

# Machine Learning in Imaging

BME 590L  
Roarke Horstmeyer

Lecture 13: CNN visualization and example applications

## **Post-processing of your results: a few options at different stages**

Options to examine your test data after processing:

- ROC curve, Precision-Recall
- Confusion matrix
- Sliding window visualization
- Layer visualizations
- Saliency maps etc.
- tSNE visualization

## ROC curve and confusion matrix

- Can set threshold for  $f(x, W)$  wherever
- Leads to sliding window between FN and FP rate
- Need to summarize both statistics as a function of sliding window

		Estimated label $f(x, W)$		
		+1	-1	Missed an event
Actual label $y$	+1	True positive	False negative	Predict event when there isn't one
	-1	False positive	True negative	

## ROC curve and confusion matrix

TP Rate =

Sensitivity =  $TP / (TP + FN) = TP / \text{Actual positives}$

False Positive Rate =  $FP / (TN + FP) = FP / \text{Actual negatives}$

Specificity =  $TN / (TN + FP) = TN / \text{Actual negatives}$   
 $= 1 - \text{False Positive Rate}$

Actual label  
y

		Estimated label $f(x, W)$		
		+1	-1	Missed an event
+1	True positive	False negative		Missed an event
	False positive	True negative		
-1	Predict event when there isn't one			
	True negative	False positive	True positive	False negative

## ROC curve and confusion matrix

TP Rate =

Sensitivity =  $TP / (TP + FN) = TP / \text{Actual positives}$

False Positive Rate =  $FP / (TN + FP) = FP / \text{Actual negatives}$

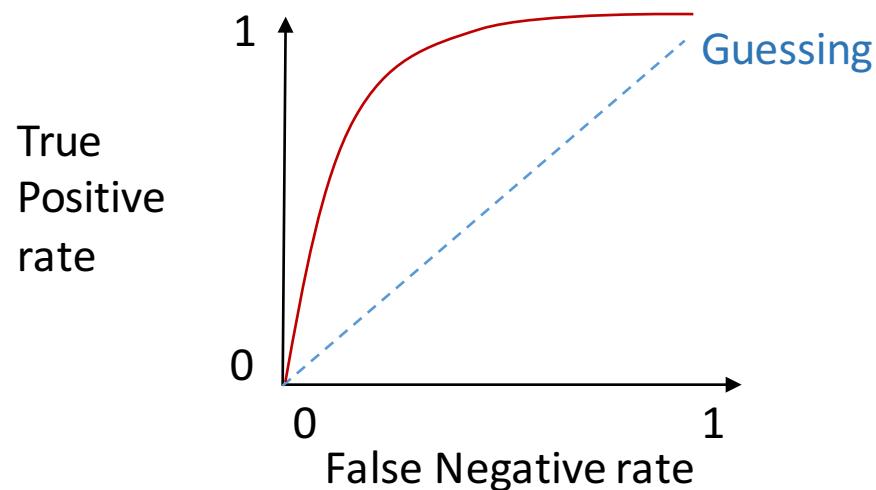
Specificity =  $TN / (TN + FP) = TN / \text{Actual negatives}$   
 $= 1 - \text{False Positive Rate}$

Actual label  
y

		Estimated label $f(x, W)$	
		+1	-1
+1	True positive	False negative	
	False positive	True negative	

Predict event when  
there isn't one

Receiver-Operator Curve



## ROC curve and confusion matrix

Recall =

Sensitivity =  $TP / (TP + FN) = TP / \text{Actual positives}$

→ **Actual label**  
*y*

Estimated label

$f(x, W)$

+1

-1

Missed an event

		+1	-1
+1	True positive	False negative	
	False positive	True negative	
-1			

Predict event when  
there isn't one

Precision =  $TP / (TP + FP) = TP / \text{Estimated positives}$

- Sometimes, you don't care about true negatives (just want to find events)
- In this case, use Precision and Recall

