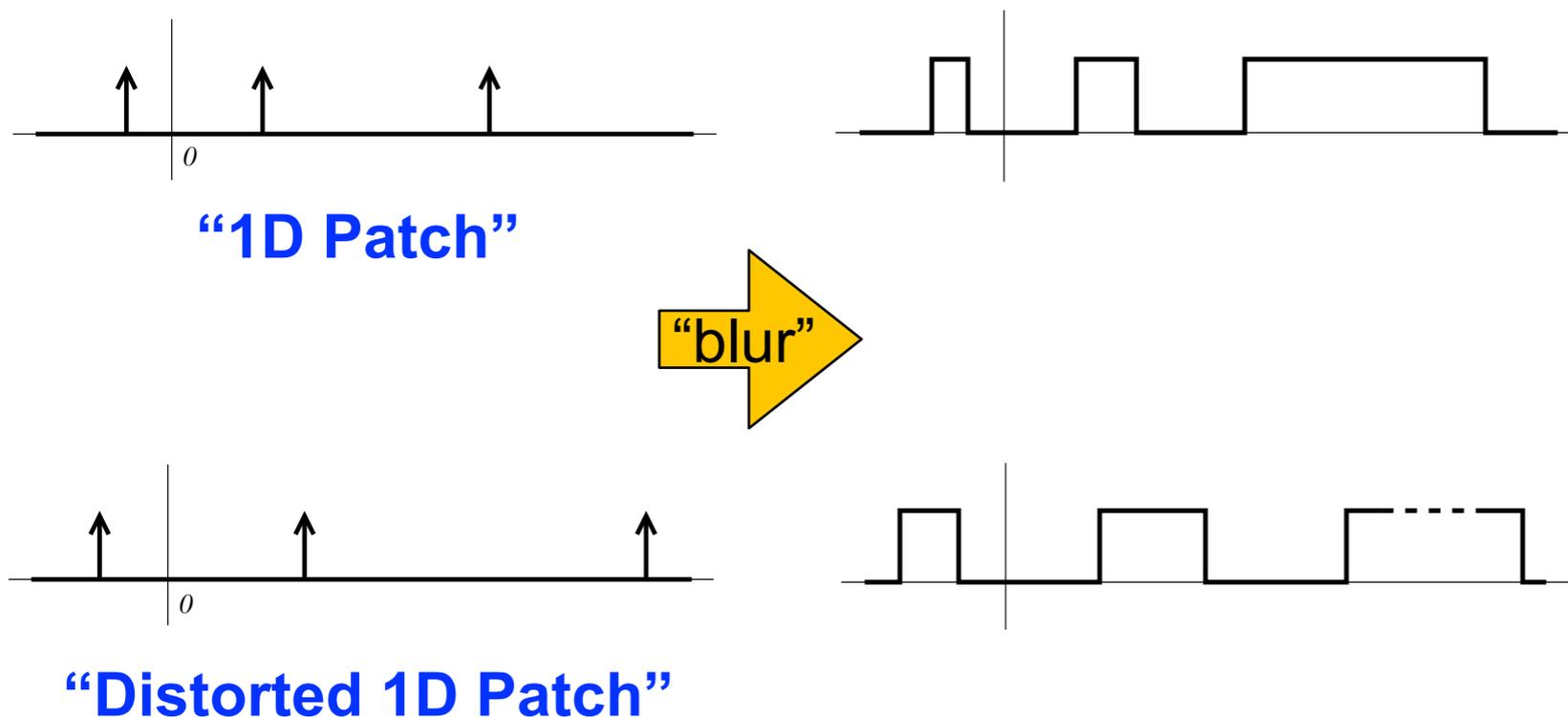
# Handling Geometric Distortion

- As pointed out in seminal work by Berg and Malik (CVPR'01) the effectiveness of SSD will degrade with significant viewpoint and/or illumination change.
- Two options to match patches:-
  1. simultaneously estimate the distortion and position of matching patch.
  2. to “blur” the template window performing matching coarse-to-fine.



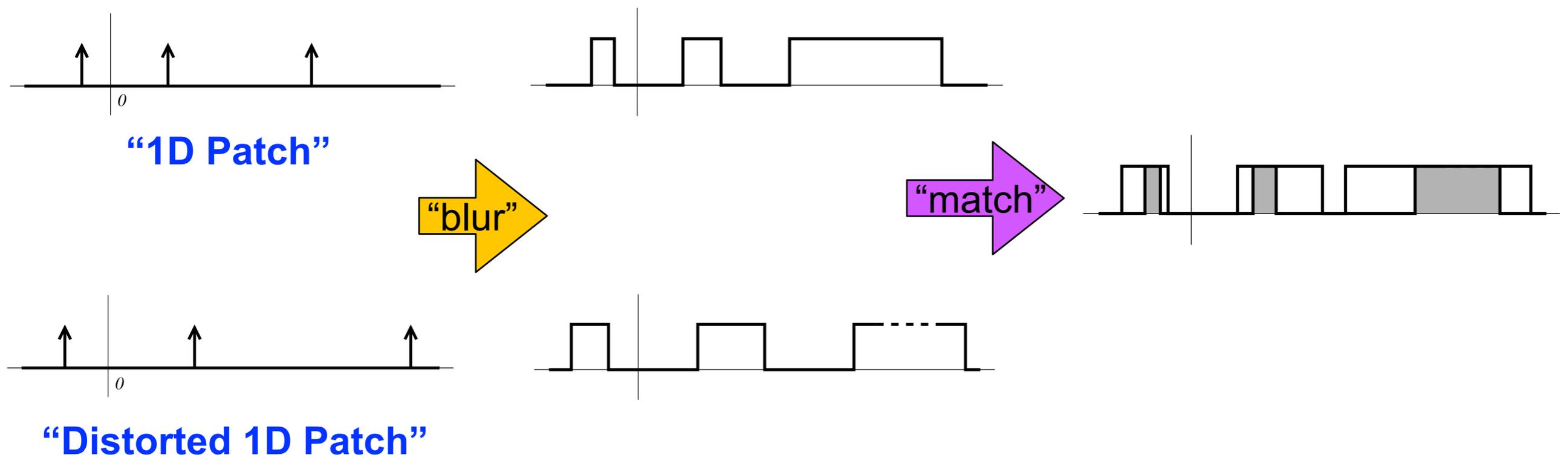
# Handling Geometric Distortion

- As pointed out in seminal work by Berg and Malik (CVPR'01) the effectiveness of SSD will degrade with significant viewpoint and/or illumination change.
- Two options to match patches:-
  1. simultaneously estimate the distortion and position of matching patch.
  2. to “blur” the template window performing matching coarse-to-fine.



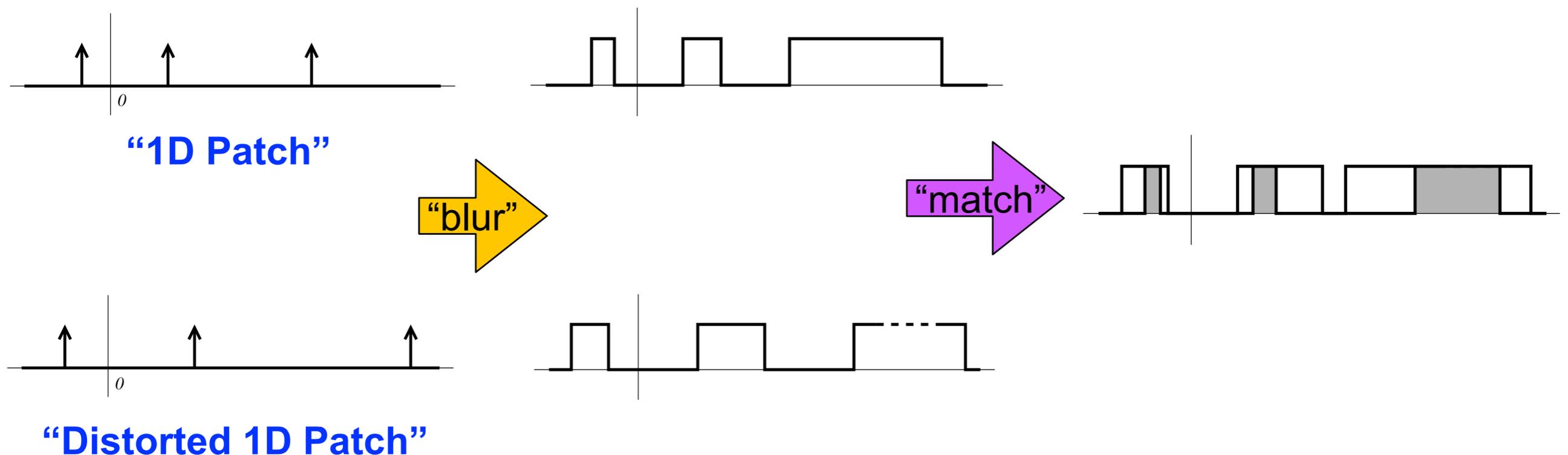
# Handling Geometric Distortion

- As pointed out in seminal work by Berg and Malik (CVPR'01) the effectiveness of SSD will degrade with significant viewpoint and/or illumination change.
- Two options to match patches:-
  1. simultaneously estimate the distortion and position of matching patch.
  2. to “blur” the template window performing matching coarse-to-fine.



# Handling Geometric Distortion

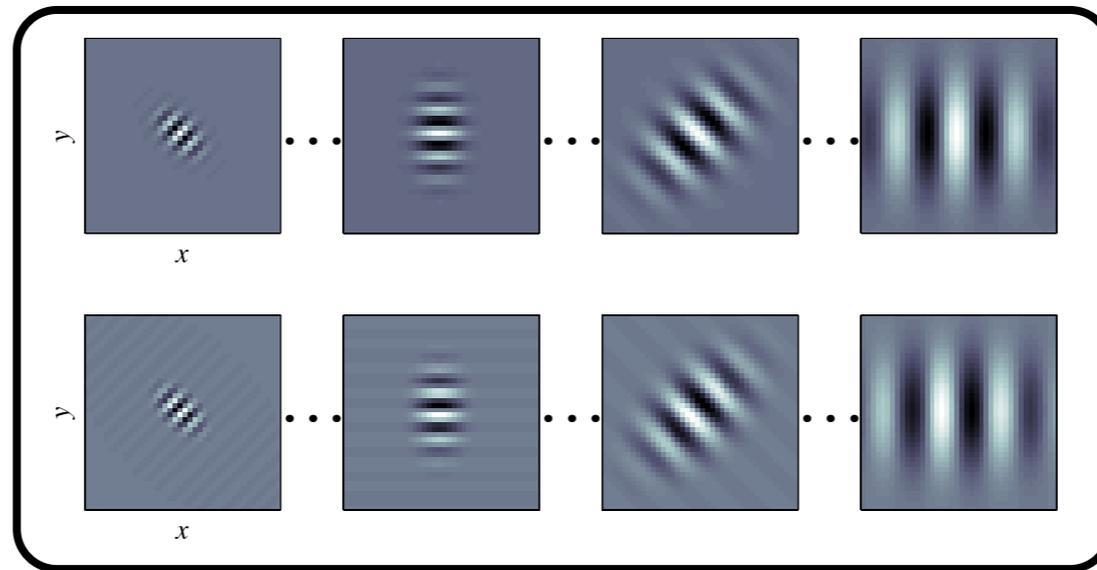
- As pointed out in seminal work by Berg and Malik (CVPR'01) the effectiveness of SSD will degrade with significant viewpoint and/or illumination change.
- Two options to match patches:-
  1. simultaneously estimate the distortion and position of matching patch.
  2. to “blur” the template window performing matching coarse-to-fine.



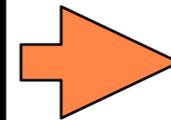
***Option 2 is attractive, low computational cost!***

# Sparseness and Positiveness

- Blurring only works if the signals being matched are sparse and positive.
- Unfortunately natural images are neither.
- Combination of oriented filter banks and rectification can remedy this problem with little loss in performance.



e.g., oriented gradients, Gabor filters

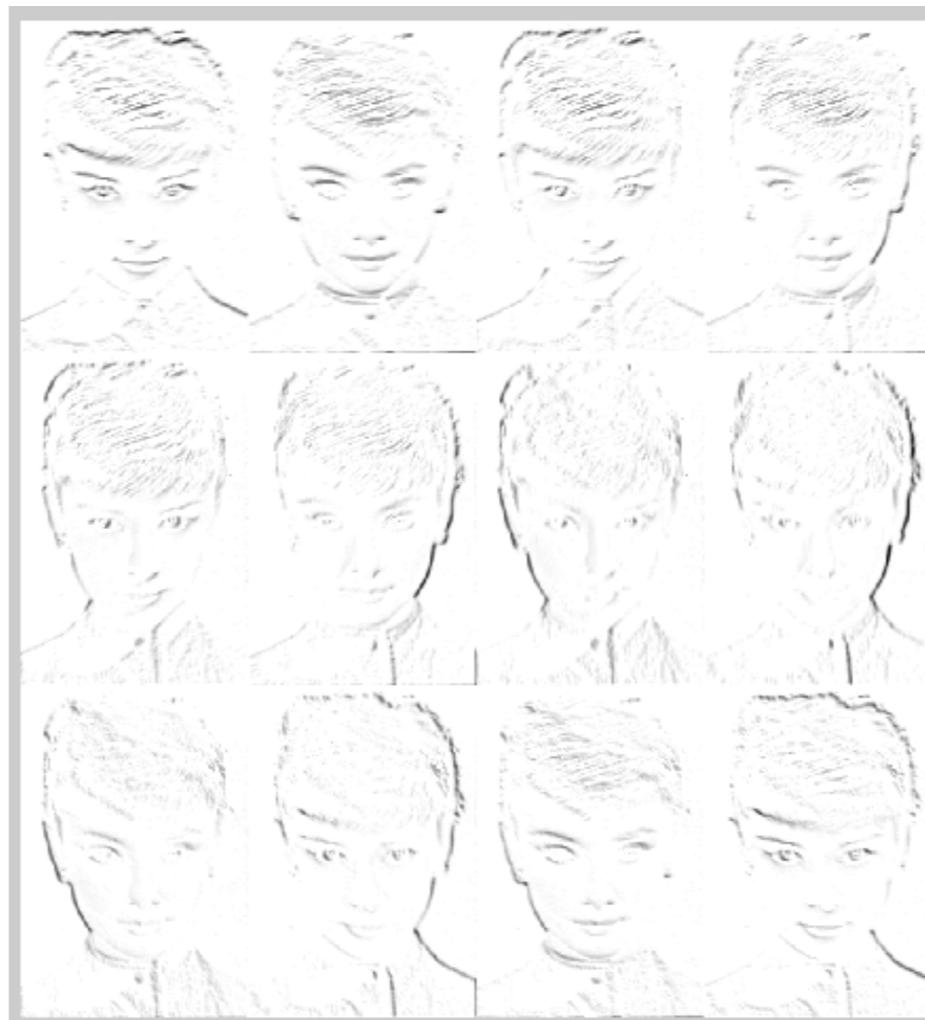
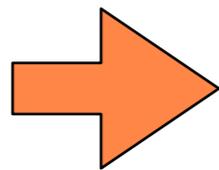


“Rectification”

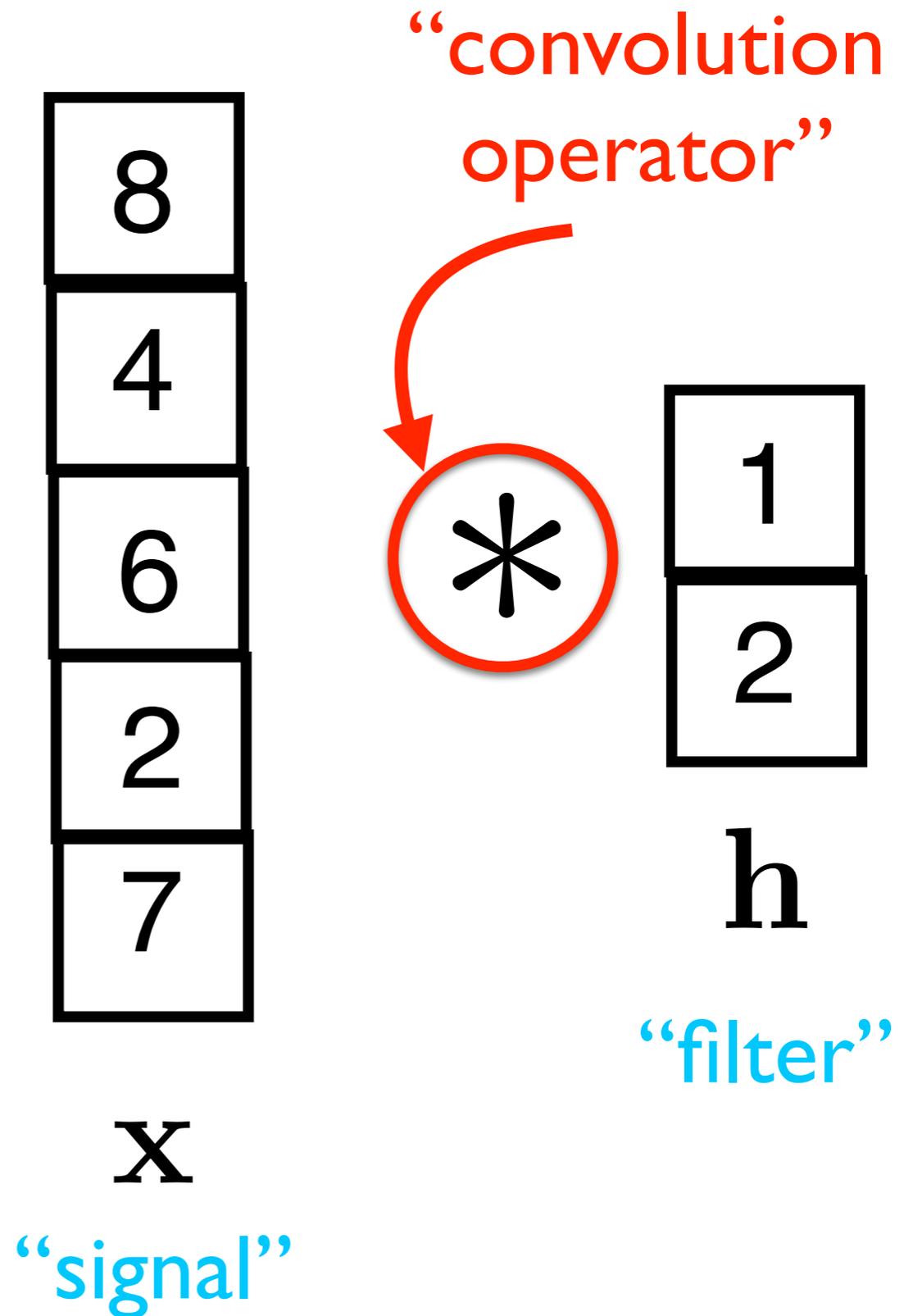
e.g., sigmoid, squared, relu, etc.

# Sparseness and Positiveness

- Blurring only works if the signals being matched are sparse and positive.
- Unfortunately natural images are neither.
- Combination of oriented filter banks and rectification can remedy this problem with little loss in performance.