# ROC curve and confusion matrix

Recall =

Sensitivity =  $TP / (TP + FN) = TP / \text{Actual positives}$

→ Actual label  
 $y$

Estimated label

$f(x, W)$

+1

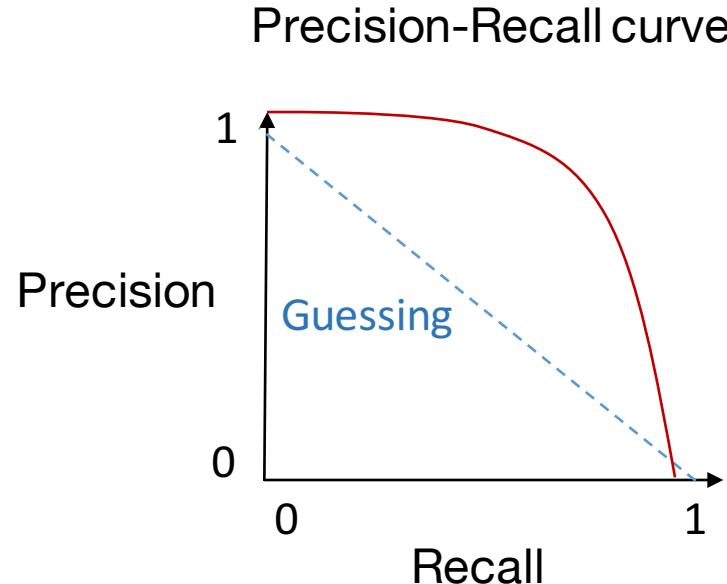
-1

Missed an event

	+1	-1
+1	True positive	False negative
-1	False positive	True negative

Predict event when  
there isn't one

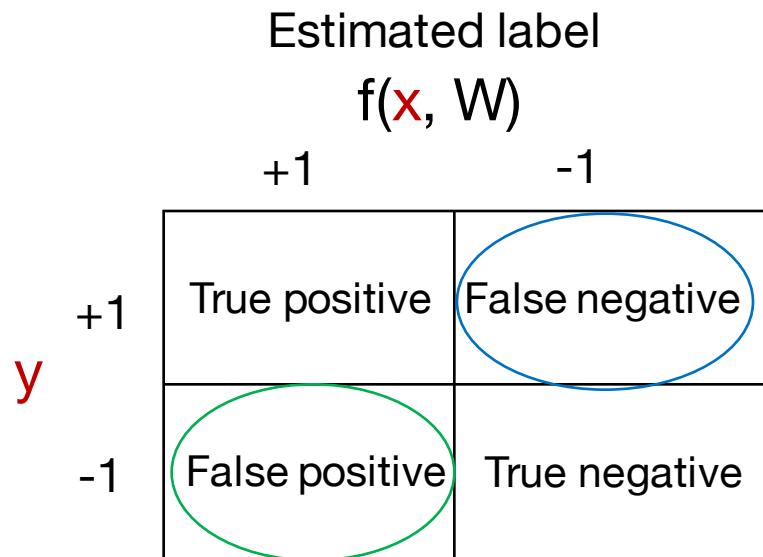
Precision =  $TP / (TP + FP) = TP / \text{Estimated positives}$