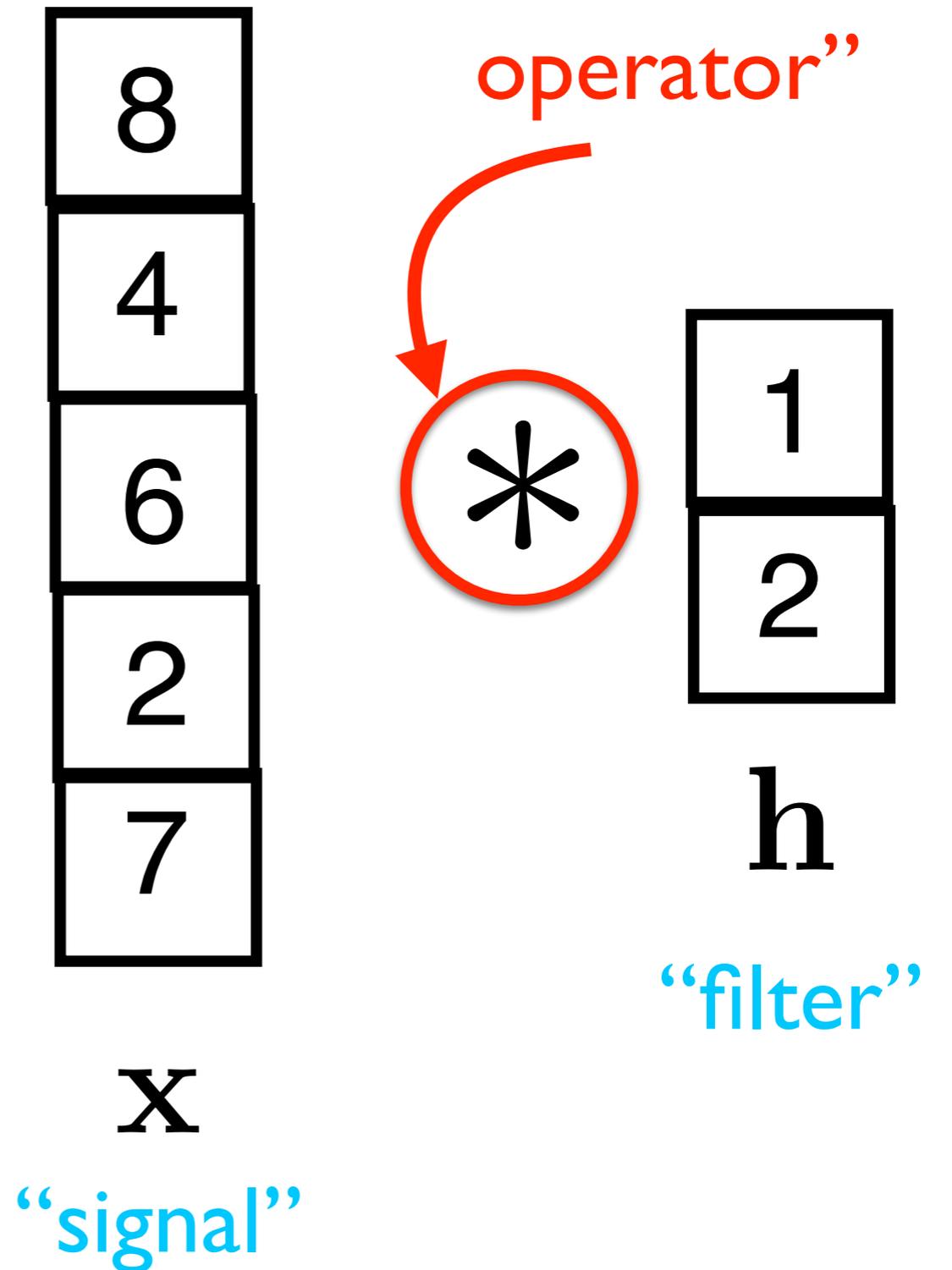


# Reminder: Convolution



# Reminder: Convolution

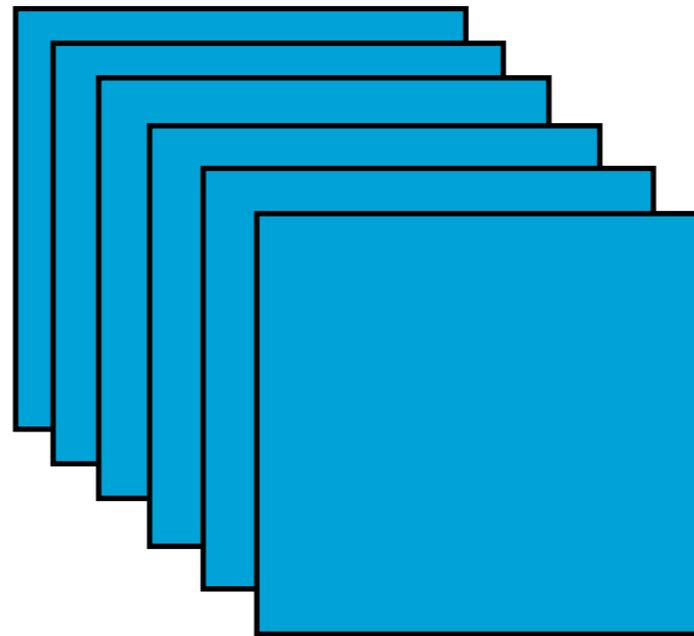
```
>> conv(x,h,'valid')  
ans =  
  
    20  
    14  
    14  
    11
```



# Multi-Channel Convolution



“multi-channel signal”  $\mathbf{x}$



“multi-channel filter”

$\mathbf{h}$

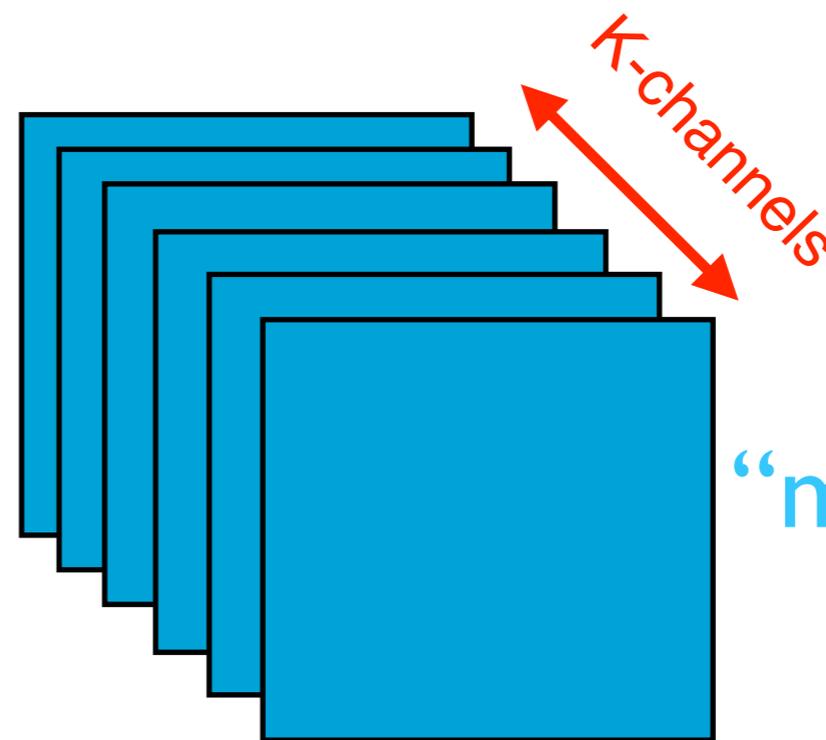


“single-channel response”  $\mathbf{y}$

# Multi-Channel Convolution



“multi-channel signal”  $\mathbf{X}$



“multi-channel filter”

$\mathbf{h}$



“single-channel response”  $\mathbf{y}$

# Multi-Channel Convolution

---

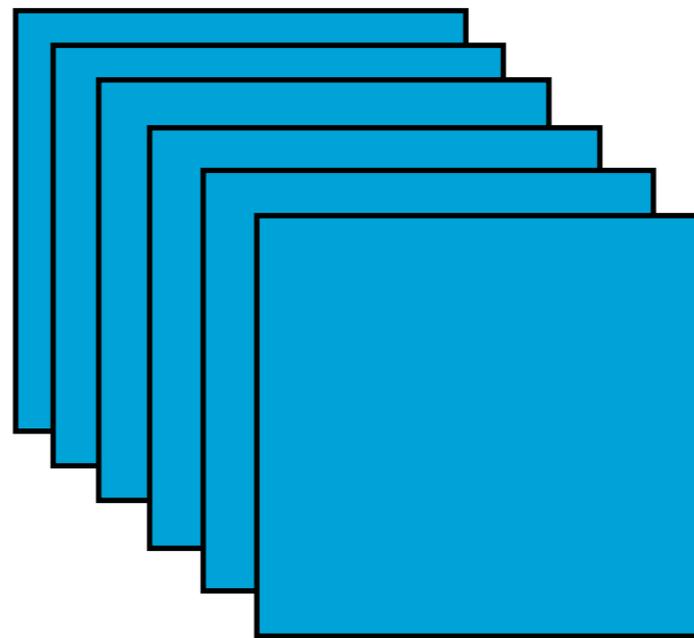
$$\mathbf{y} = \sum_{k=1}^K \mathbf{x}^{(k)} * \mathbf{h}^{(k)}$$

# Multi-Channel Convolution



“multi-channel signal”  $\mathbf{x}$

\*



“multi-channel filter”

$\mathbf{h}$



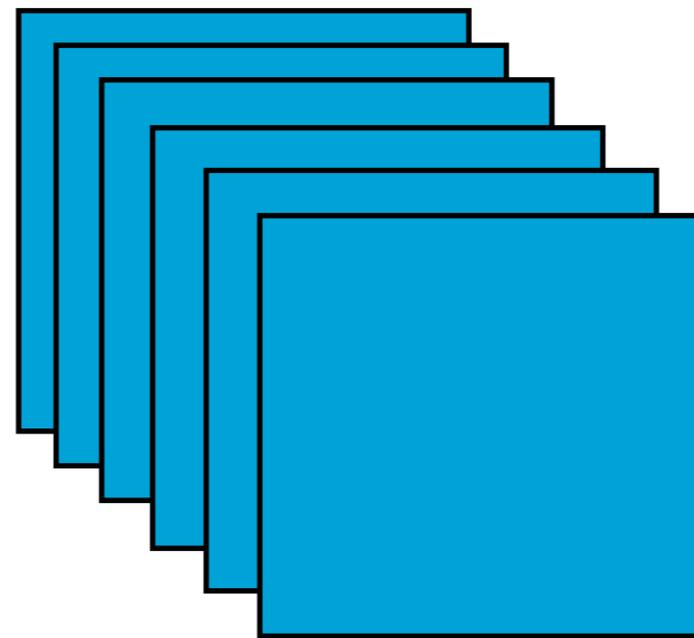
“multi-channel response”  $\mathbf{y}$

# Multi-Channel Convolution



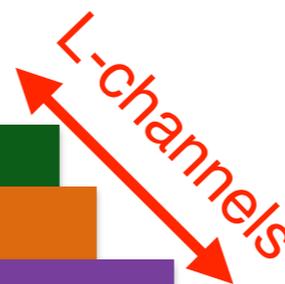
“multi-channel signal”

$\mathbf{X}$



“multi-channel filter”

$\mathbf{h}$



“multi-channel response”

$\mathbf{y}$

# Multi-Channel Convolution

---

$$\mathbf{y}^{(l)} = \sum_{k=1}^K \mathbf{x}^{(k)} * \mathbf{h}^{(k,l)} \quad \text{for } l = 1 : L$$

# CNNs for Object Detection

---



# CNNs for Object Detection

---

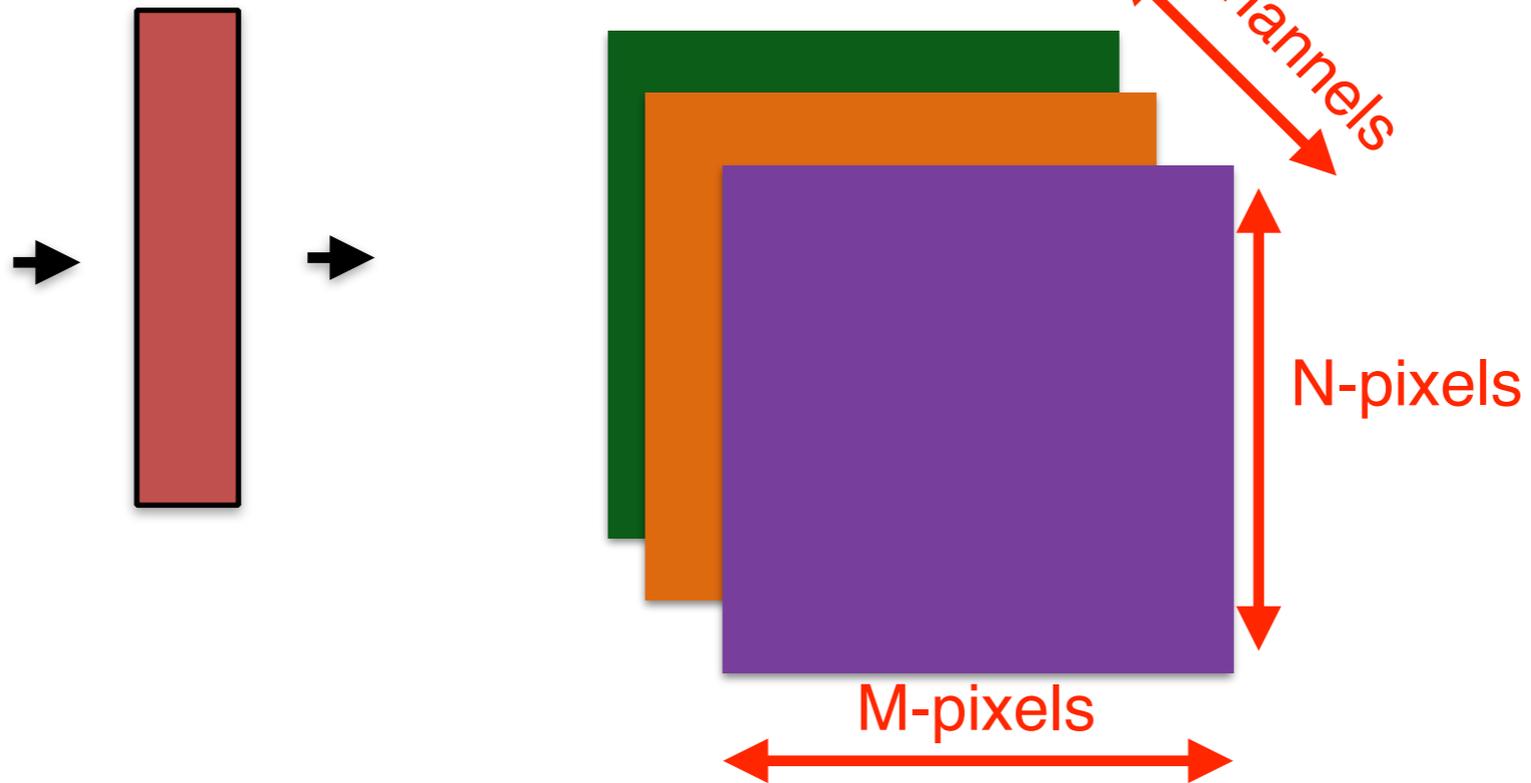
image patch  
3@  
(224x224)



# CNNs for Object Detection

image patch  
3@  
(224x224)

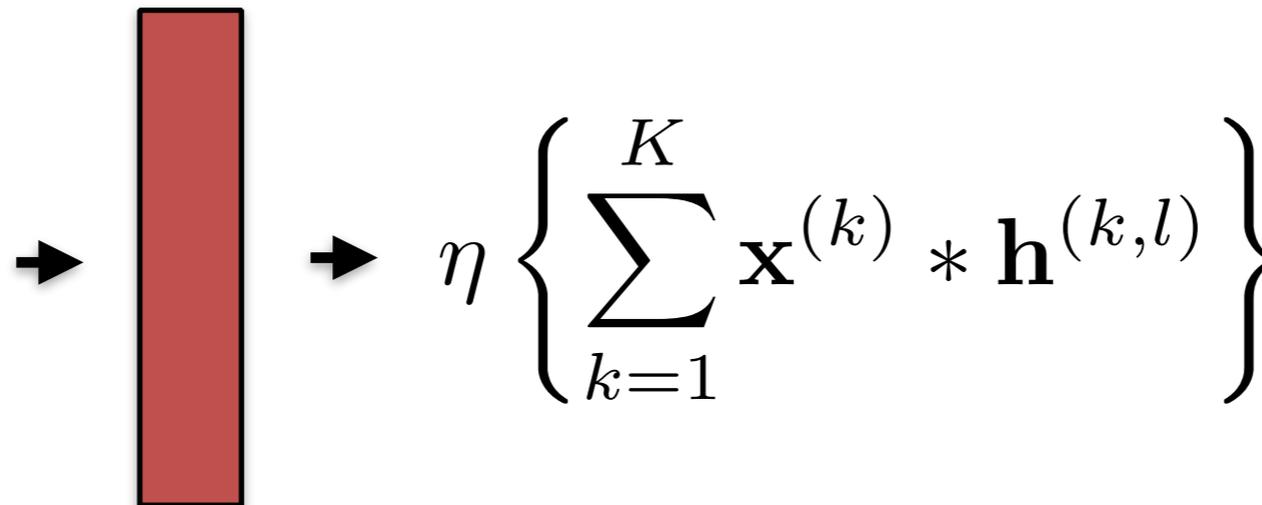
conv  
L@  
(NxM)



# CNNs for Object Detection

image patch  
3@  
(224x224)

conv  
L@  
(NxM)

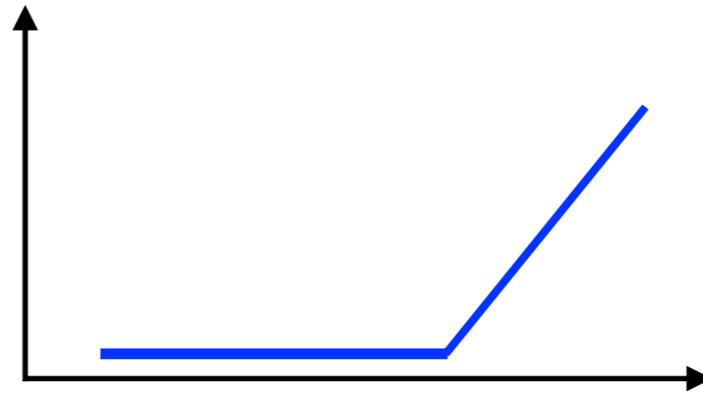


$\eta\{\}$  → non-linear function (relu, max pooling)

# ReLU - Sparse and Positive

- Rectified Linear Unit

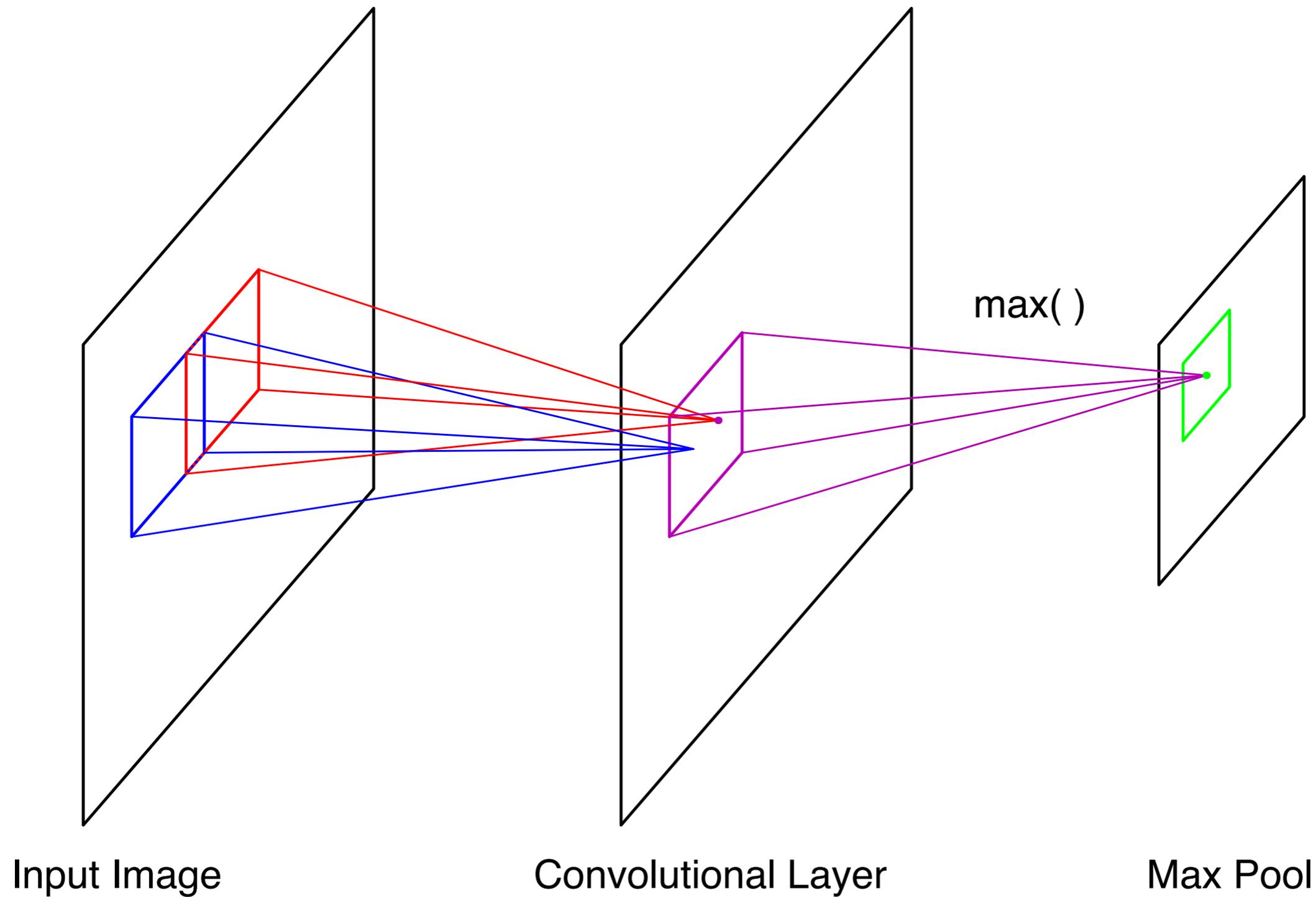
$$\text{relu}\{x\} = \max(0, x)$$



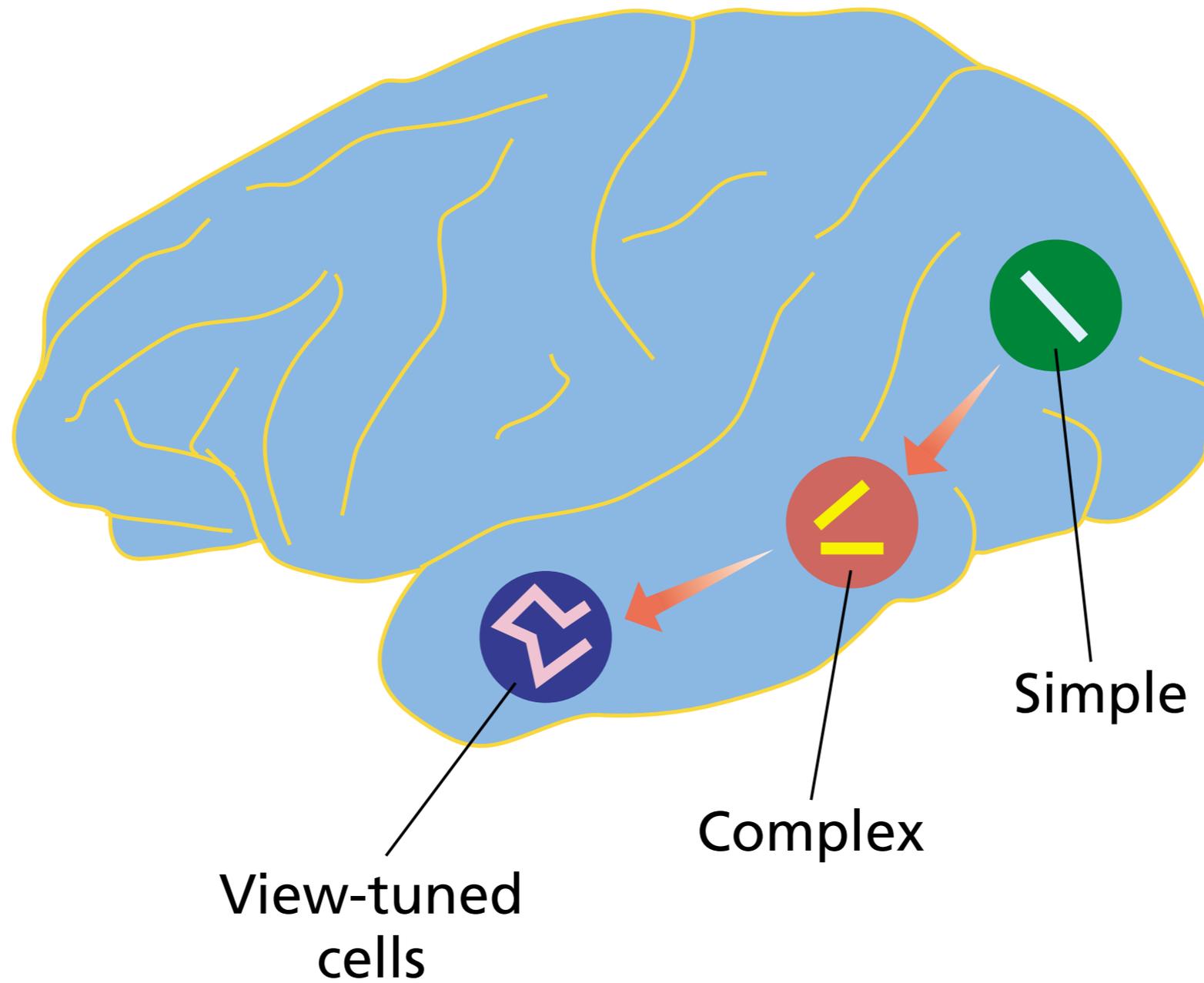
- Connection to LASSO and sparsity??

$$\|y - \mathbf{Ax}\|_2^2 + \frac{\lambda}{2} \|\mathbf{x}\|_1$$

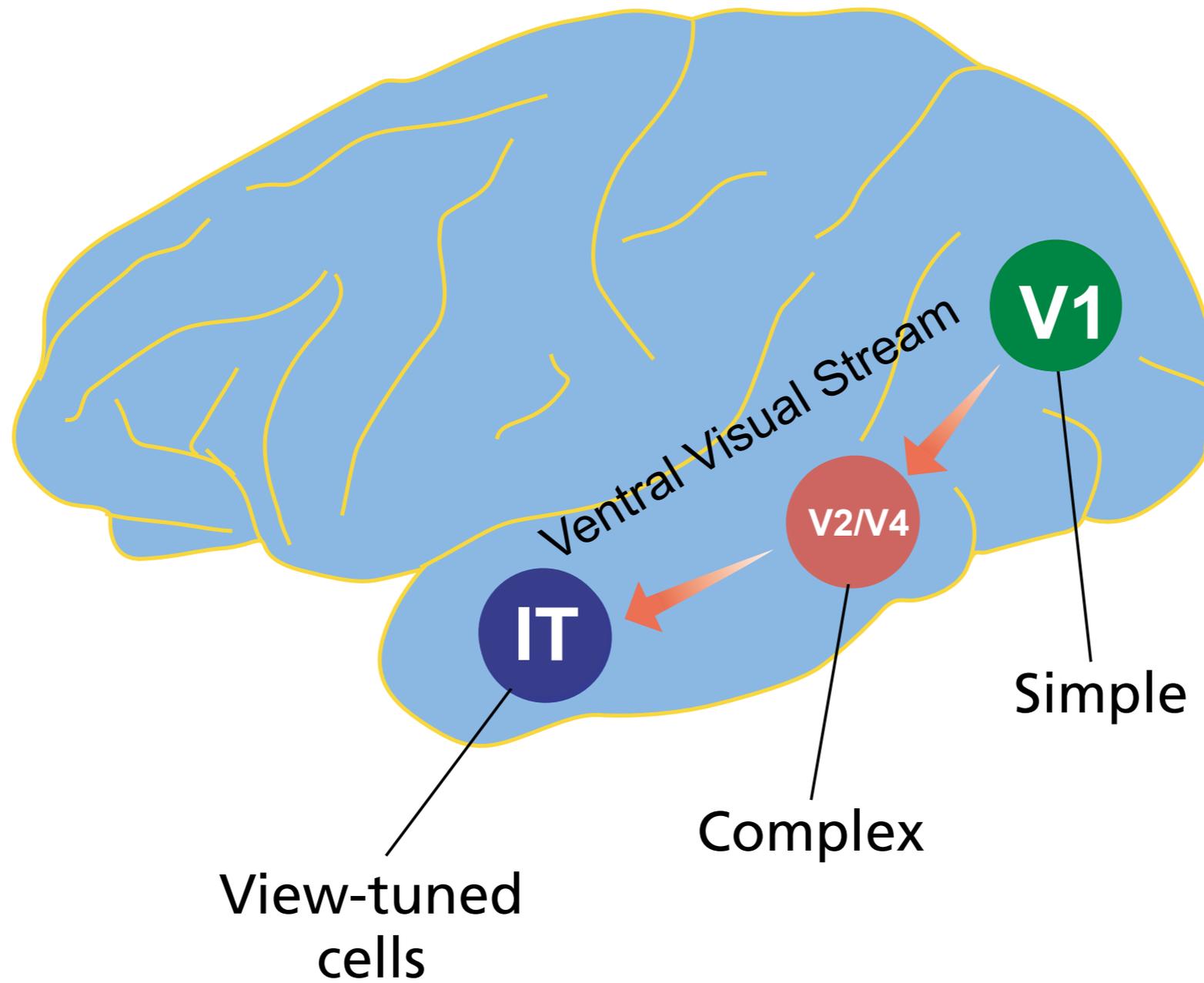
# Max Pooling - Down Sampling



# Hierarchical Learning

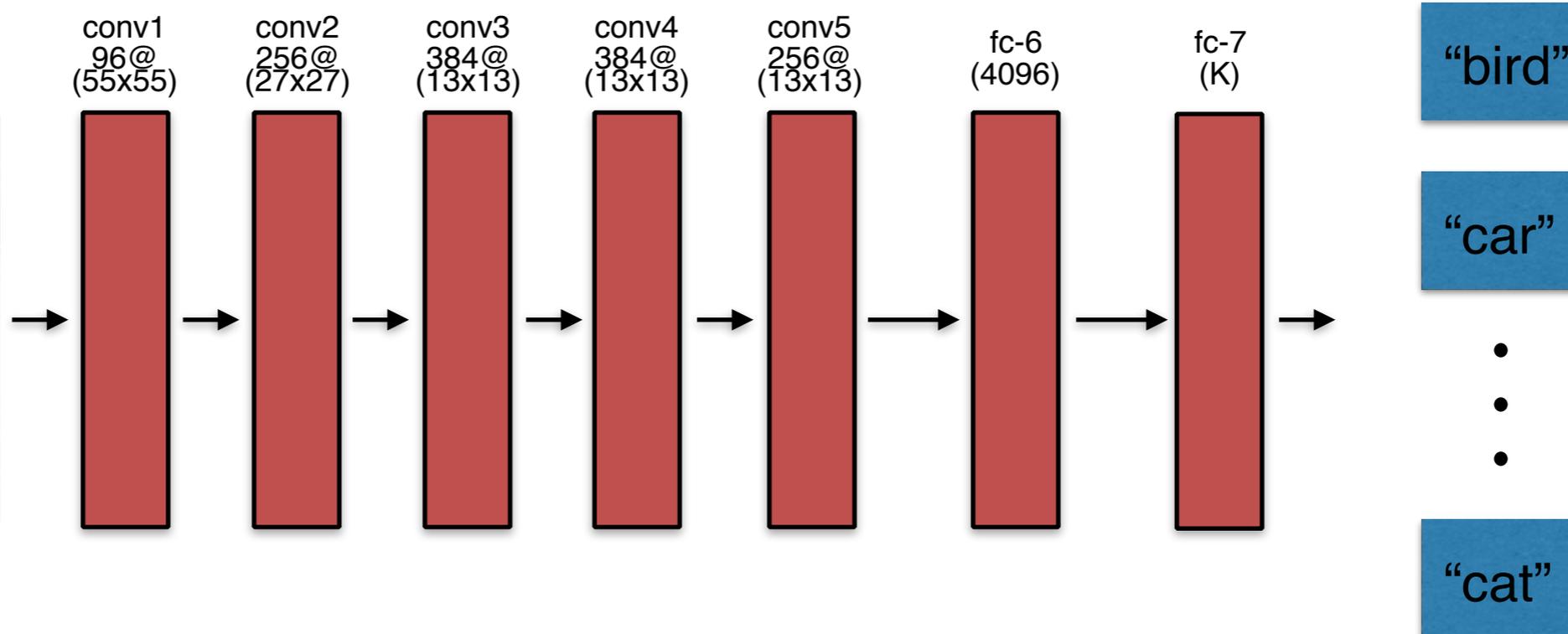


# Hierarchical Learning



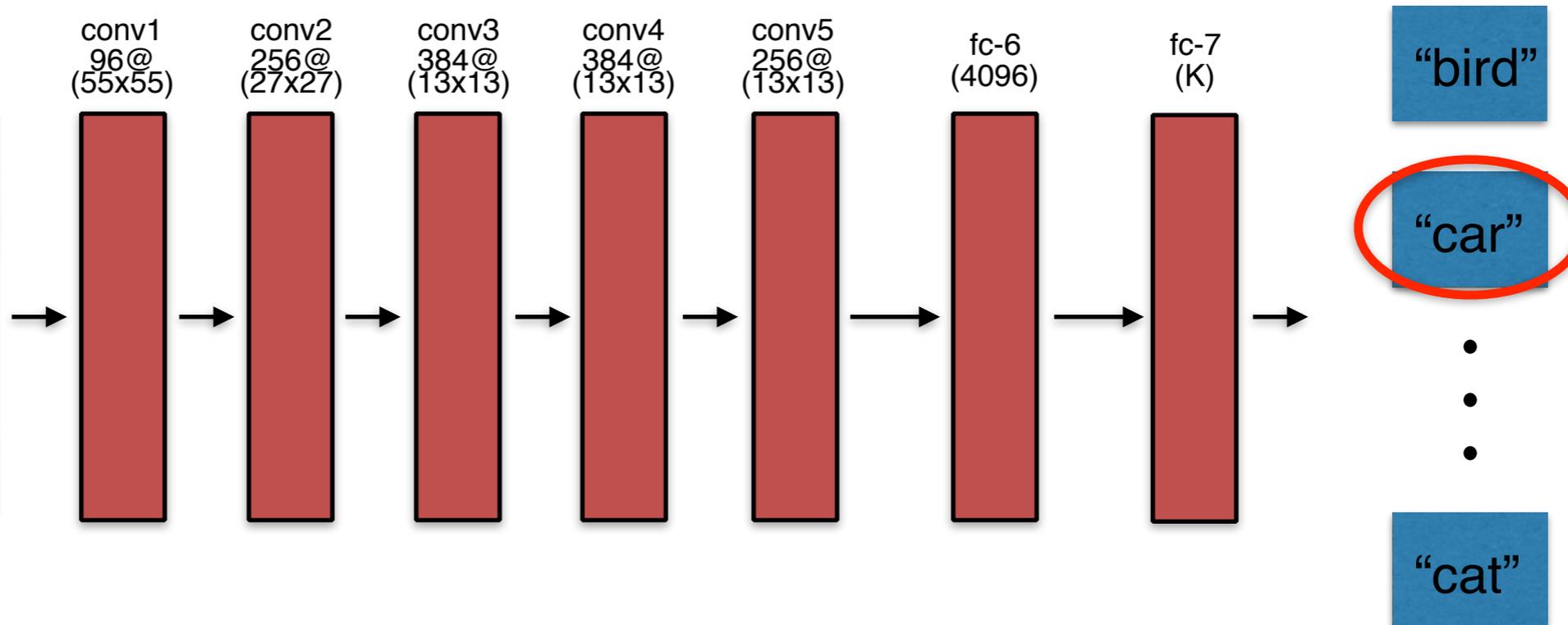
# Current State of the Art

image patch  
3@  
(227x227)



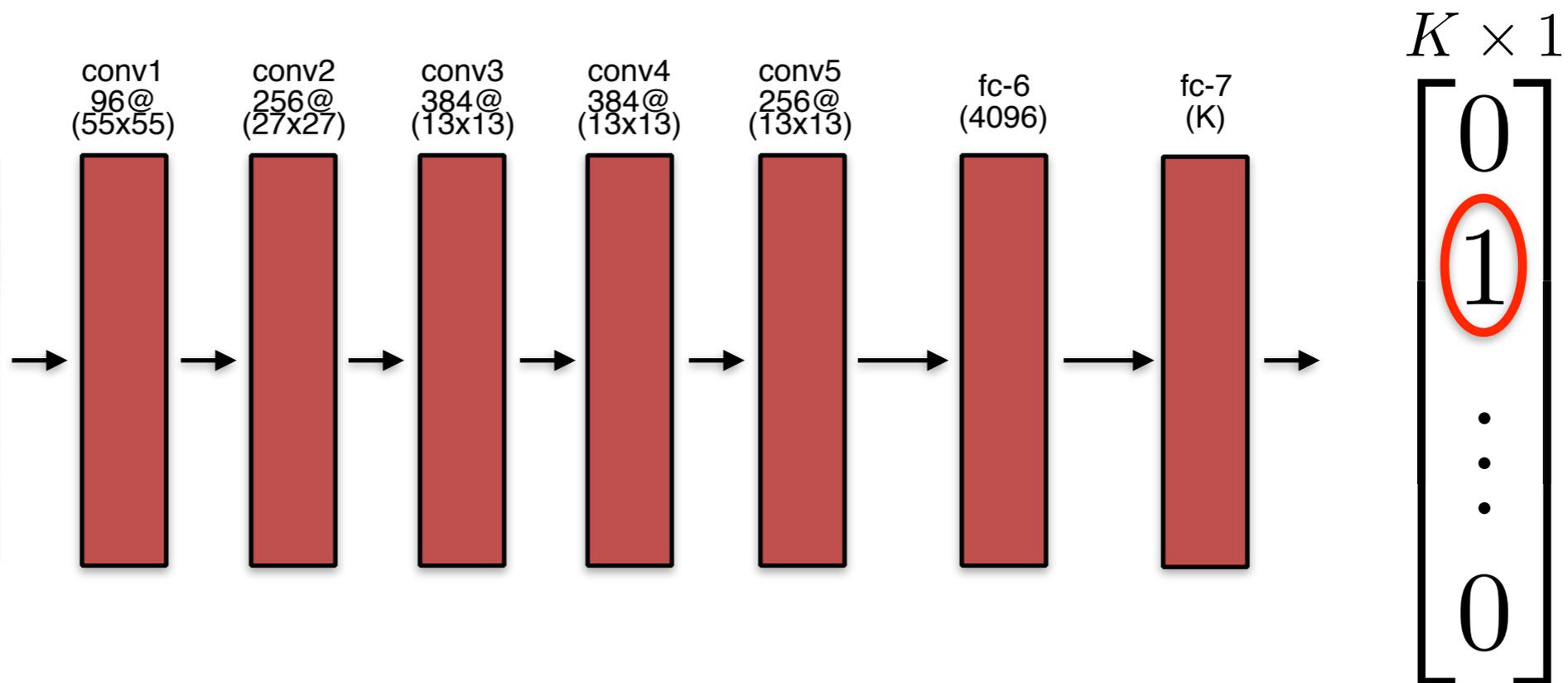
# Current State of the Art

image patch  
3@  
(227x227)



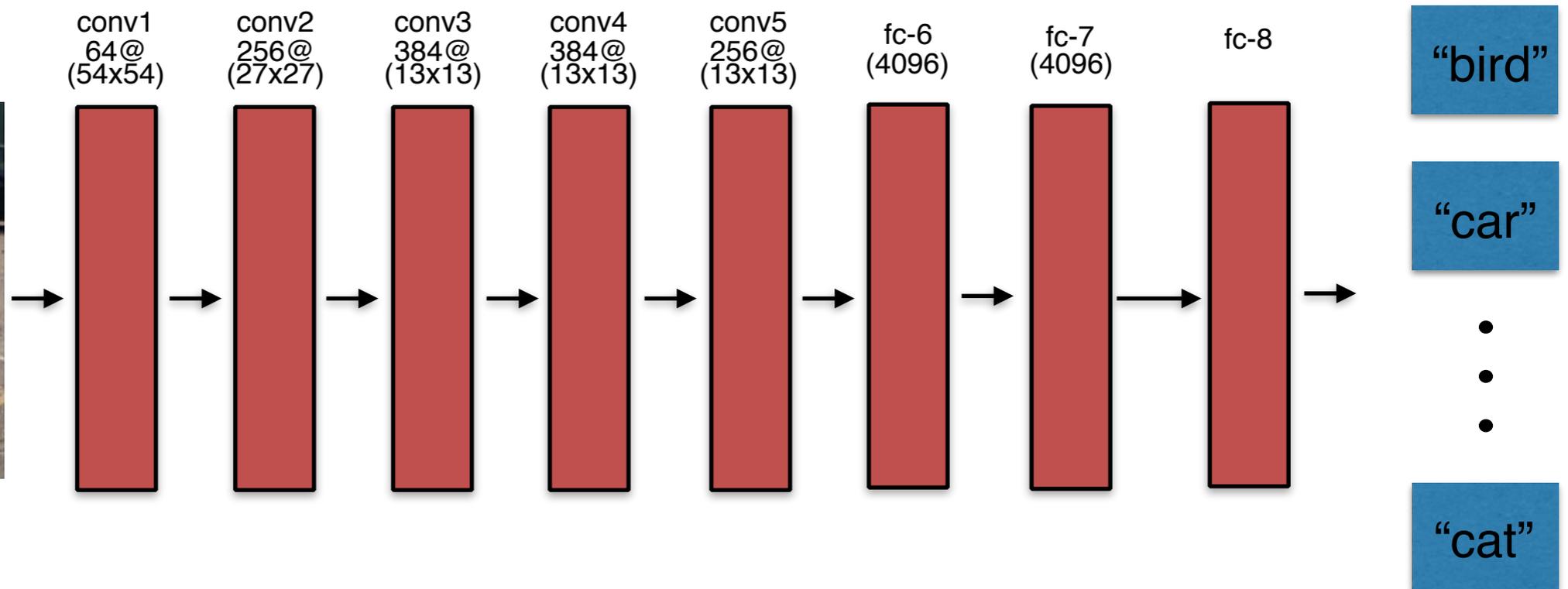
# Current State of the Art

image patch  
3@  
(227x227)



# Current State of the Art - Pose Selection

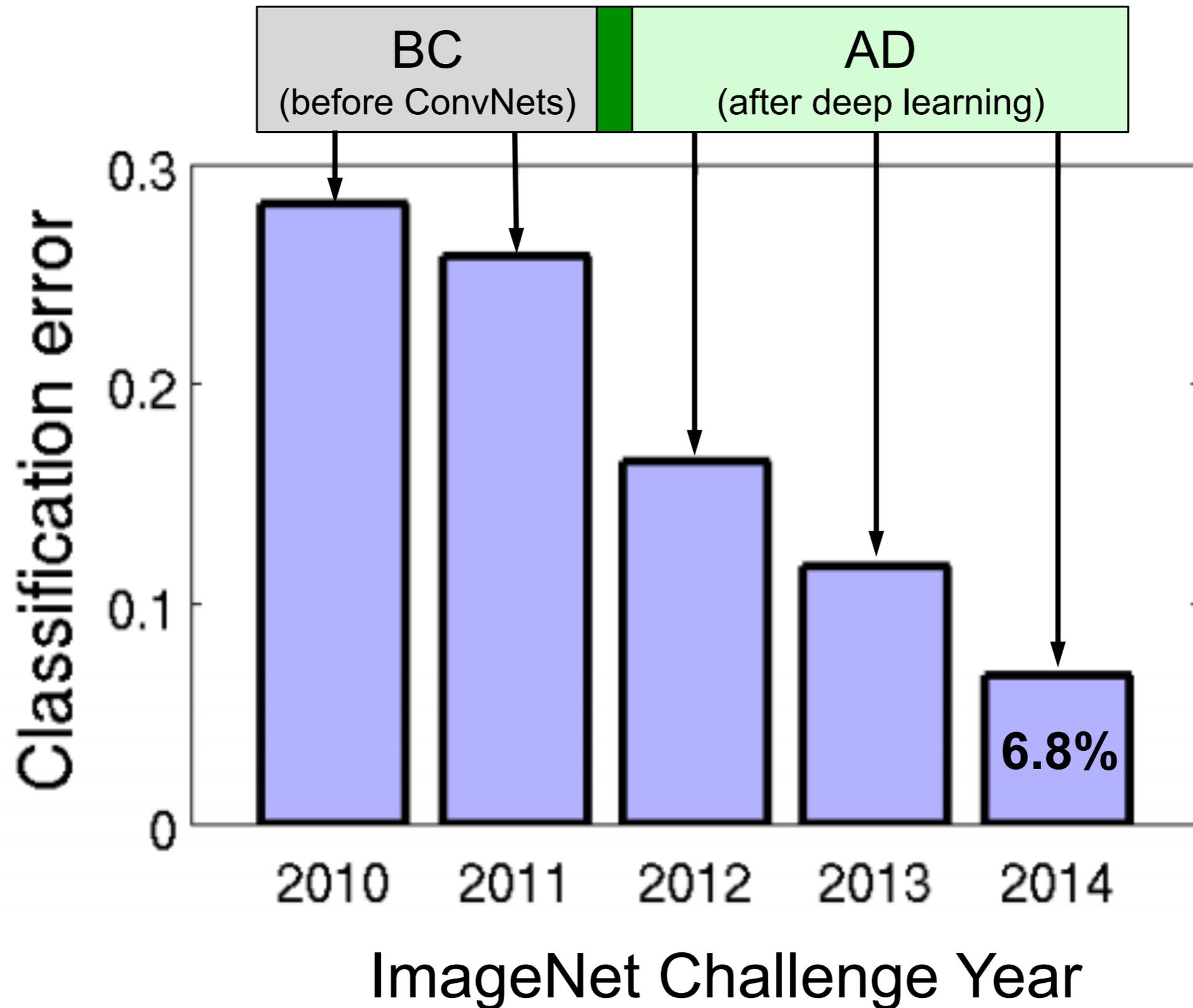
image patch  
3@  
(224x224)



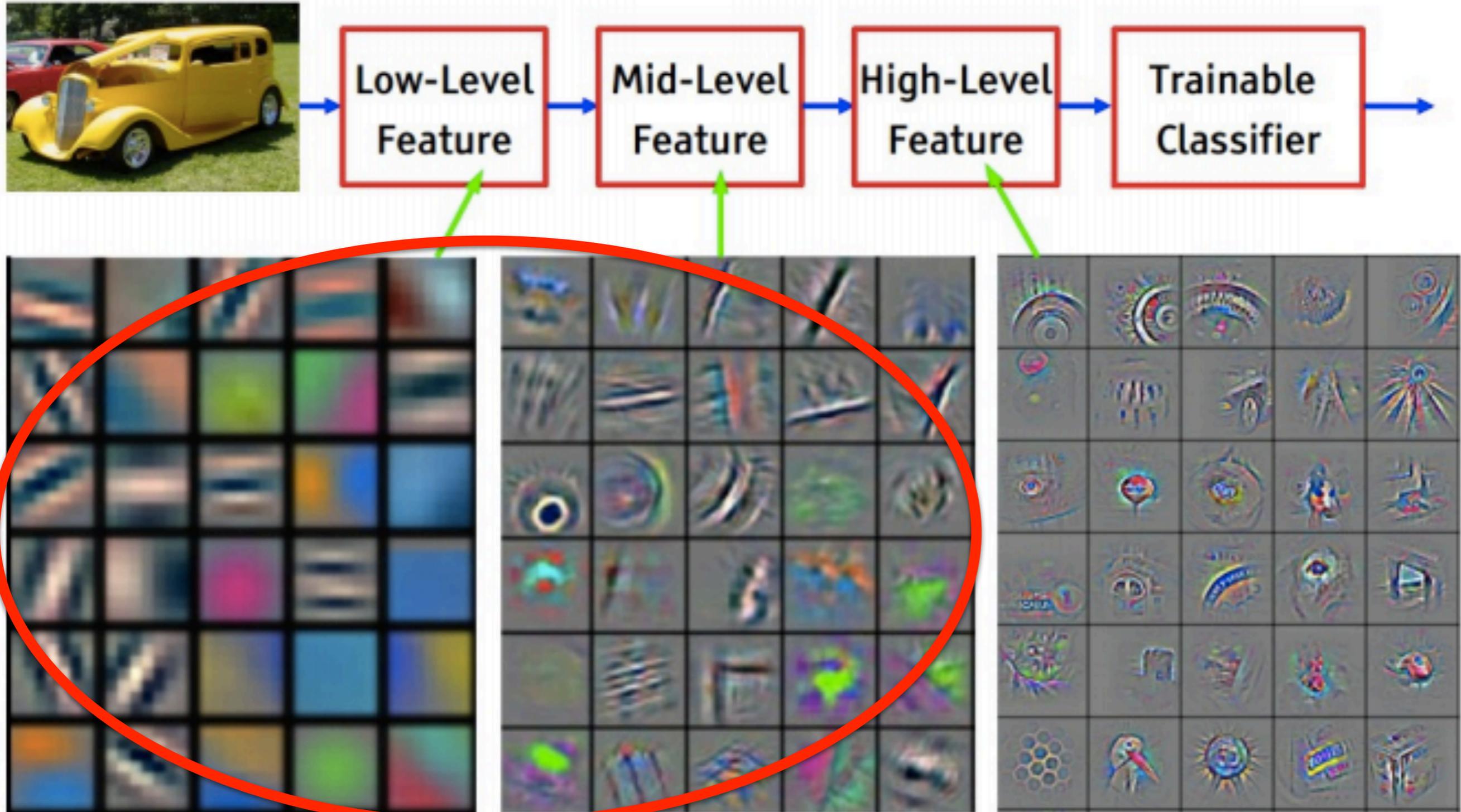
A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In NIPS, 2012.