F1 Metric:  $(1/\text{precision} + 1/\text{recall})^{-1}$

# ROC curve and confusion matrix

Just 2 categories



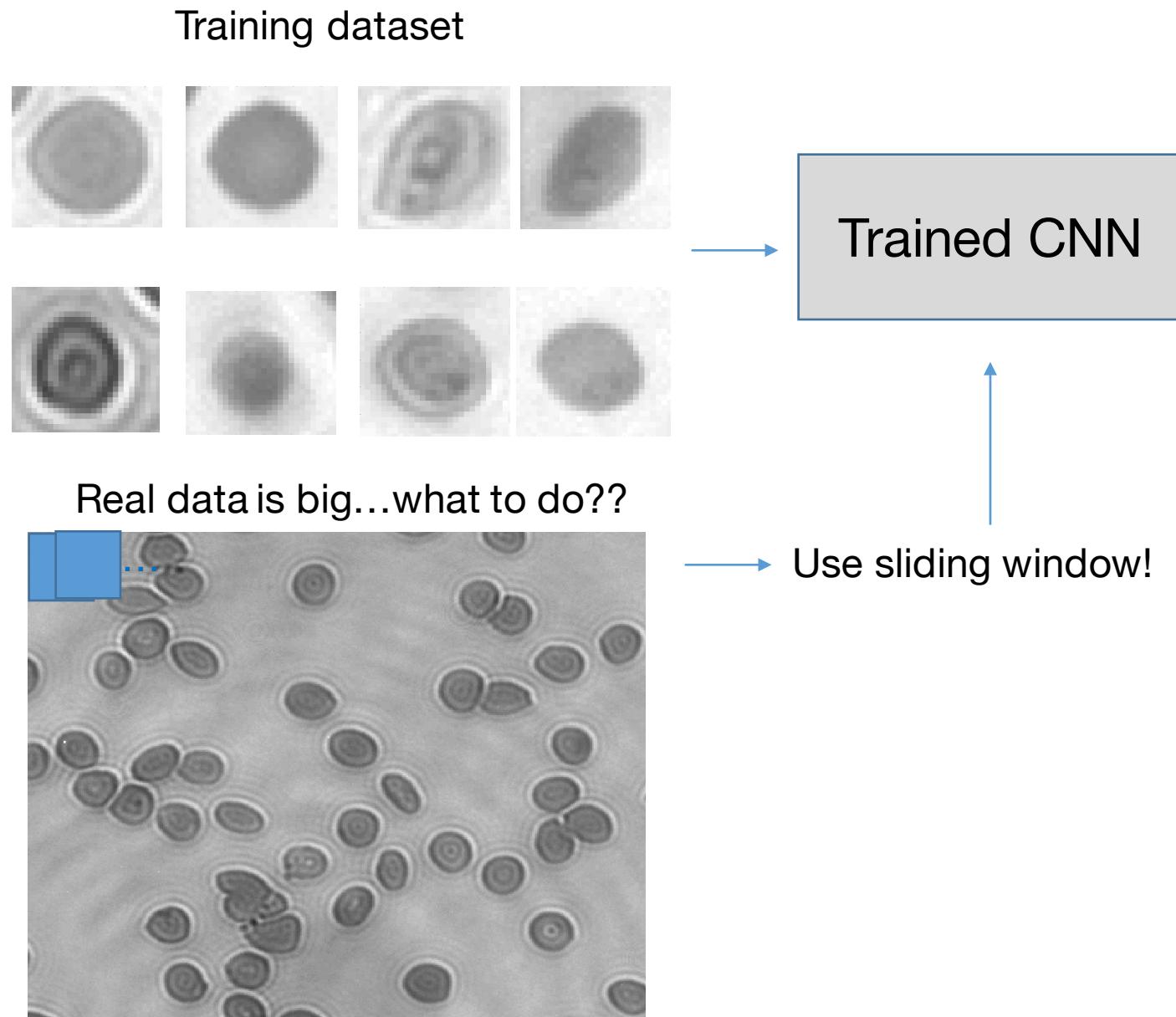
Confusion Matrix: 2+ categories

Estimated label

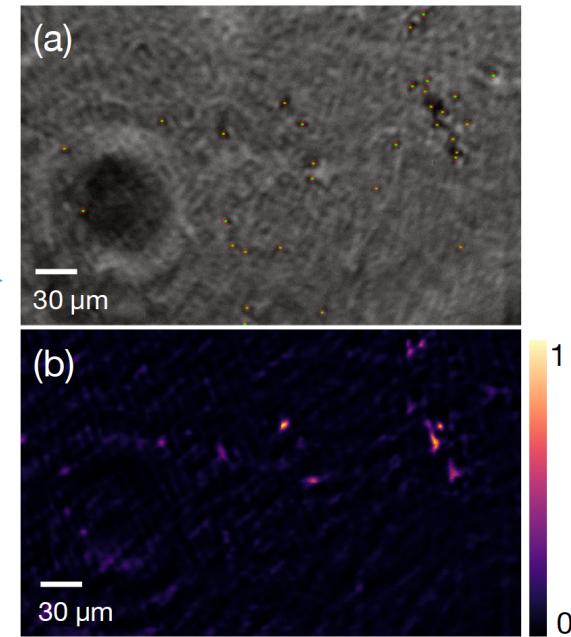
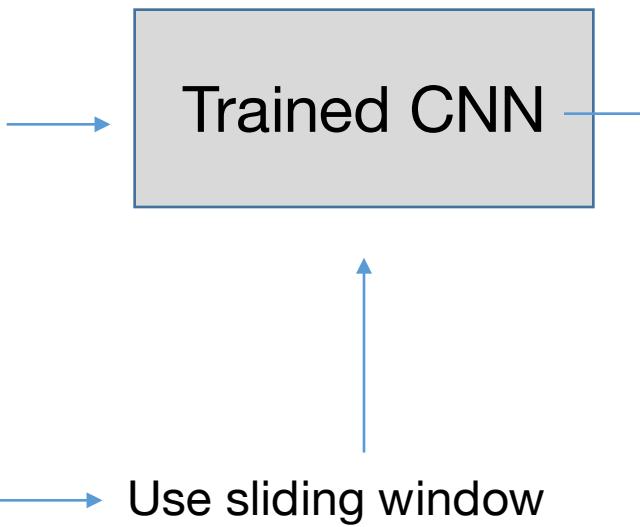
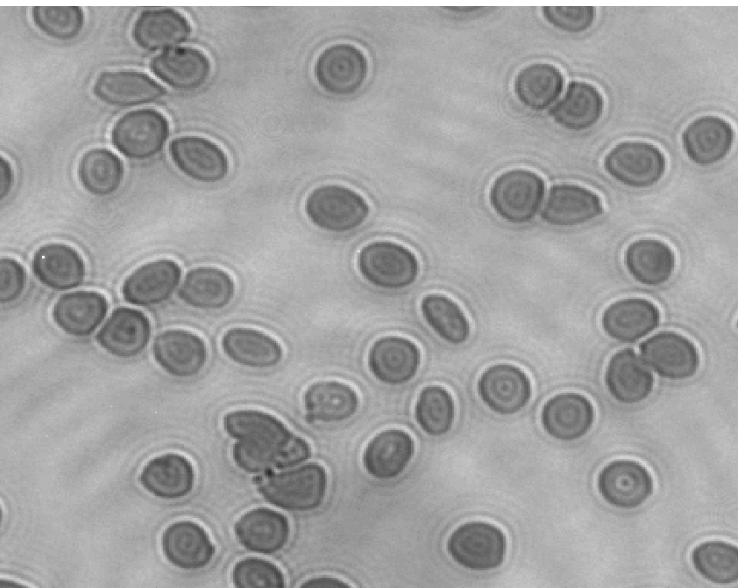
$f(\mathbf{x}, \mathbf{W})$