K. Chatfield, V. Lempitsky, A. Vedaldi and A. Zisserman. "Return of the Devil in the Details: Delving Deep into Convolutional Networks." In BMVC, 2014.

# Impact on Object Recognition

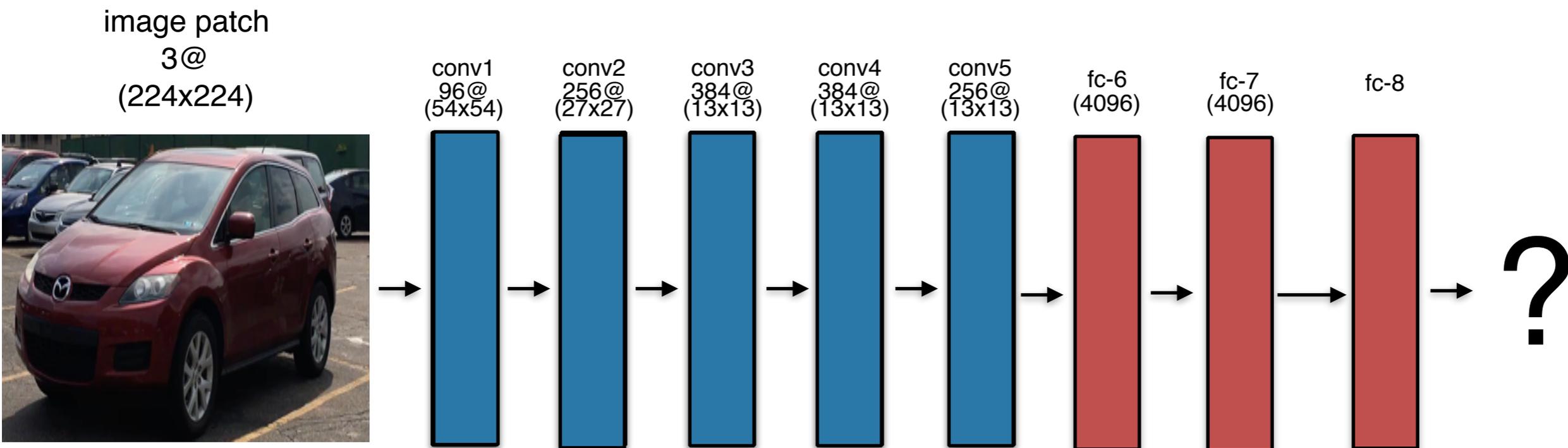


# Visualizing CNNs



Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

# CNNs as Feature Extraction



parameters to learn



pre-learned parameters (VGG)

A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In NIPS, 2012.

K. Chatfield, V. Lempitsky, A. Vedaldi and A. Zisserman. "Return of the Devil in the Details: Delving Deep into Convolutional Networks." In BMVC, 2014.

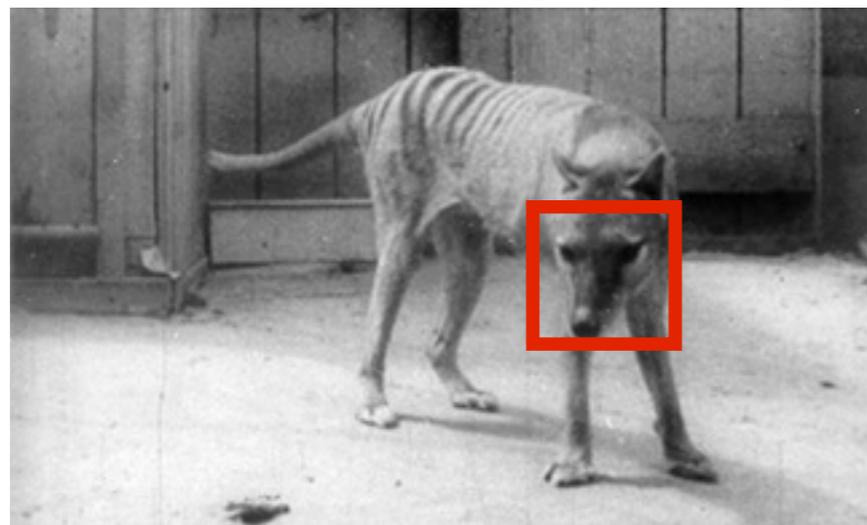
# Today

---

- Deep Features
- Deep Tracking
- Deep Flow

# Drawback to Conventional Methods

- Most methods for object tracking employ “online” learning.
- Online methods are expensive, have to make simplifying assumptions (e.g. circulant Toeplitz) to make things efficient.



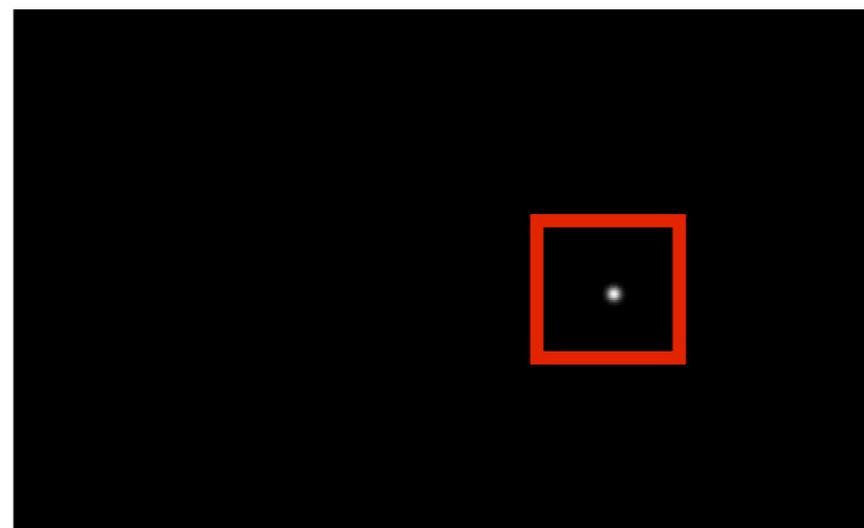
“known signal”  $\mathbf{x}$

\*



$\mathbf{h}$

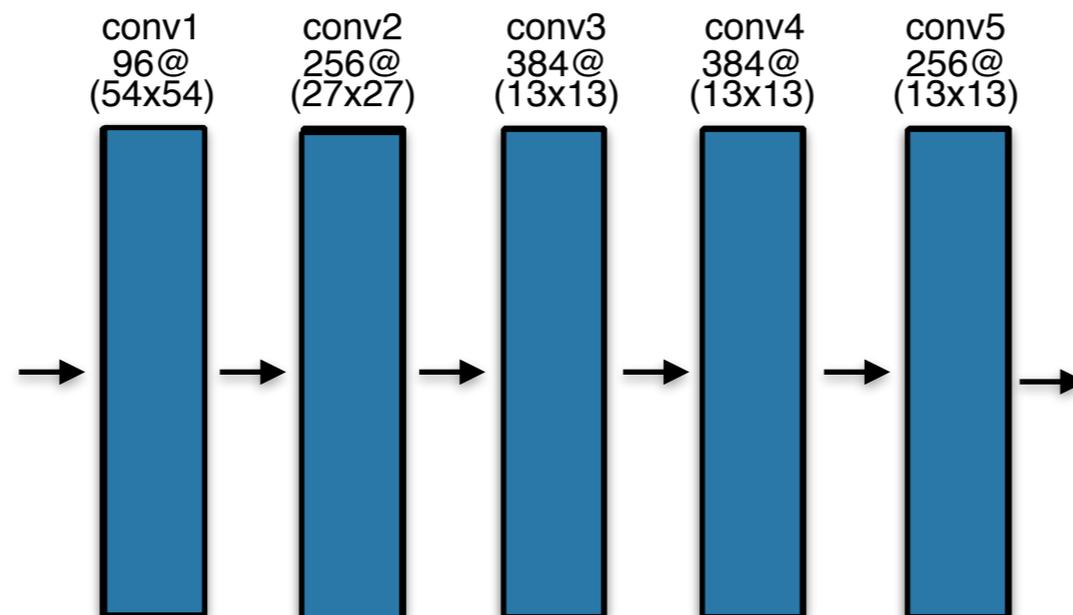
“unknown  
filter”



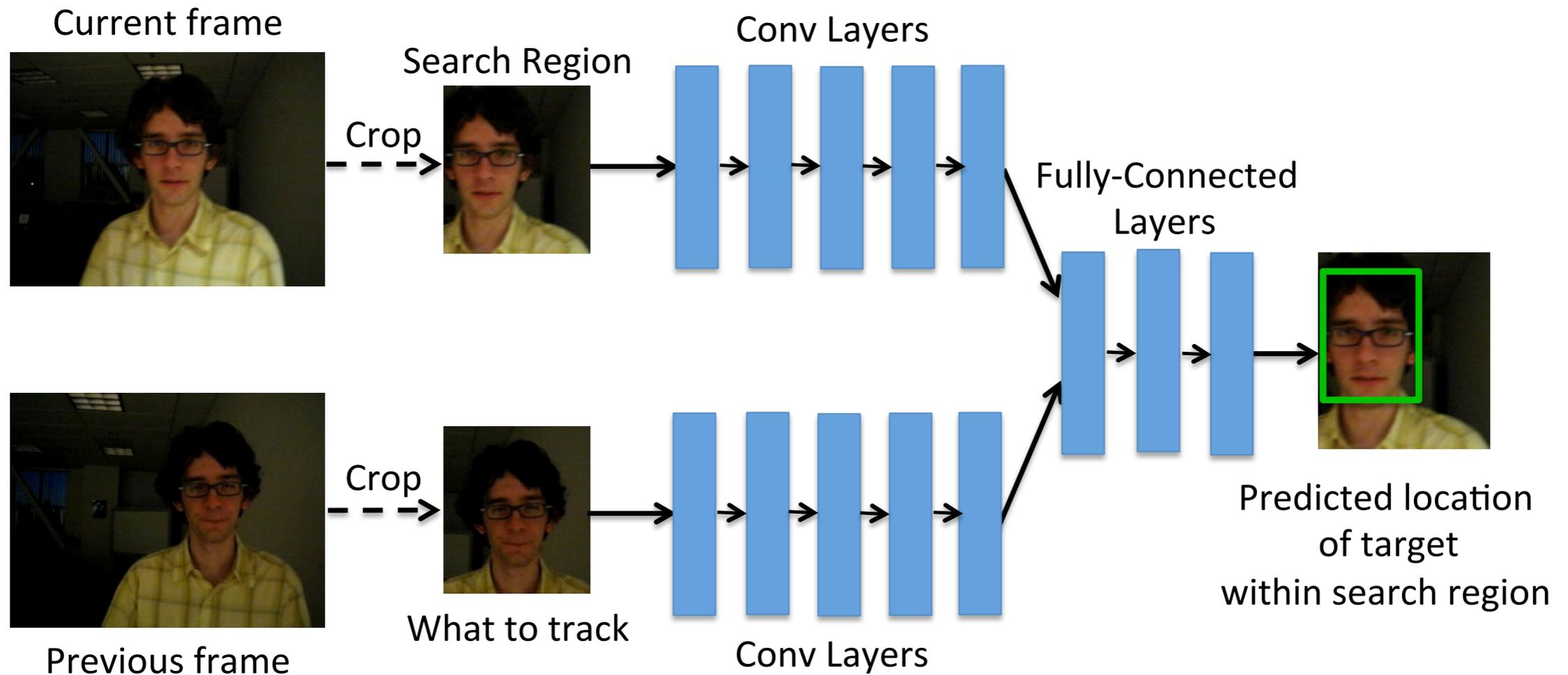
“known response”  $\mathbf{y}$

# Deep Tracking Methods

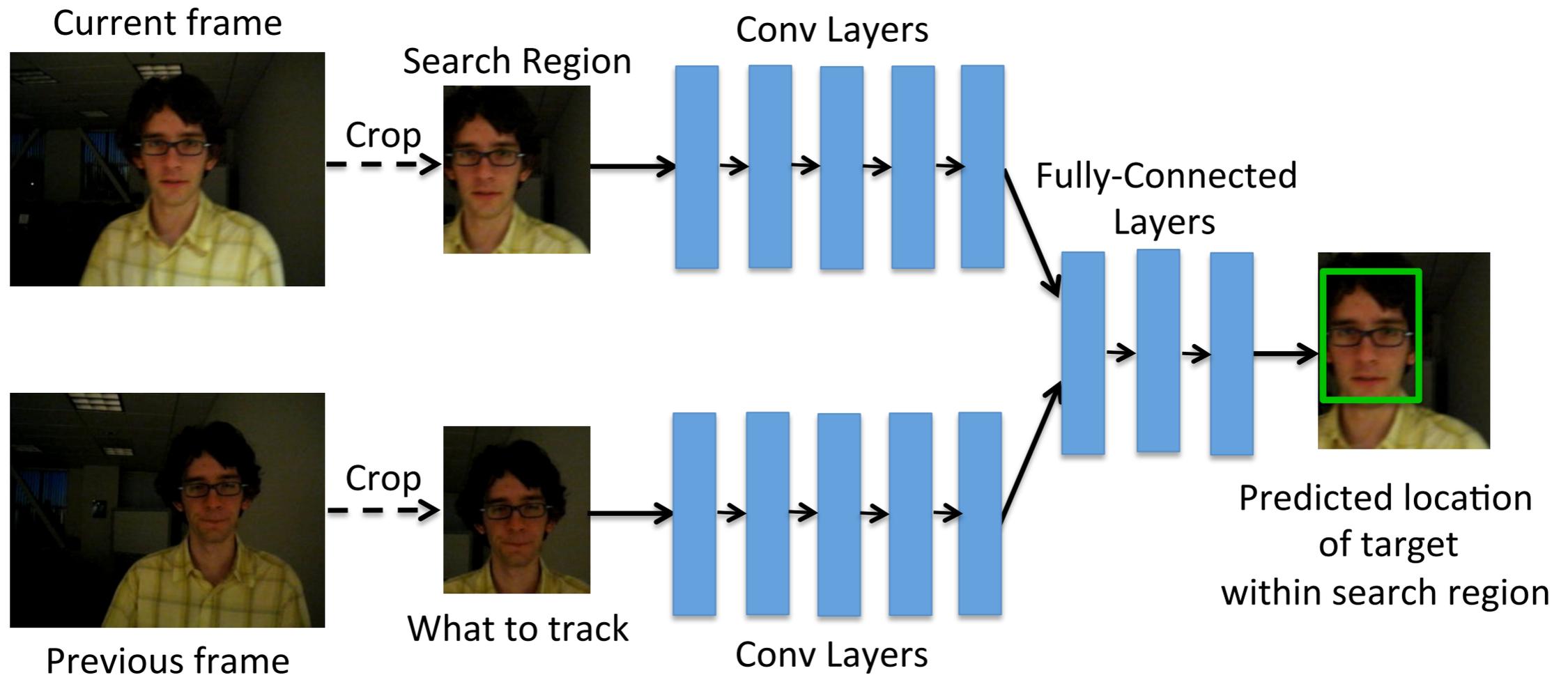
- Recently, there have been works that have tried to explore the employment of tracking using deep learning features.
- As efficiency is key, strategy is to learn from a large ensemble of labeled offline videos.
- Of particular interest are two papers,
  1. D. Held, S. Thrun, and S. Savarese “Learning to Track at 100 FPS with Deep Regression Networks”, ECCV 2016.
  2. L. Bertinetto J. Valmadre J. F. Henriques, A. Vedaldi, P. H. S. Torr “Fully-Convolutional Siamese Networks for Object Tracking”, ArXiv 2016.



# Deep Regression Networks

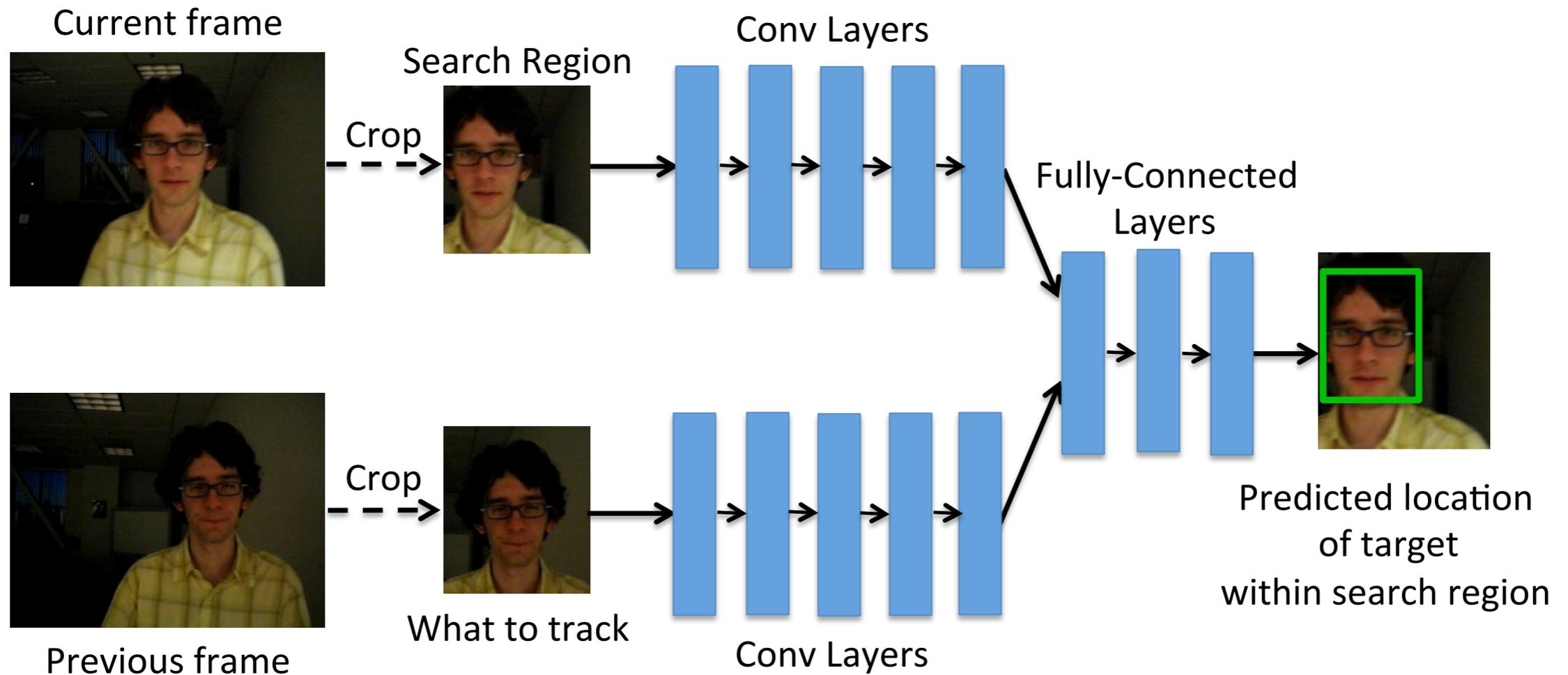


# Deep Regression Networks



**What does this method remind you of?**

# Deep Regression Networks



**What does this method remind you of?**

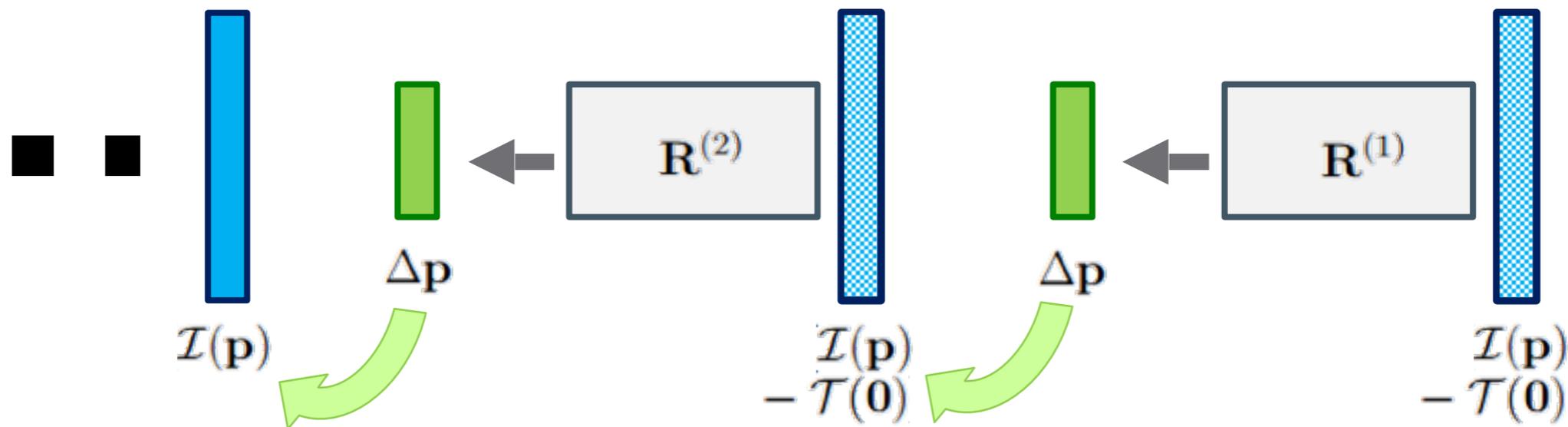
**Why is it fast?**

# Supervised Descent Method (SDM)

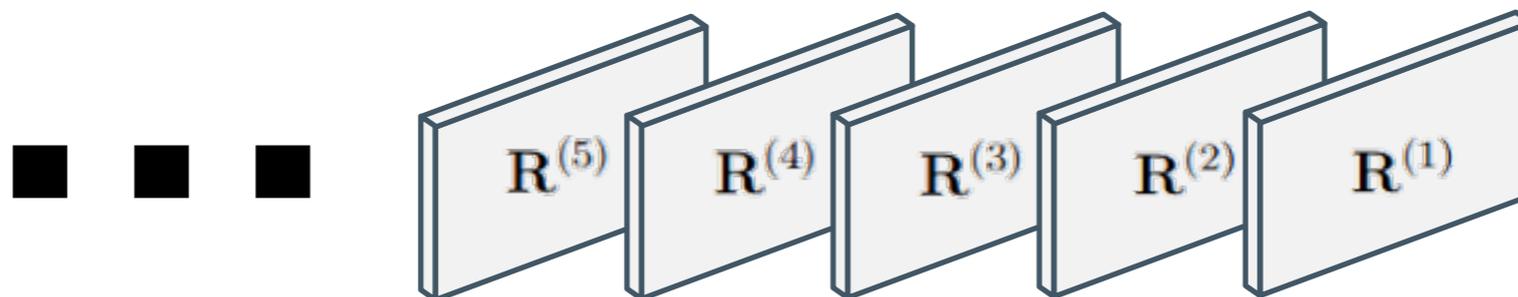
- SDM assumes a linear relationship between appearance and geometry:

$$\Delta \mathbf{p} = \mathbf{R}[\mathcal{I}(\mathbf{p}) - \mathcal{T}(\mathbf{0})]$$

- Iteratively updates until convergence



- However,  $\mathbf{R}^{(k)}$  is iteration specific.

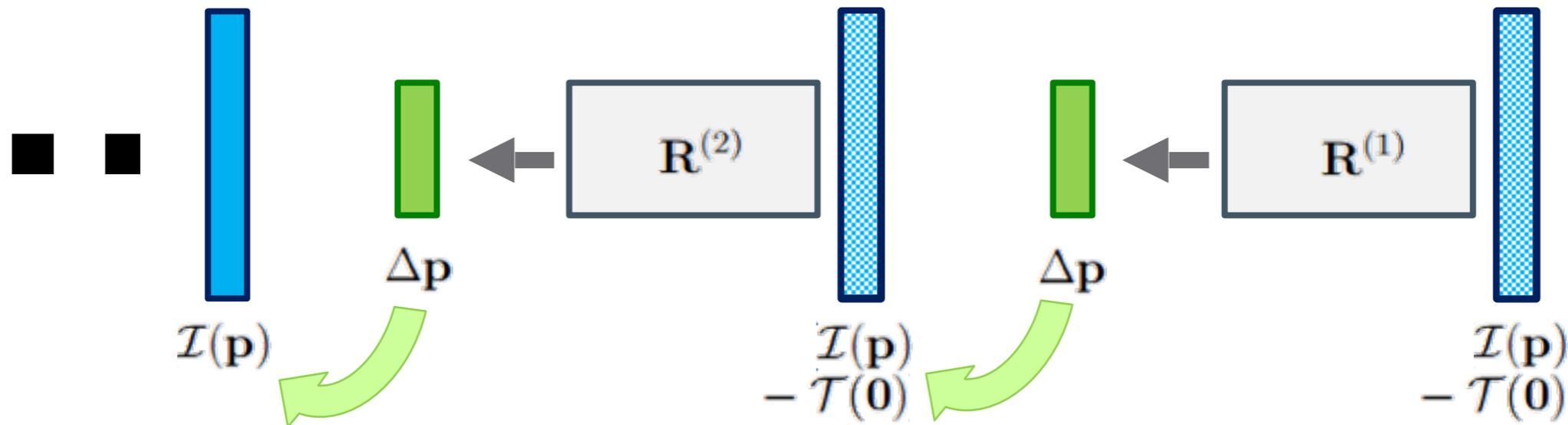


# Supervised Descent Method (SDM)

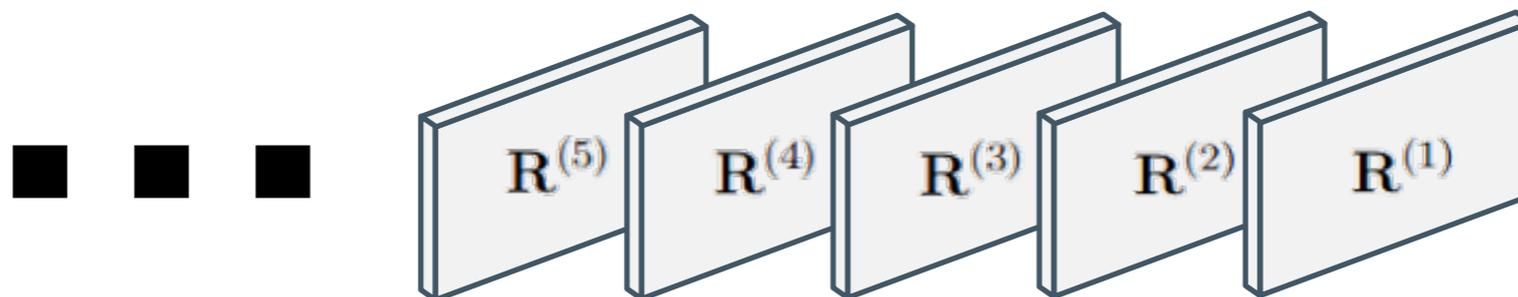
- SDM assumes a linear relationship between appearance and geometry:

$$\Delta \mathbf{p} = \mathbf{R}[\mathcal{I}(\mathbf{p}) - \mathcal{T}(\mathbf{0})]$$

- Iteratively updates until convergence



- However,  $\mathbf{R}^{(k)}$  is iteration specific.



**What is a potential issue here?**

# Deep Regression Networks

Previous video frame centered on object



Current video frame, shifted, with ground-truth bounding box

# Deep Regression Networks

Image  
centered on  
object



Shifted image  
with ground-truth  
bounding box



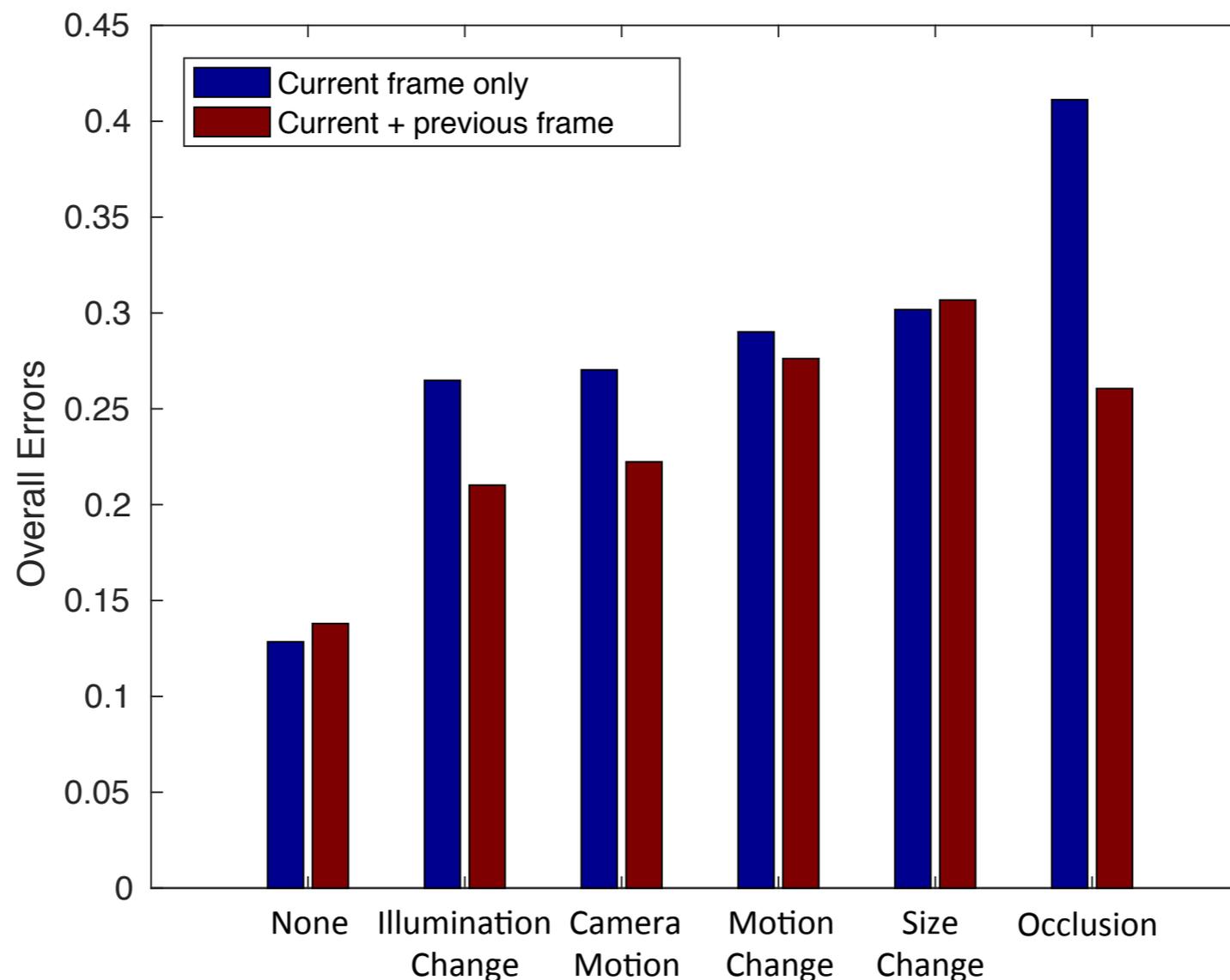
# How does it work?

---

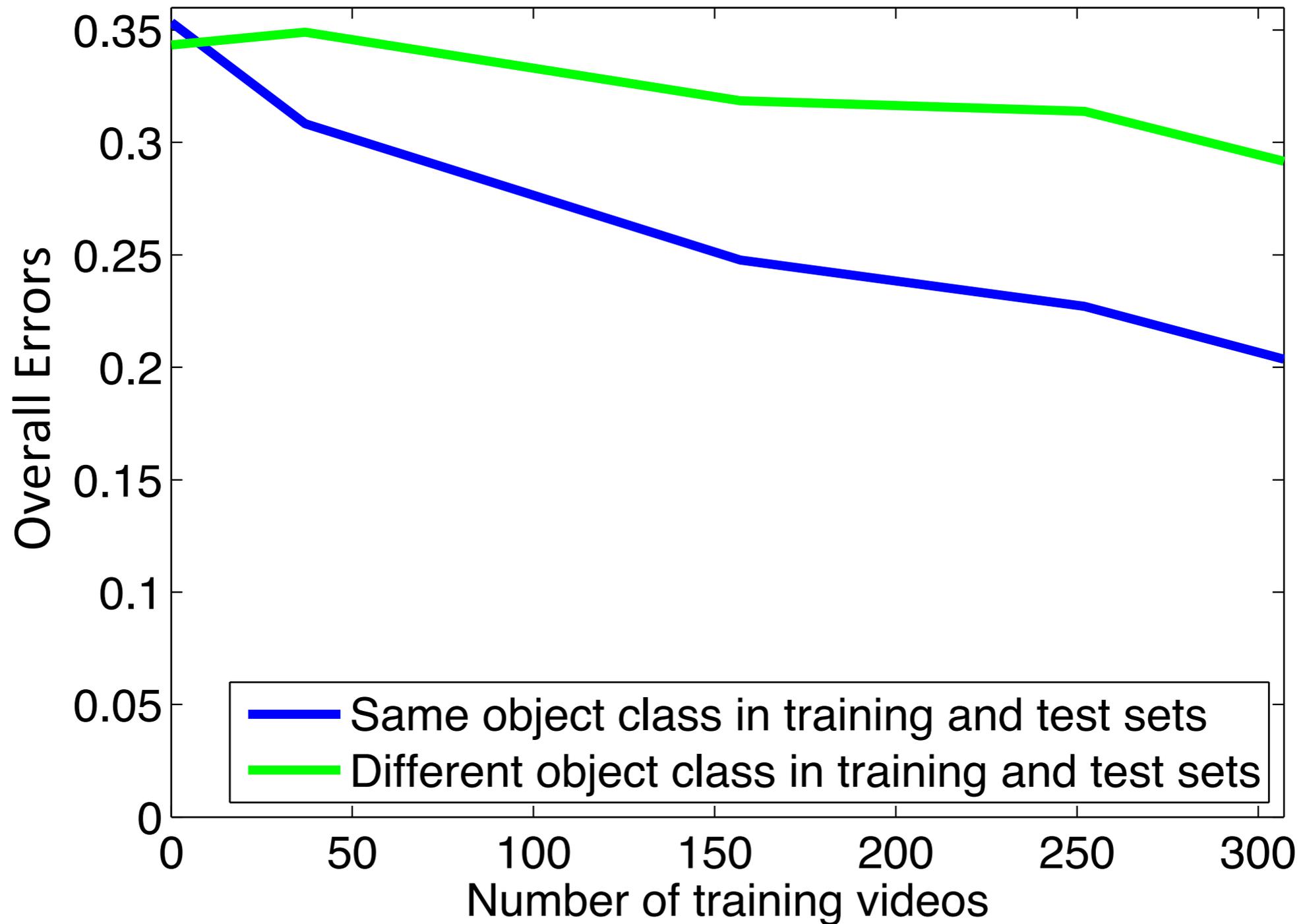
- Two hypotheses,
  1. The network compares the previous frame to the current frame to find the target object in the current frame.
  2. The network acts as a local generic “object detector” and simply locates the nearest “object.”

# How does it work?

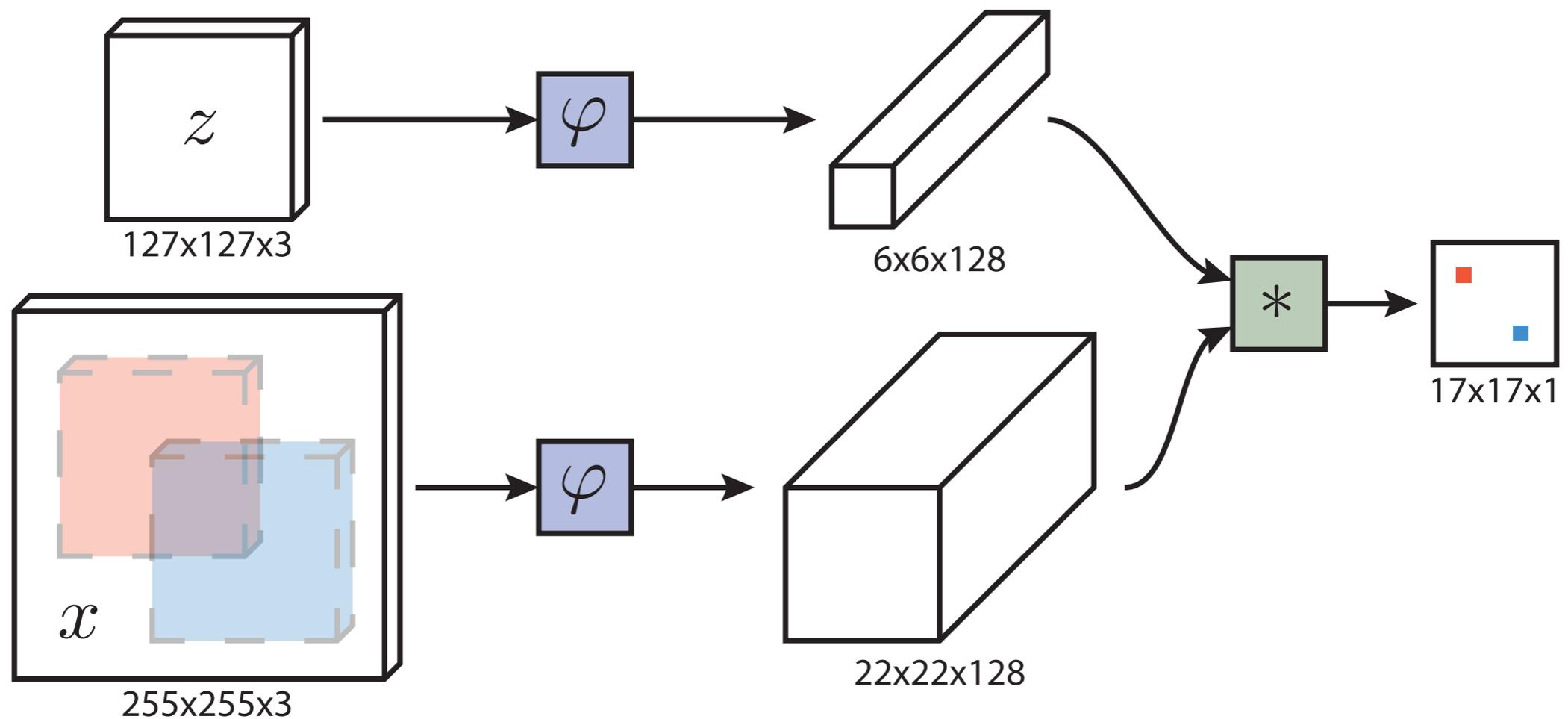
- Two hypotheses,
  1. The network compares the previous frame to the current frame to find the target object in the current frame.
  2. The network acts as a local generic “object detector” and simply locates the nearest “object.”