Actual ↓	State1 (Predicted)	State2 (Predicted)	State3 (Predicted)	State4 (Predicted)	State5 (Predicted)	State6 (Predicted)	State7 (Predicted)	State8 (Predicted)
State1 (Actual)	90.12 %	0.00 %	9.88 %	0.00 %	0.00 %	0.00 %	0.00 %	0.00 %
State2 (Actual)	0.00 %	100.00 %	0.00 %	0.00 %	0.00 %	0.00 %	0.00 %	0.00 %
State3 (Actual)	0.00 %	0.00 %	92.66 %	0.00 %	0.00 %	7.34 %	0.00 %	0.00 %
State4 (Actual)	0.00 %	0.00 %	0.00 %	100.00 %	0.00 %	0.00 %	0.00 %	0.00 %

# How can we visualize performance?



# How can we visualize performance?

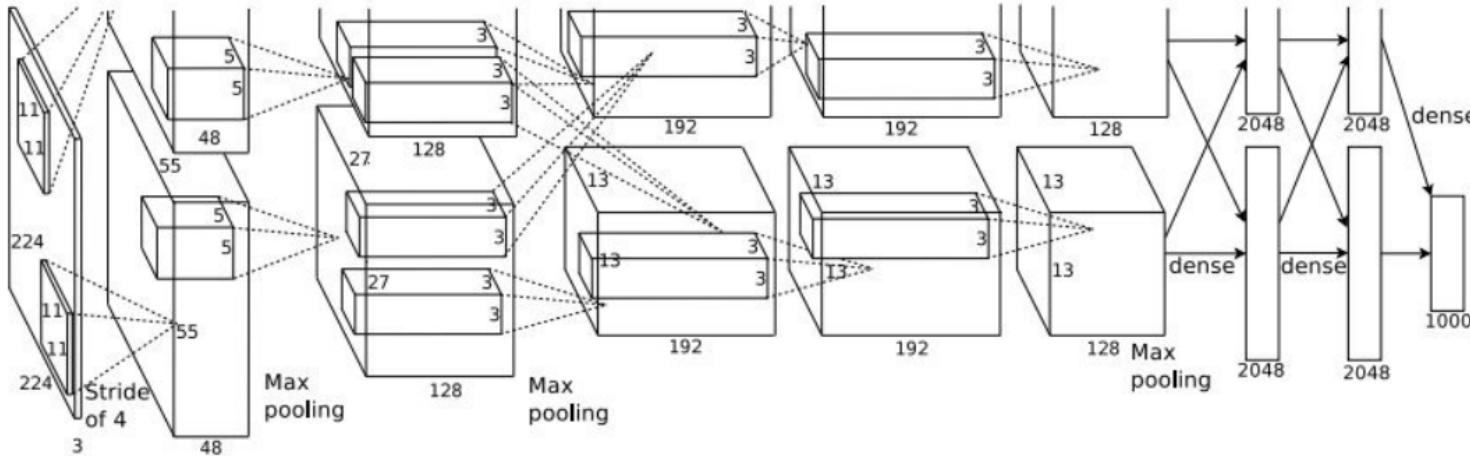


# How can we visualize what's in the network?

This image is CC0 public domain



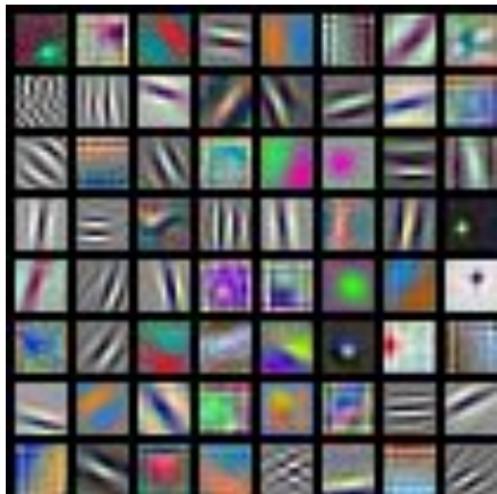
Input Image:  
 $3 \times 224 \times 224$



What are the intermediate features looking for?