# Generality vs. Specificity

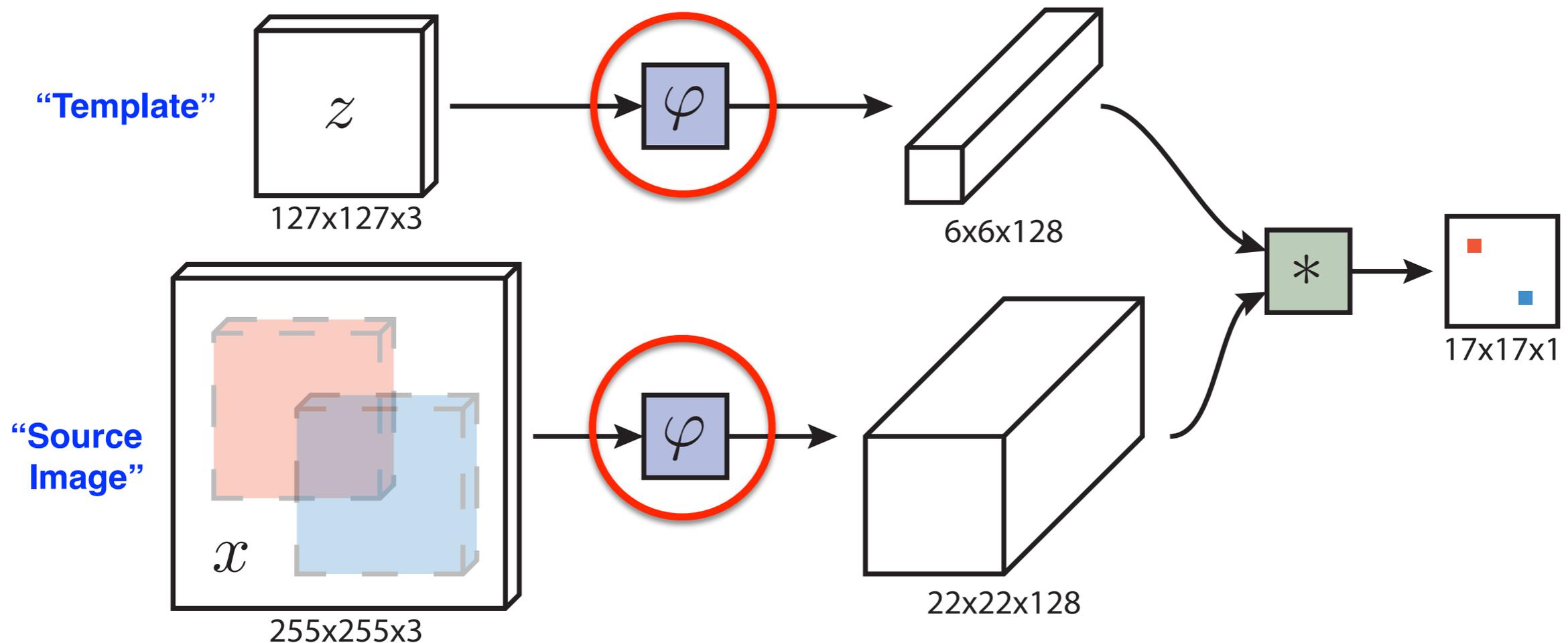


# Fully Convolutional Siamese Networks

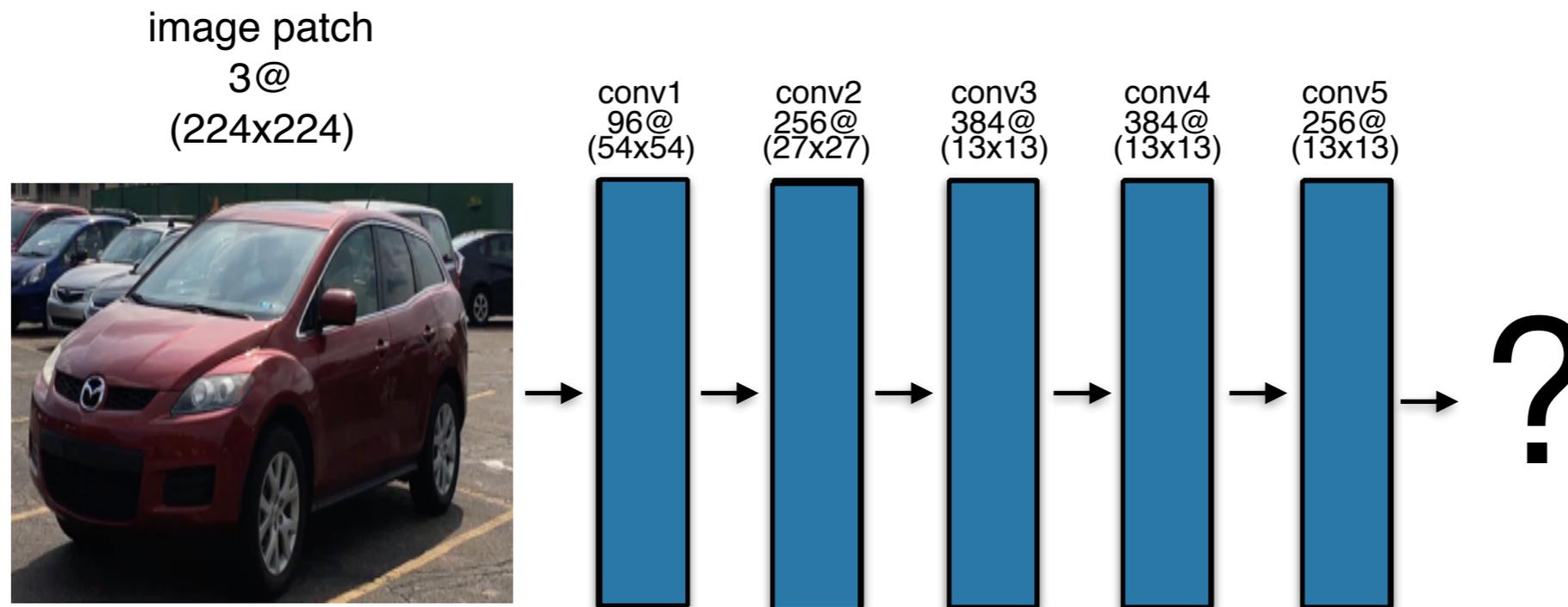


# Fully Convolutional Siamese Networks

“Siamese” as they apply an identical transformation to both inputs

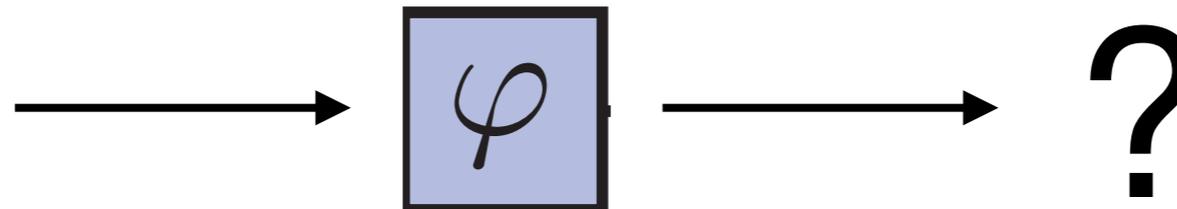


# Fully Convolutional



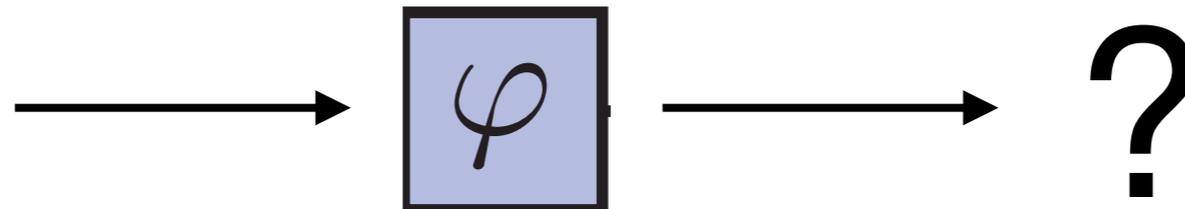
# Fully Convolutional

image patch  
3@  
(224x224)



# Fully Convolutional

image patch  
3@  
(224x224)



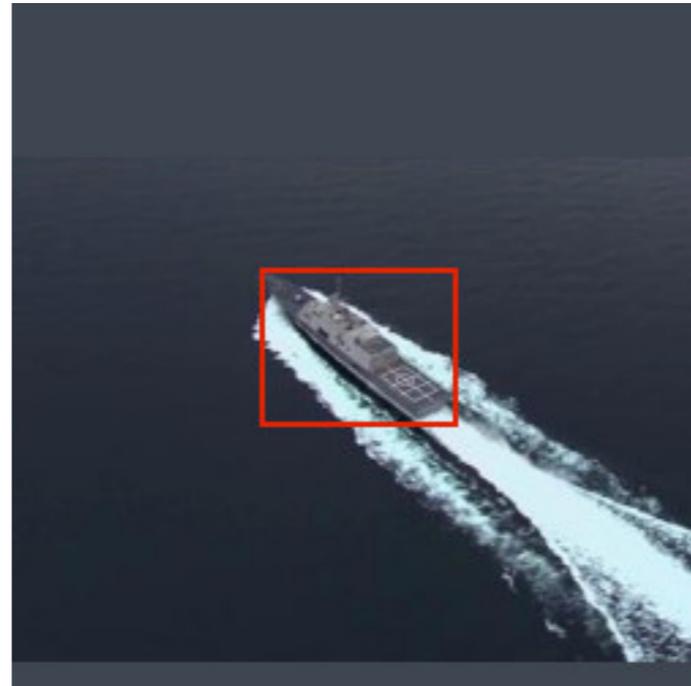
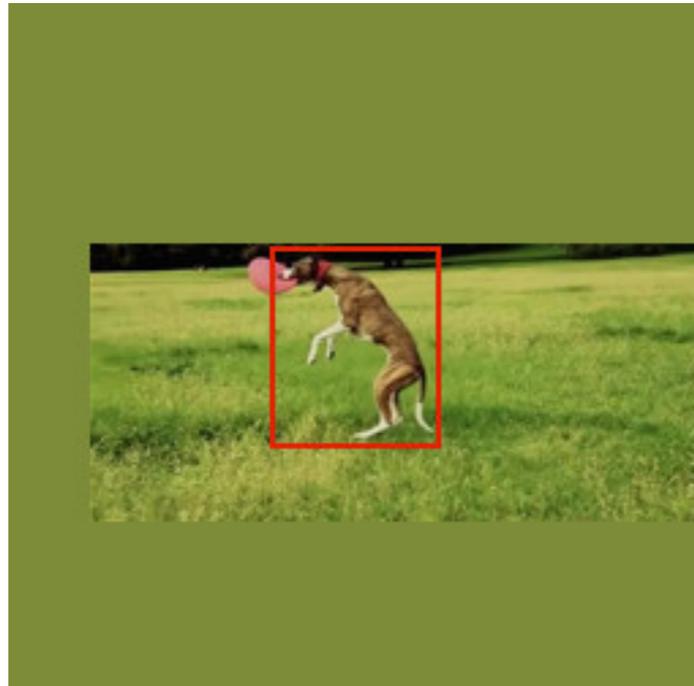
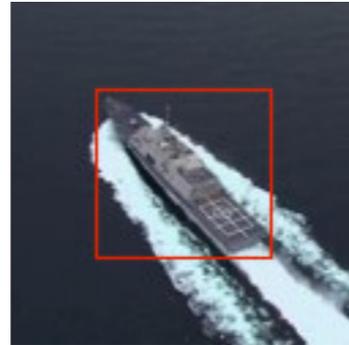
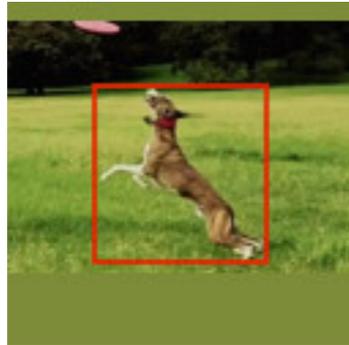
$$\varphi\{\mathbf{e}_i * \mathbf{x}\} = \mathbf{e}_i * \varphi\{\mathbf{x}\}$$

$$\mathbf{e}_i = [0, \dots, 1, \dots, 0]^T$$

$i$ -th index

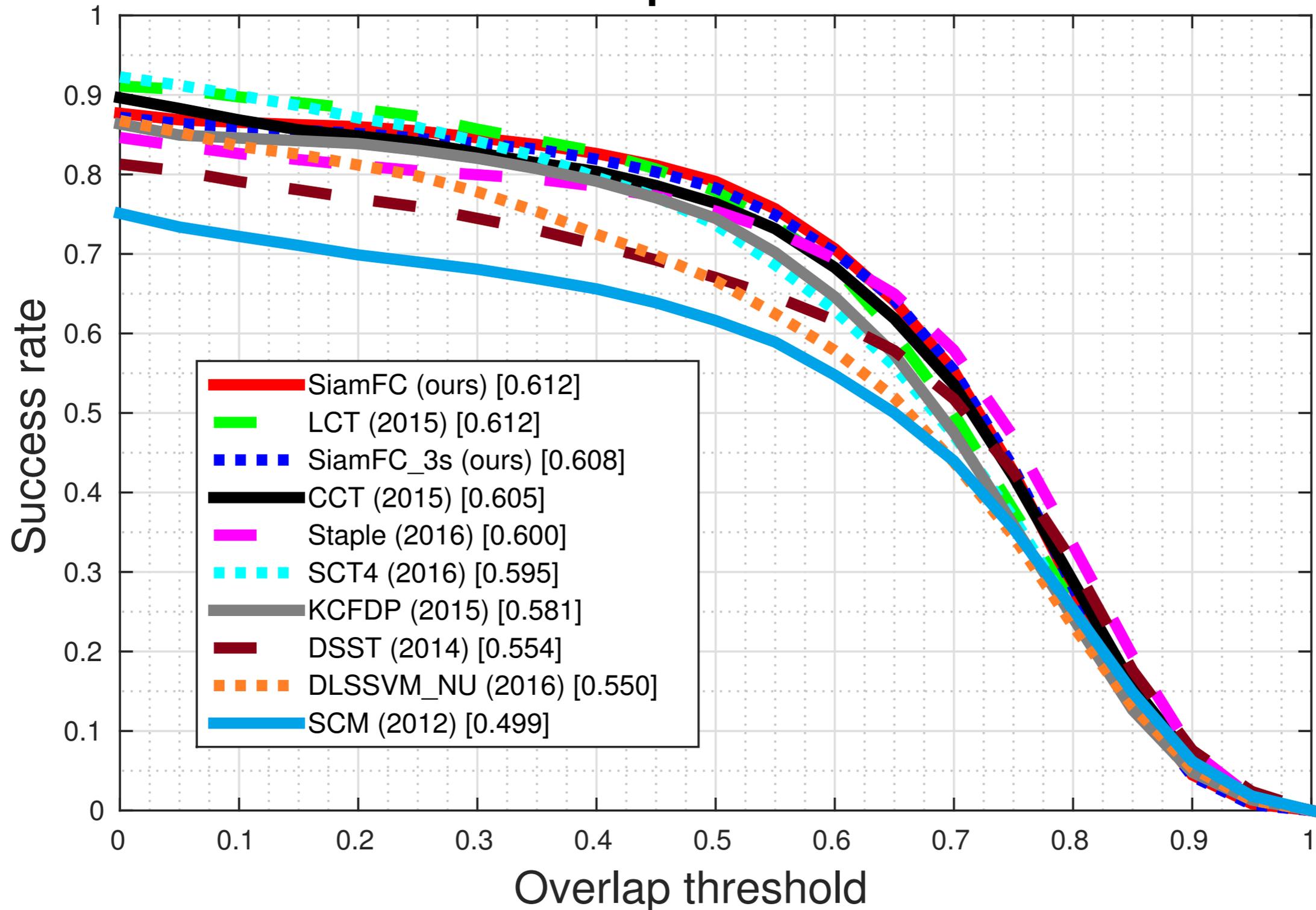
# Fully Convolutional Siamese Networks

- Example training sequences.

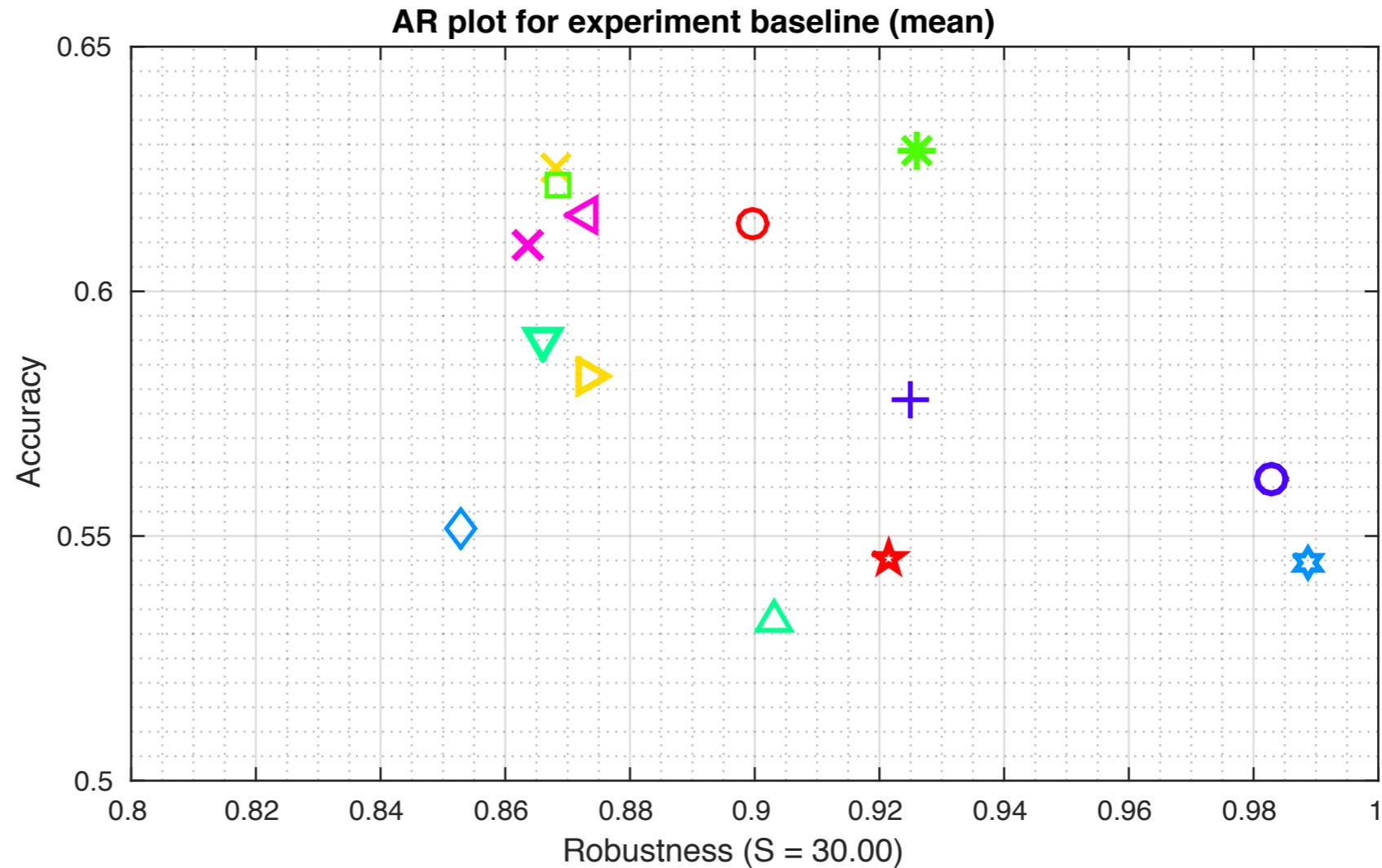
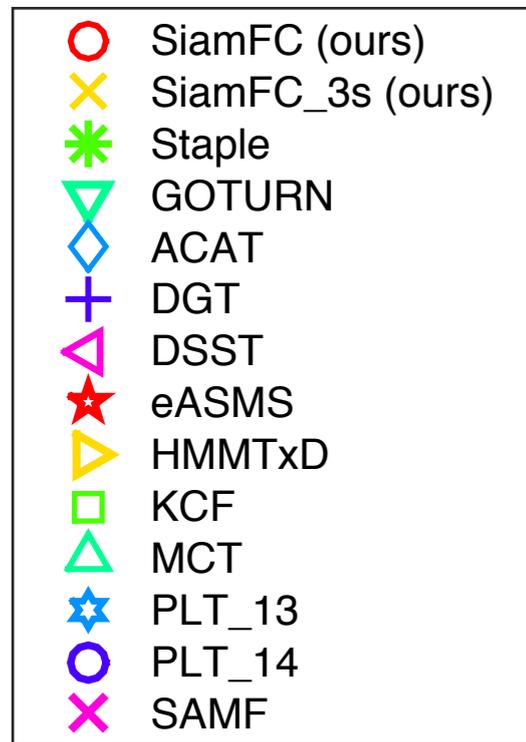


# Fully Convolutional Siamese Networks

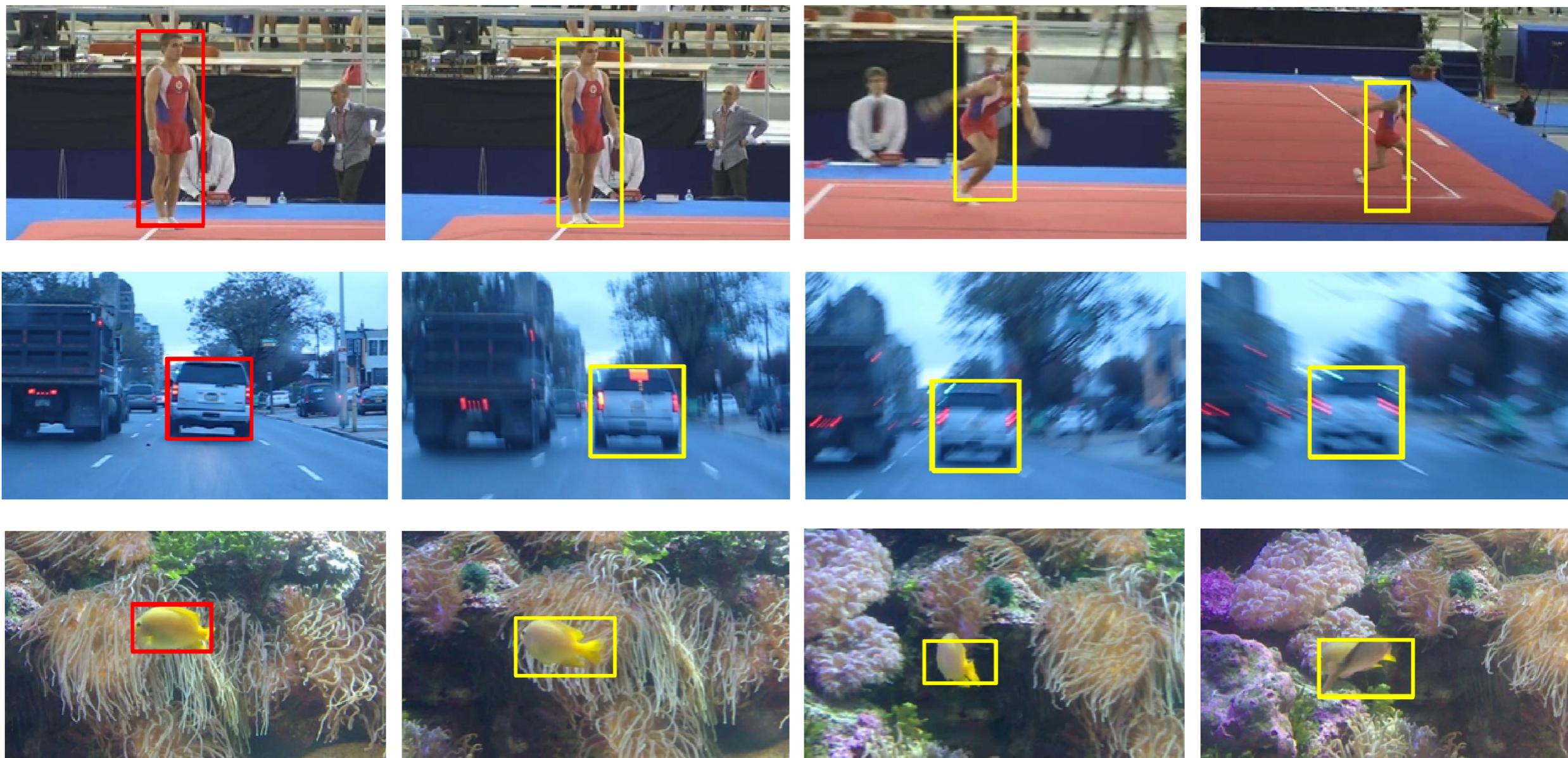
Success plots of OPE



# Fully Convolutional Siamese Networks



# Fully Convolutional Siamese Networks



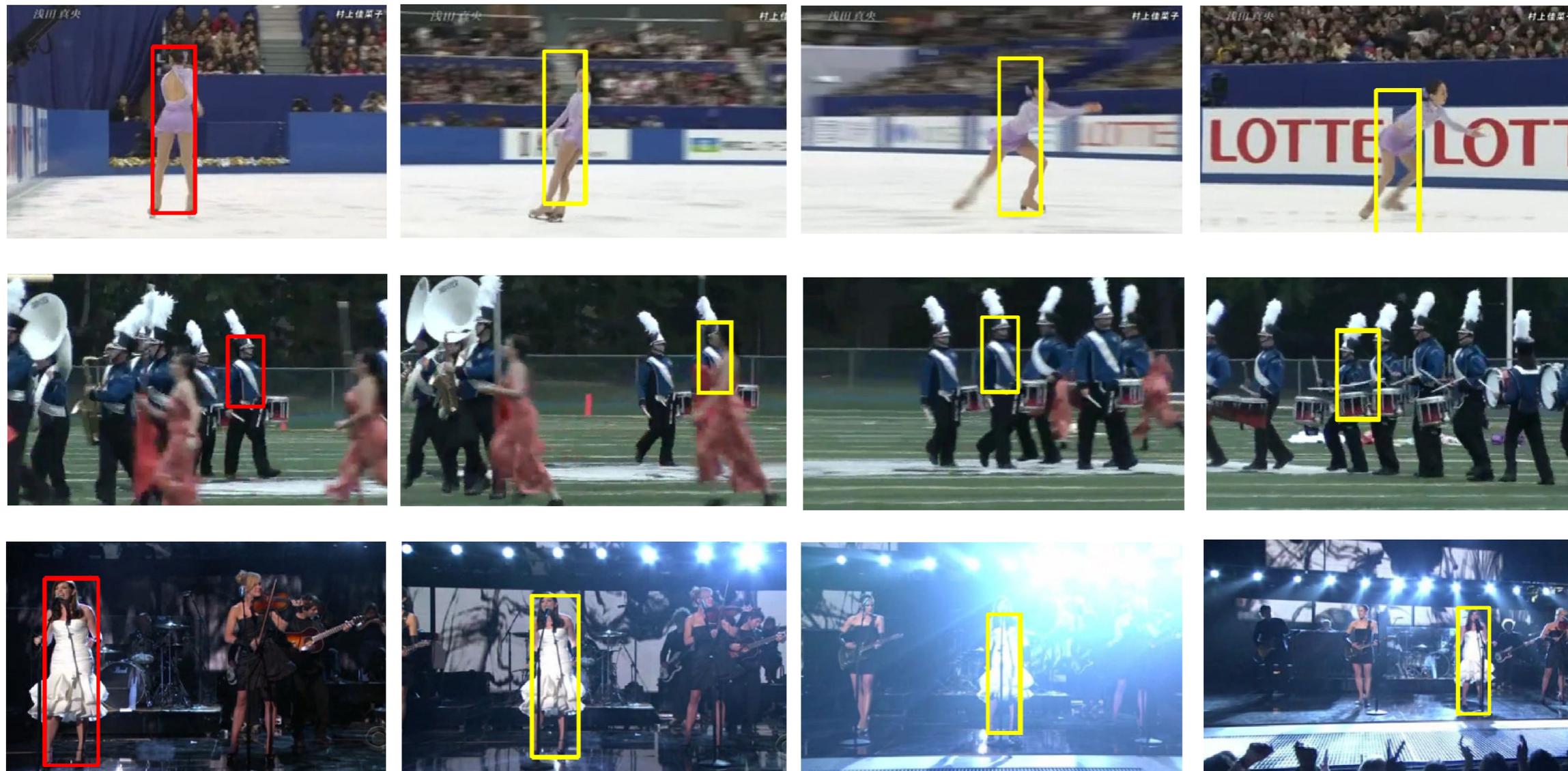
Frame 1 (init.)