Class Scores:  
1000 numbers

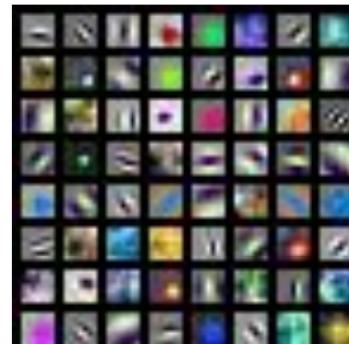
# First Layer: Visualize Filters



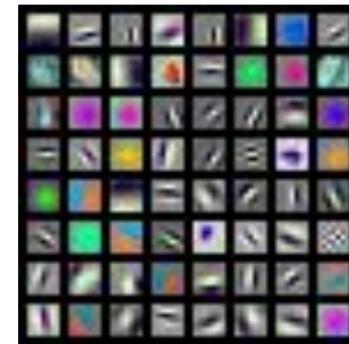
AlexNet:  
 $64 \times 3 \times 11 \times 11$



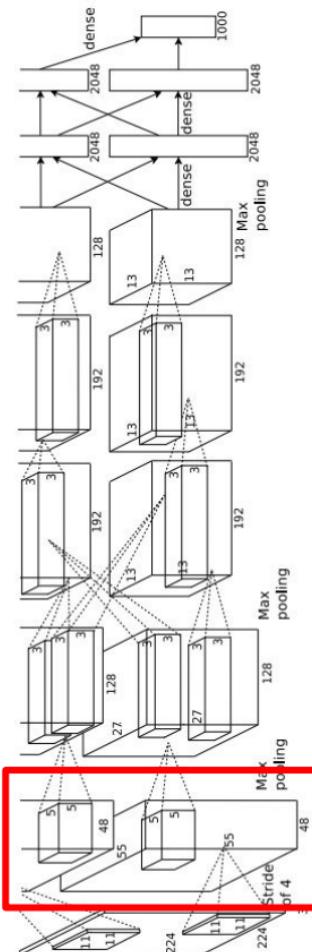
ResNet-18:  
 $64 \times 3 \times 7 \times 7$



ResNet-101:  
 $64 \times 3 \times 7 \times 7$



DenseNet-121:  
 $64 \times 3 \times 7 \times 7$



Krizhevsky, "One weird trick for parallelizing convolutional neural networks", arXiv 2014  
He et al, "Deep Residual Learning for Image Recognition", CVPR 2016  
Huang et al, "Densely Connected Convolutional Networks", CVPR 2017

# Visualize the filters/kernels (raw weights)

We can visualize filters at higher layers, but not that interesting

(these are taken from ConvNetJS CIFAR-10 demo)