Frame 50

Frame 100

Frame 200

# Fully Convolutional Siamese Networks



Frame 1 (init.)

Frame 50

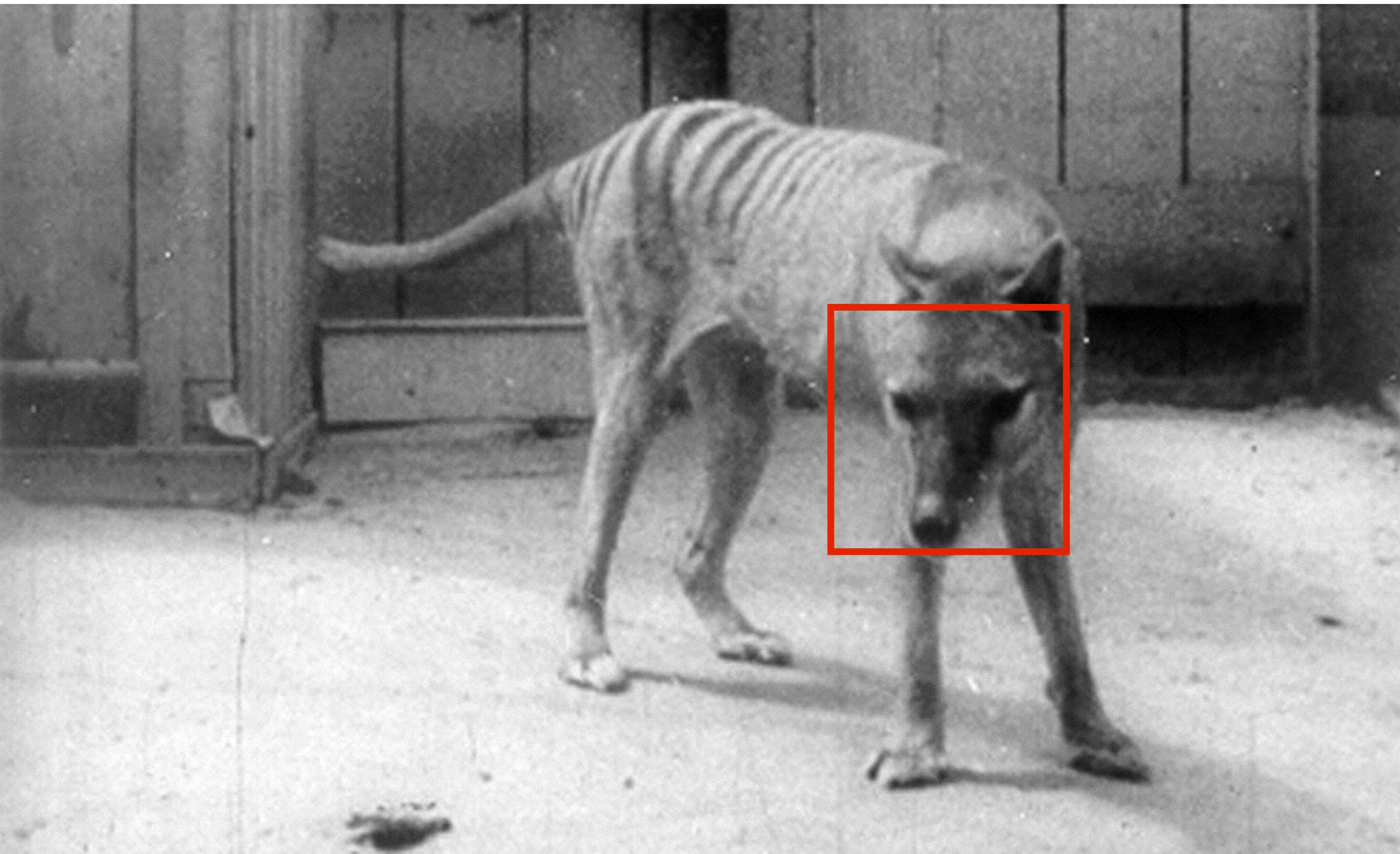
Frame 100

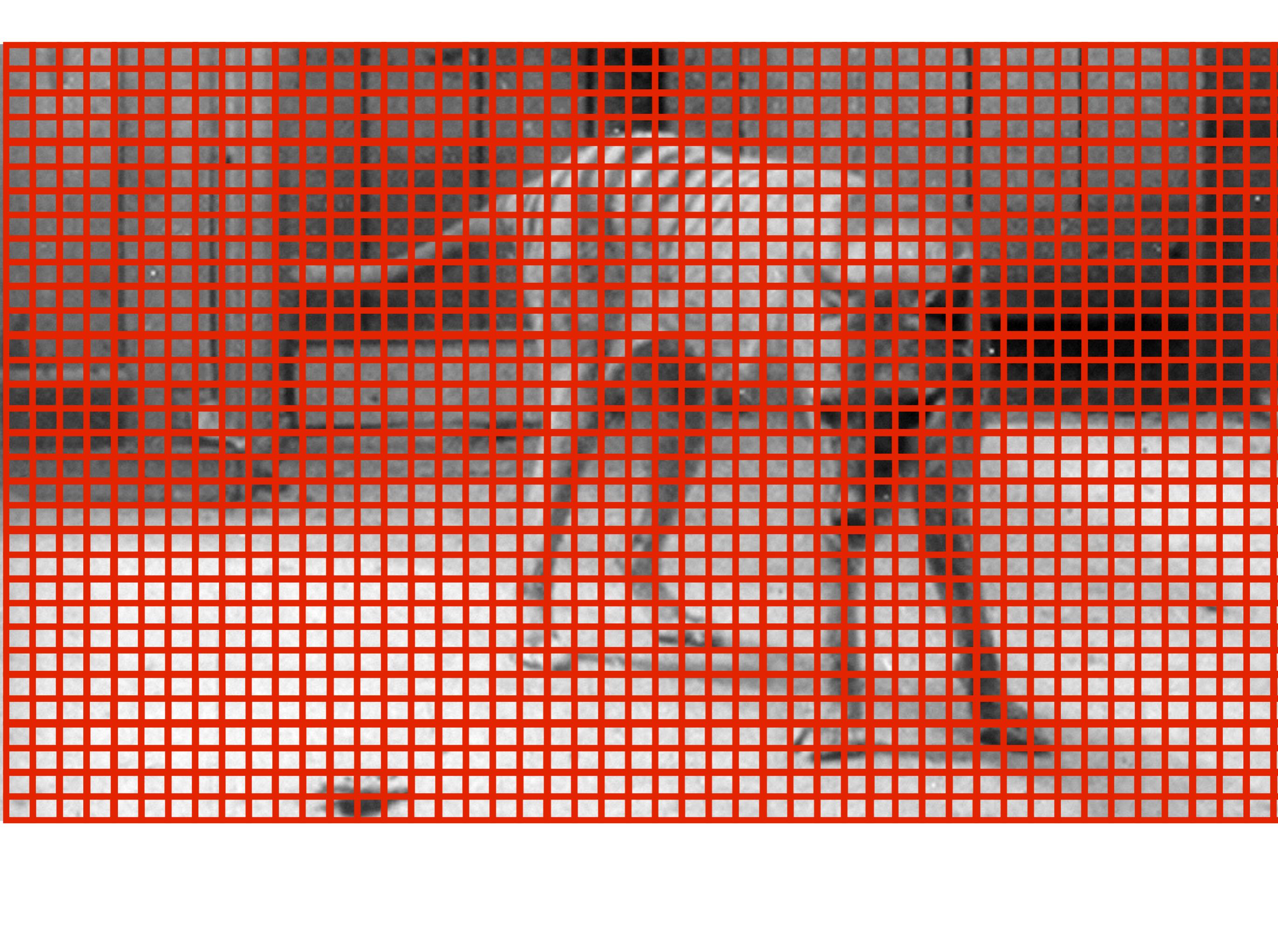
Frame 200

# Today

- Deep Features
- Deep Tracking
- Deep Flow







# Flow = Parts Based Registration

$$\min_{\mathbf{x}} \sum_{i=1}^N D_i(\mathbf{x}_i) + \lambda R(\mathbf{x})$$

# Flow = Parts Based Registration

Does the image at  $\mathbf{x}_i$  look like the  $i^{\text{th}}$  part?

$$\min_{\mathbf{x}} \sum_{i=1}^N D_i(\mathbf{x}_i) + \lambda R(\mathbf{x})$$

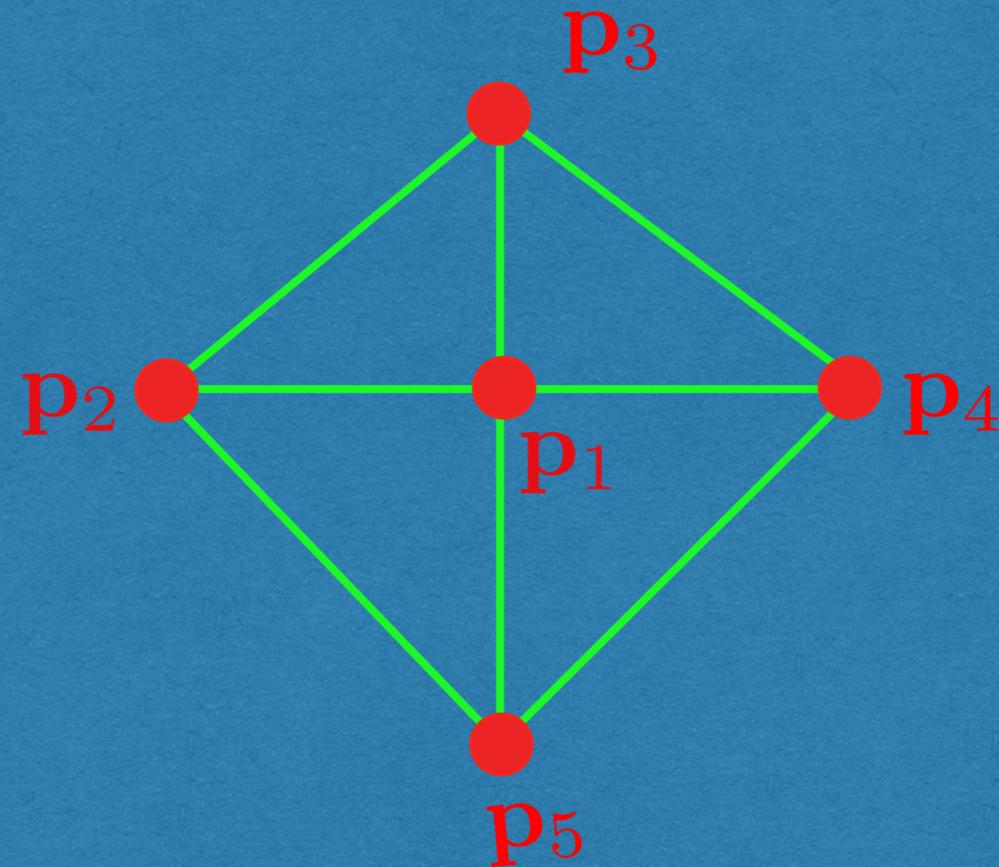
# Flow = Parts Based Registration

Does the image at  $\mathbf{x}_i$  look like the  $i^{\text{th}}$  part?

$$\min_{\mathbf{x}} \sum_{i=1}^N D_i(\mathbf{x}_i) + \lambda R(\mathbf{x})$$

Do the joint locations of the parts match the object?

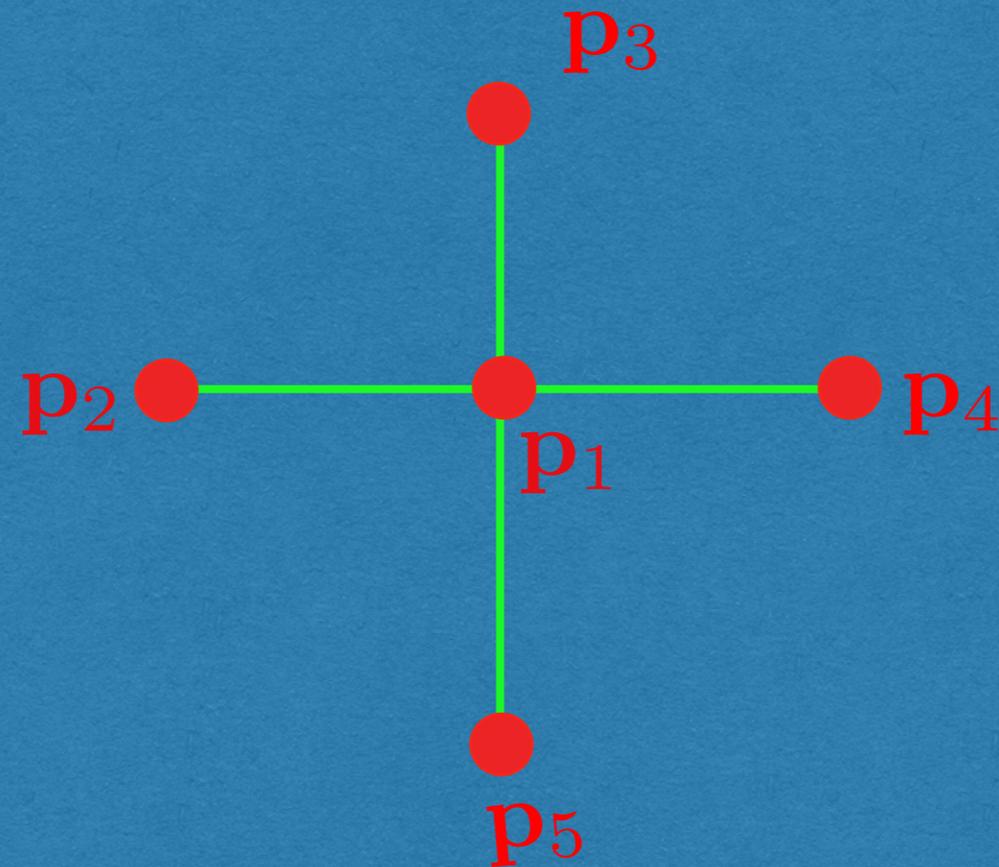
# Reminder - Exhaustive Search



$$\mathcal{O}(M^N)$$

*“We can do much better than this if the graph is sparse.”*

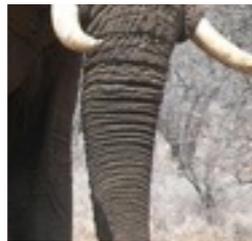
# Reminder - Exhaustive Search



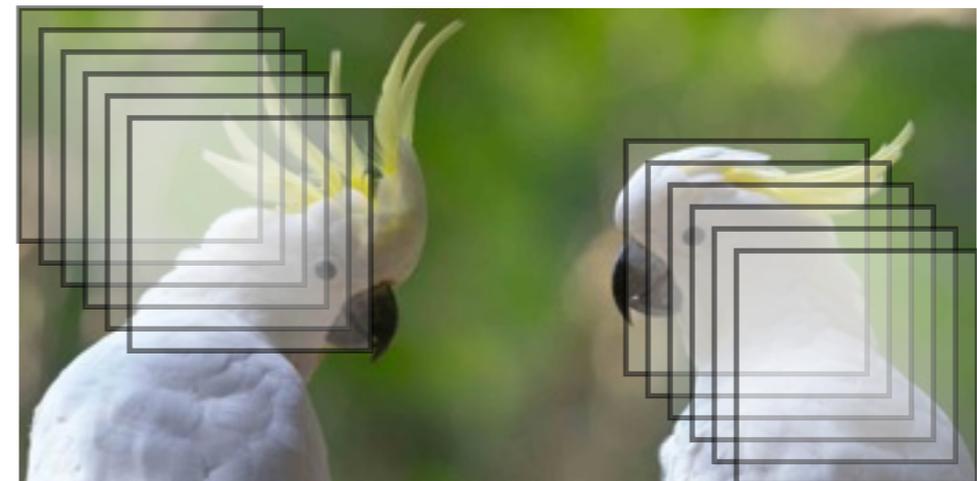
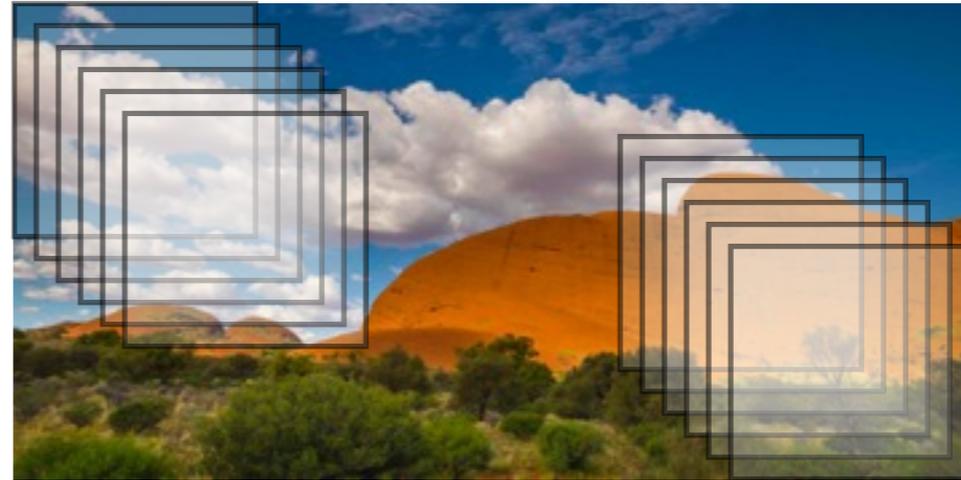
$$\mathcal{O}(NM^2)$$