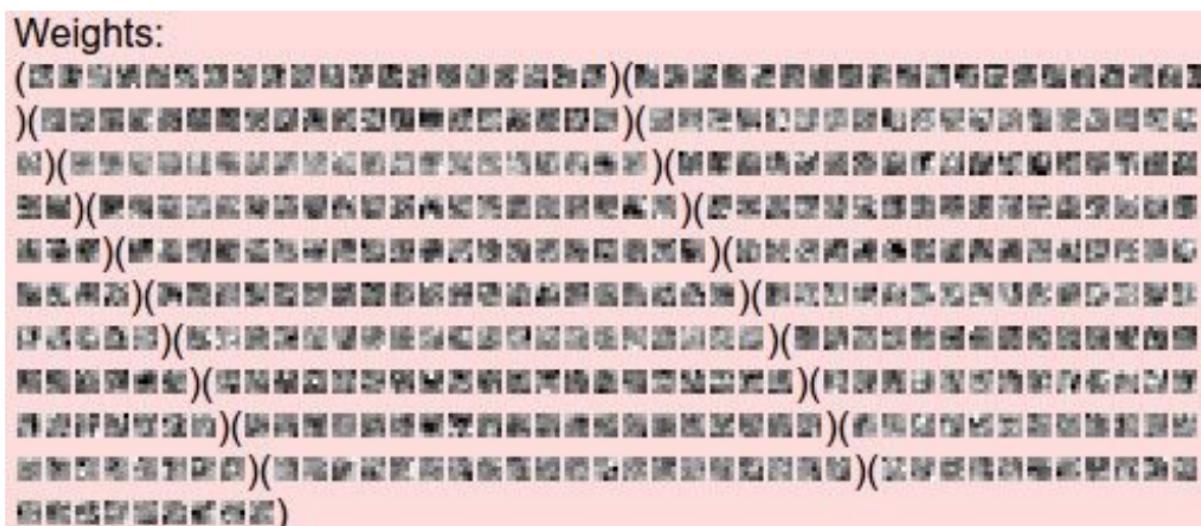
layer 1 weights

$16 \times 3 \times 7 \times 7$



layer 2 weights

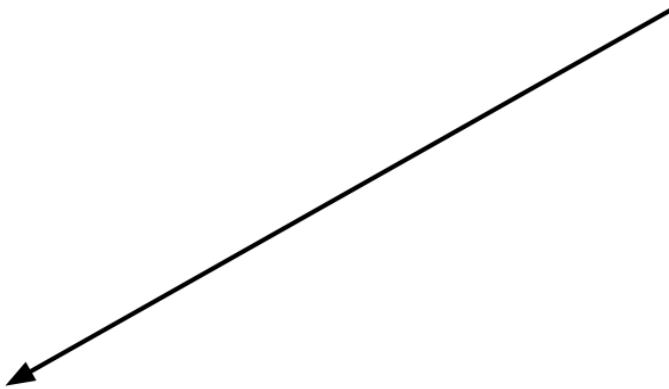
$20 \times 16 \times 7 \times 7$



layer 3 weights

$20 \times 20 \times 7 \times 7$

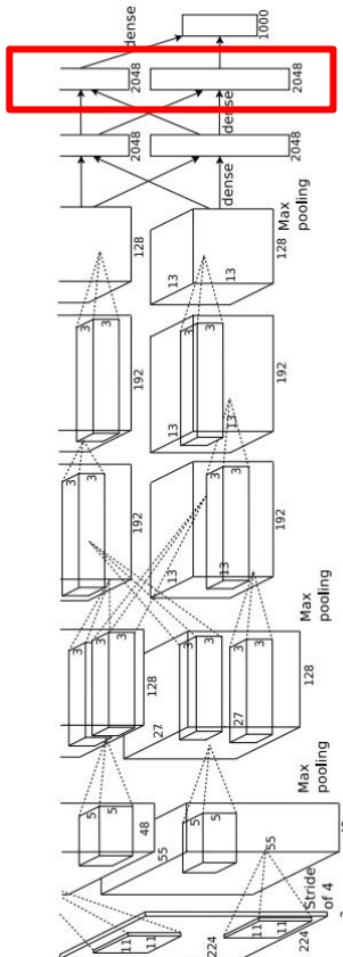
# Last Layer



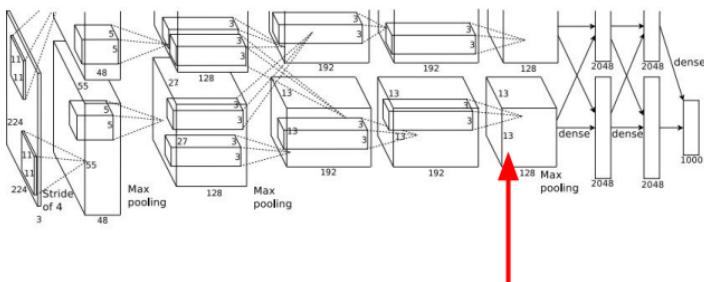
4096-dimensional feature vector for an image  
(layer immediately before the classifier)

Run the network on many images, collect the  
feature vectors

## FC7 layer



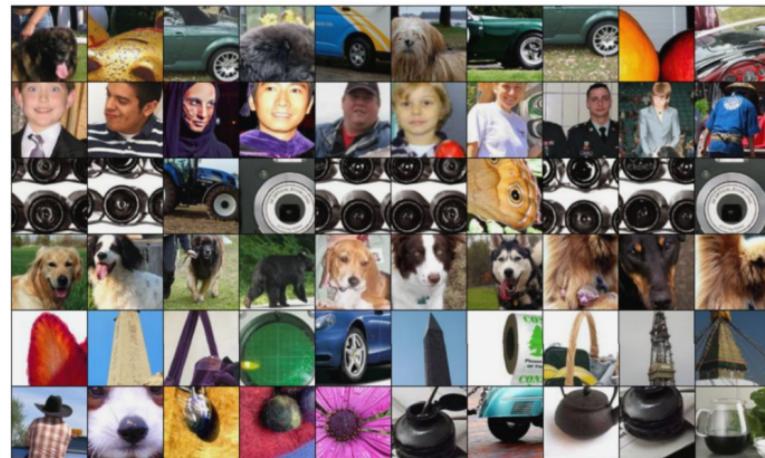
# Maximally Activating Patches



Pick a layer and a channel; e.g. conv5 is  $128 \times 13 \times 13$ , pick channel 17/128