*“We can do much better than this if the graph is sparse.”*

# Learning $\{D_i(\mathbf{x}_i)\}_{i=1}^N$



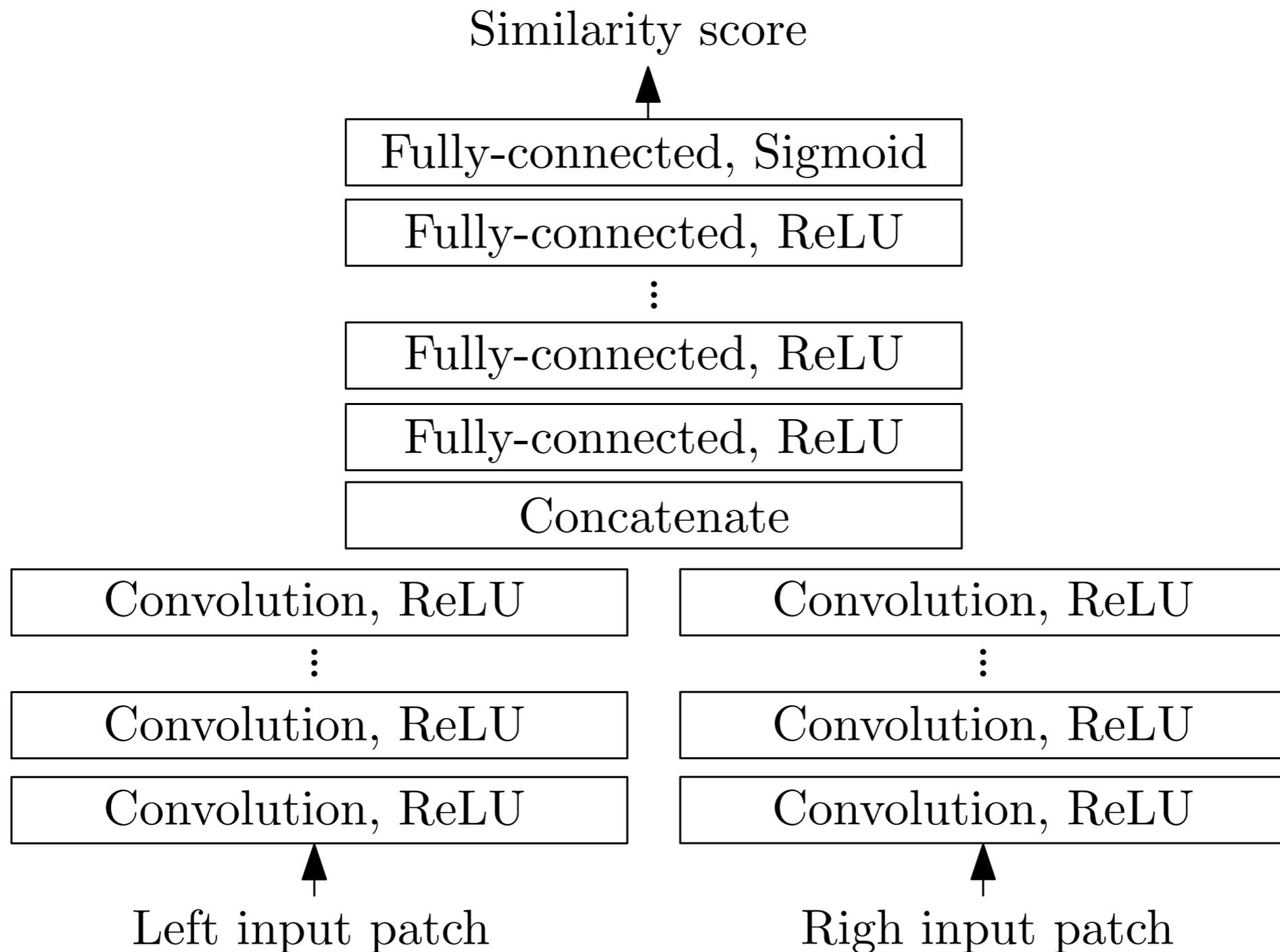
patch at pixel  $u$   
(positive)



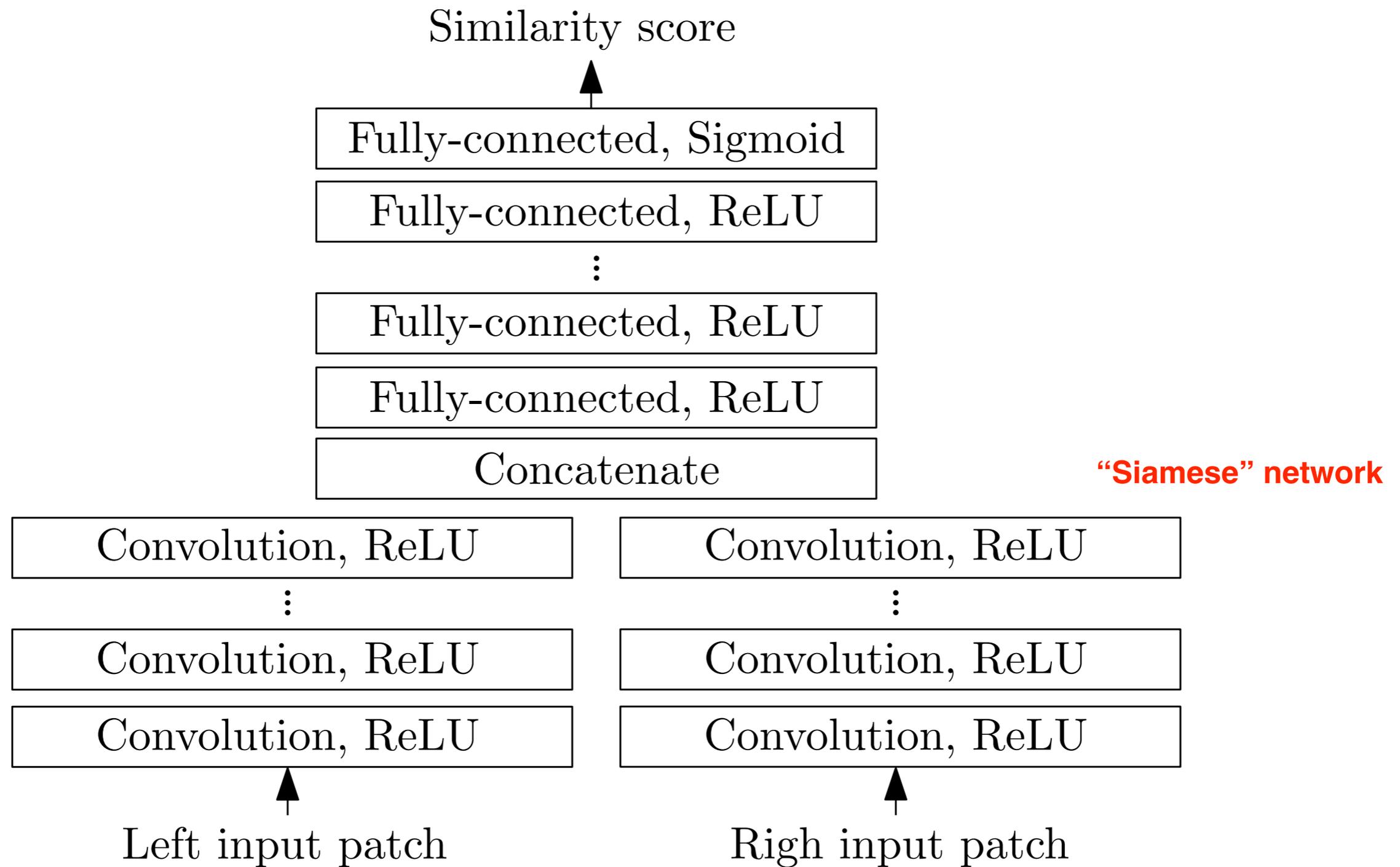
all other patches  
(negatives)

for every pixel  $u$

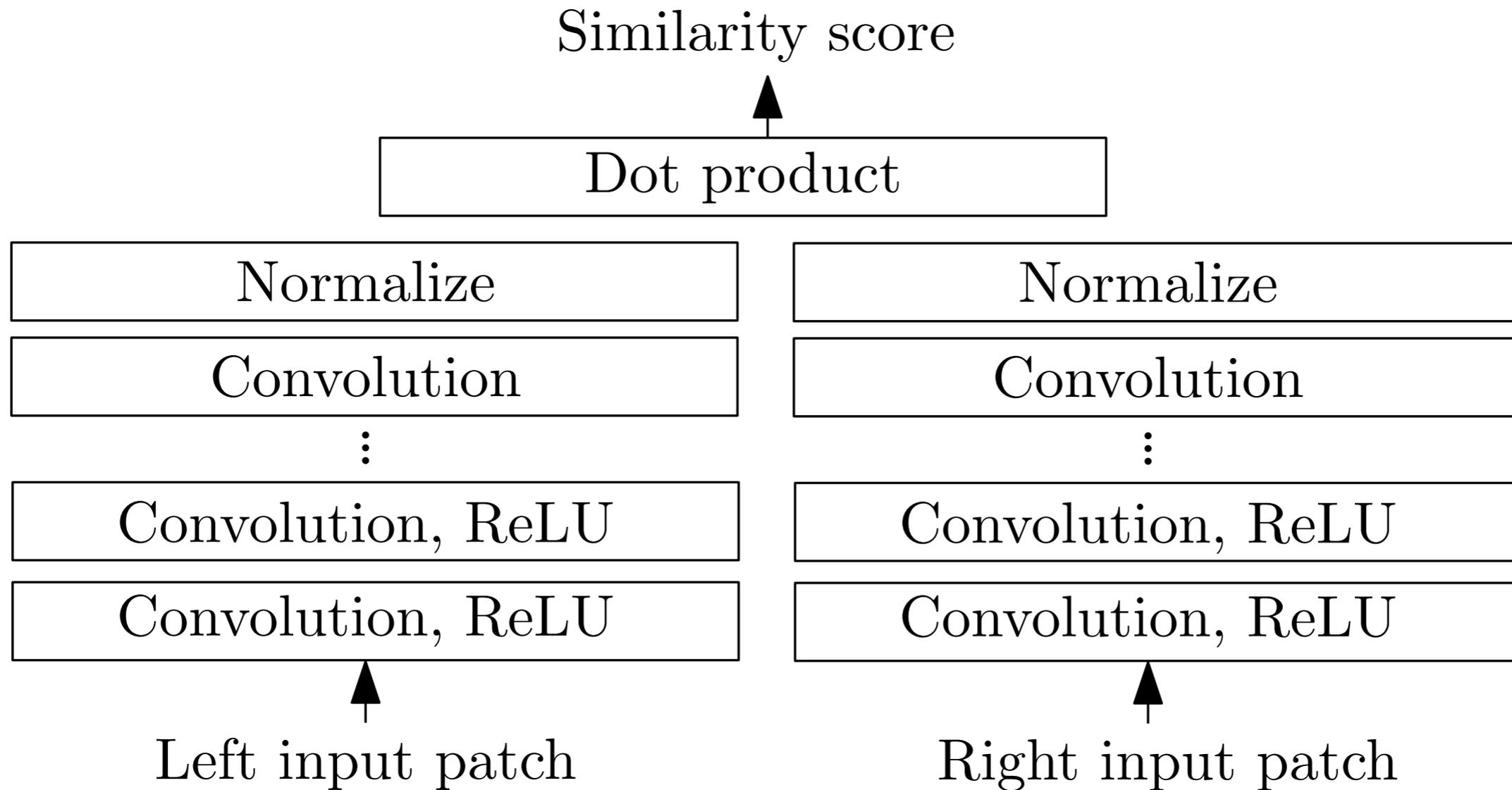
# Learning $\{D_i(\mathbf{x}_i)\}_{i=1}^N$



# Learning $\{D_i(\mathbf{x}_i)\}_{i=1}^N$



# Fast Architecture



# Results - KITTI 2015

Left input image



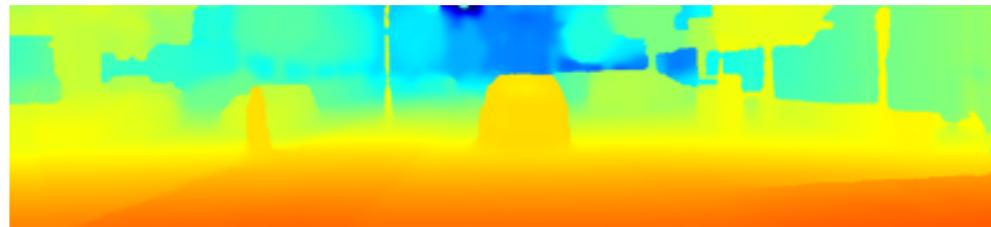
Right input image



Ground truth



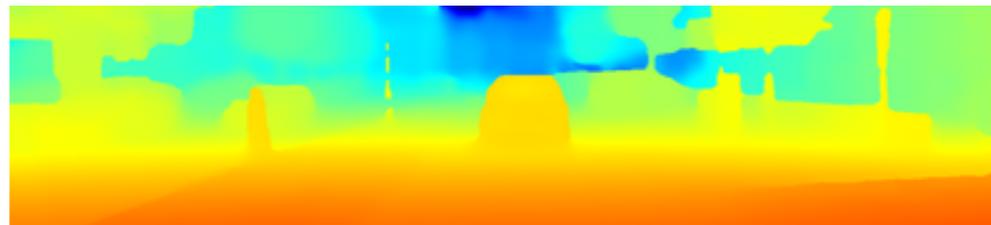
Fast architecture



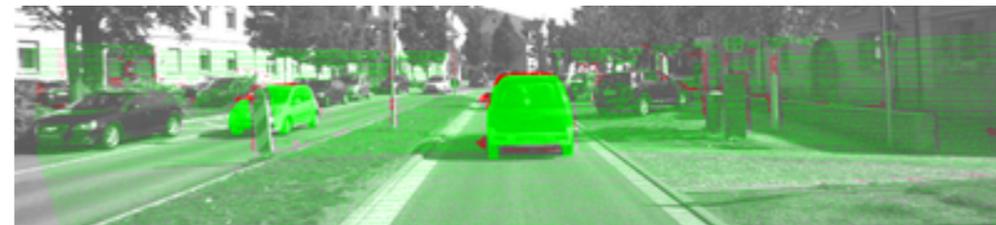
Error: 2.79 %



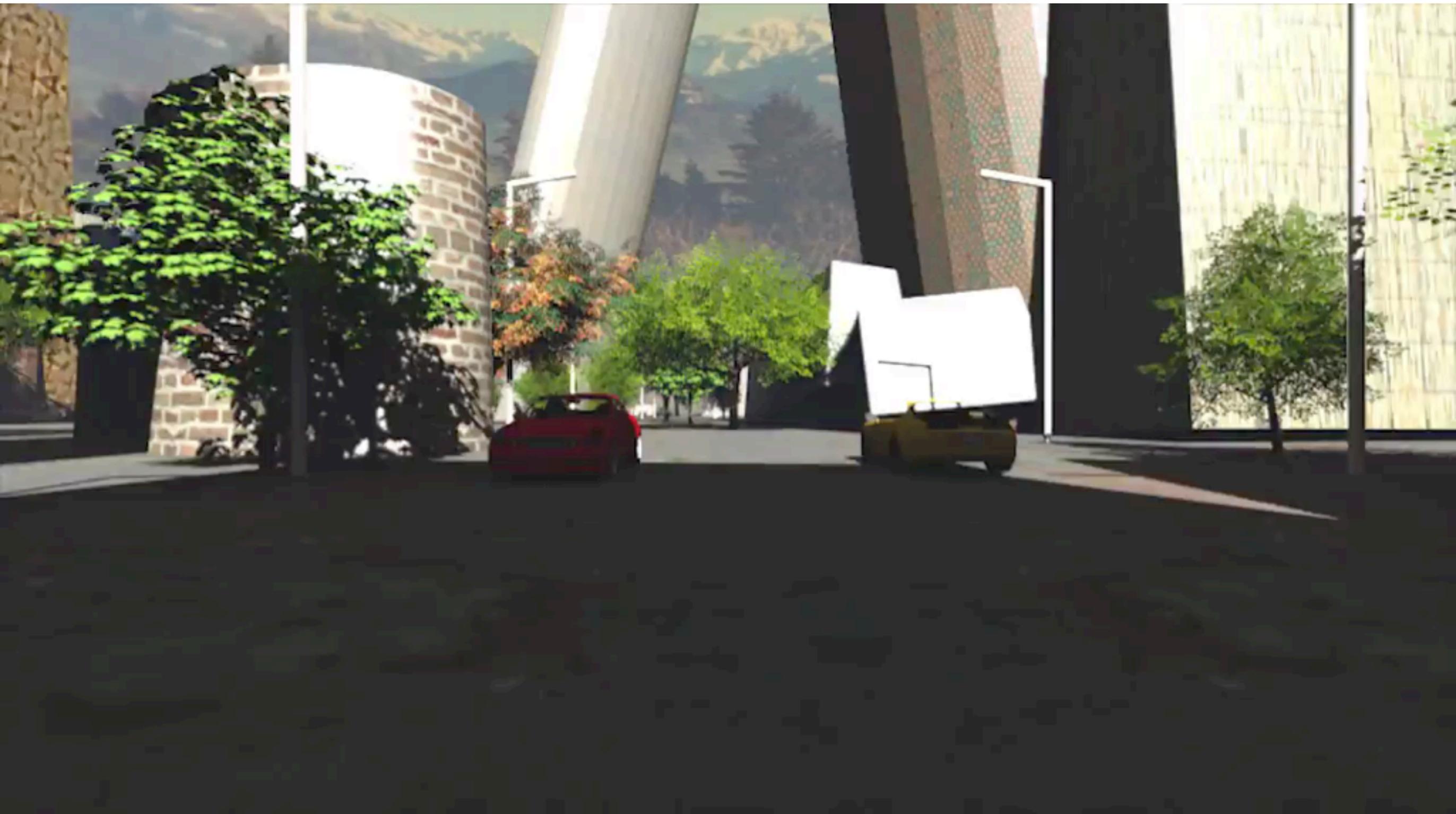
Accurate architecture



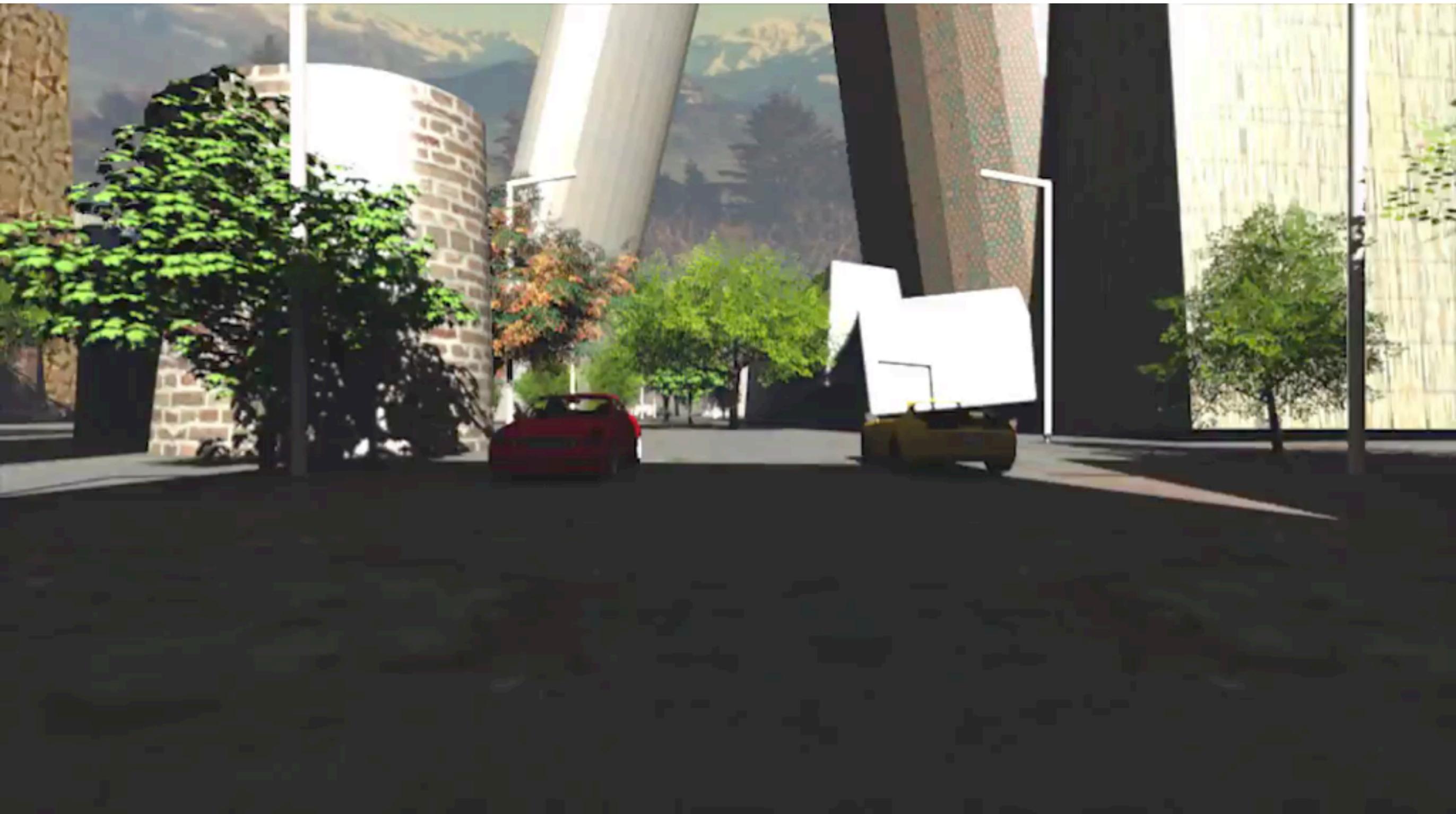
Error: 2.36 %



# More Data?

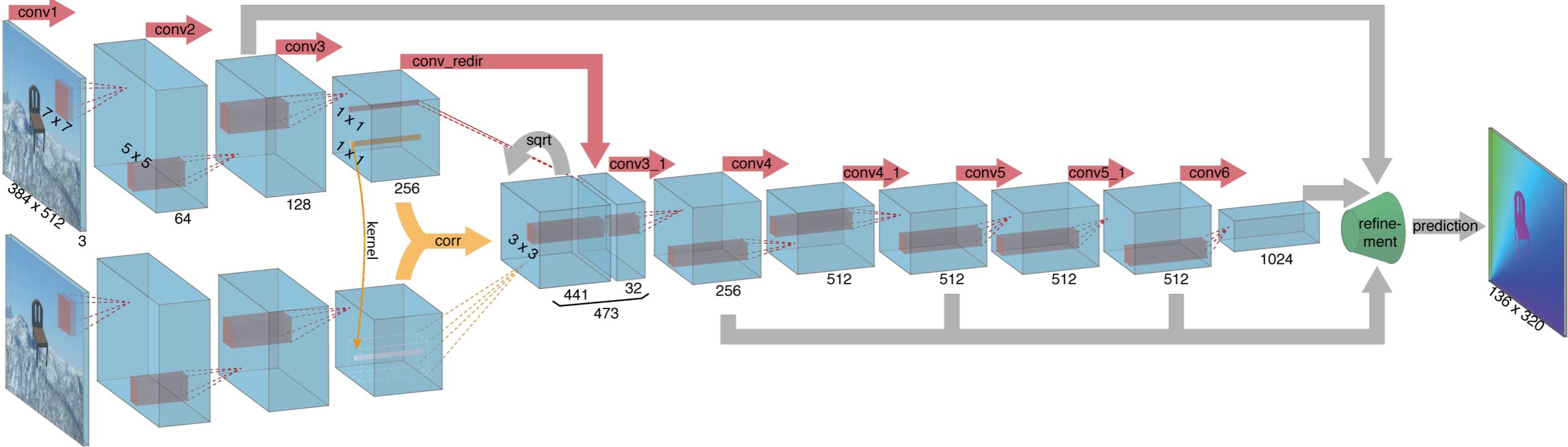


# More Data?



# FlowNet

## FlowNetCorr





# FlowNet - Results