Run many images through the network,  
record values of chosen channel

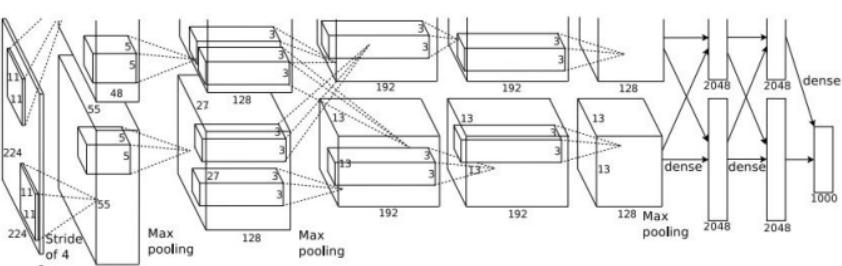
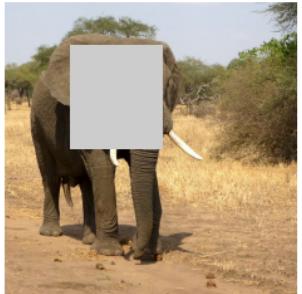
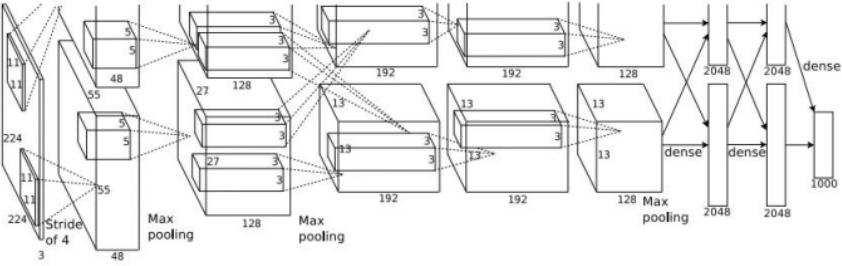
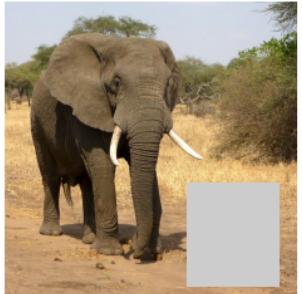
Visualize image patches that correspond  
to maximal activations



Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015  
Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015;  
reproduced with permission.

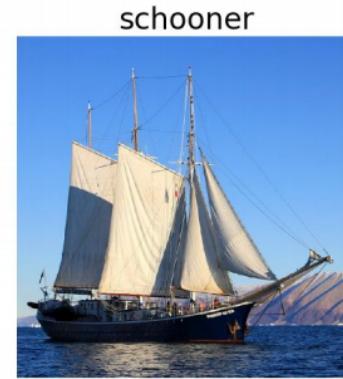
# Which pixels matter: Saliency vs Occlusion

Mask part of the image before feeding to CNN,  
check how much predicted probabilities change

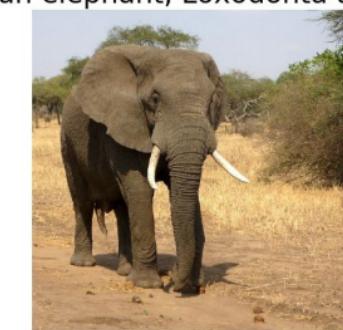
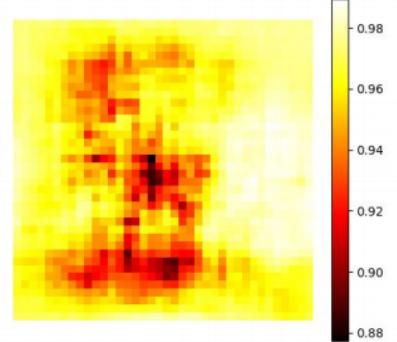


Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014

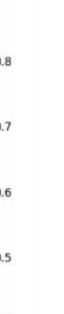
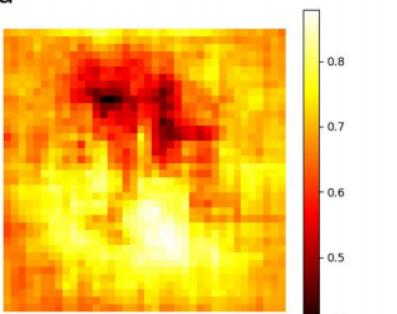
[Boat image is CC0 public domain](#)  
[Elephant image is CC0 public domain](#)  
[Go-Karts image is CC0 public domain](#)



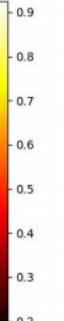
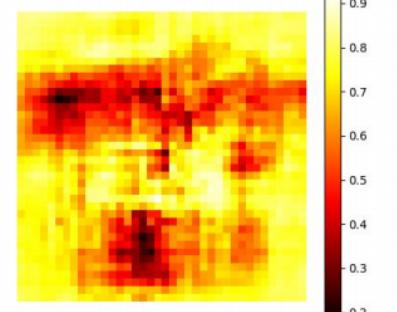
schooner



African elephant, Loxodonta africana

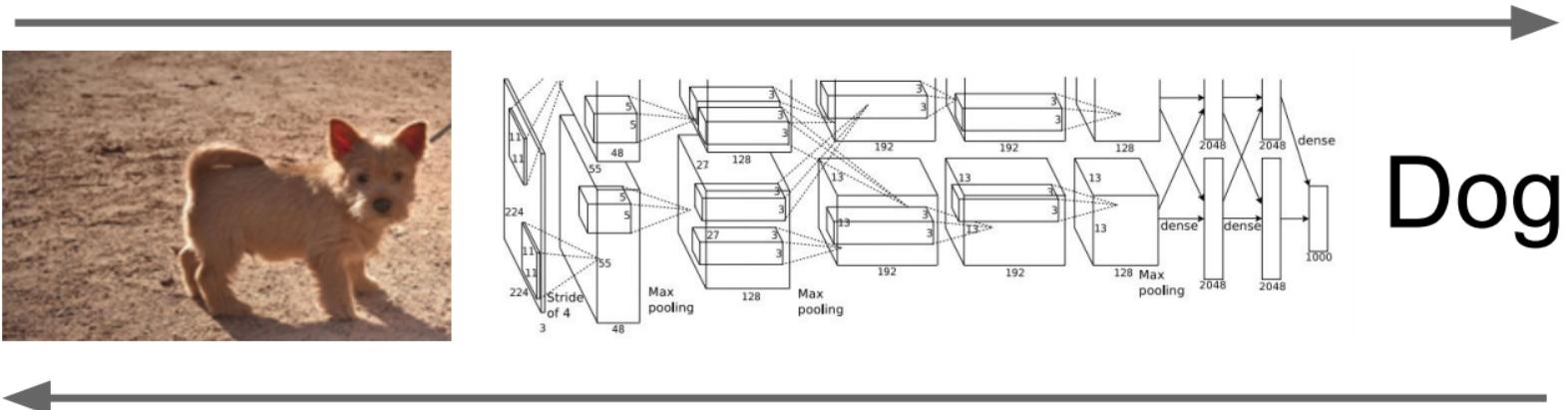


go-kart

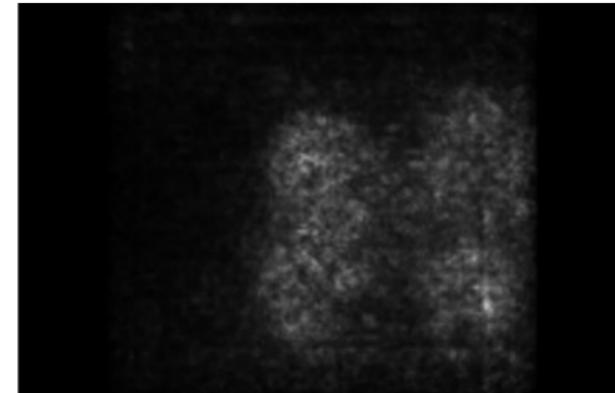


# Which pixels matter: Saliency via Backprop

Forward pass: Compute probabilities



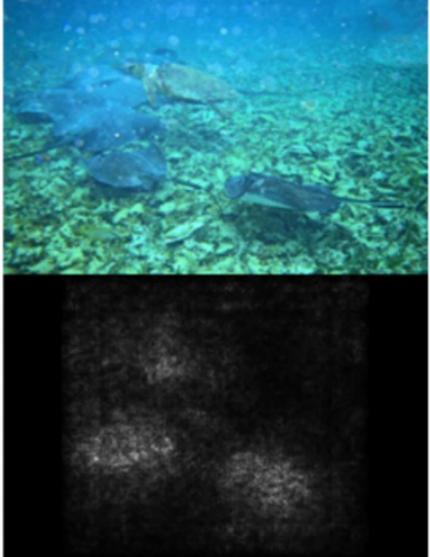
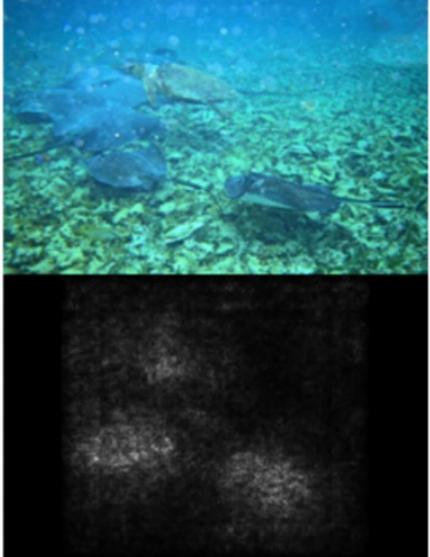
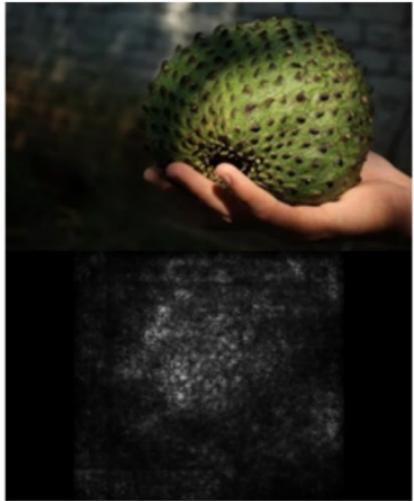
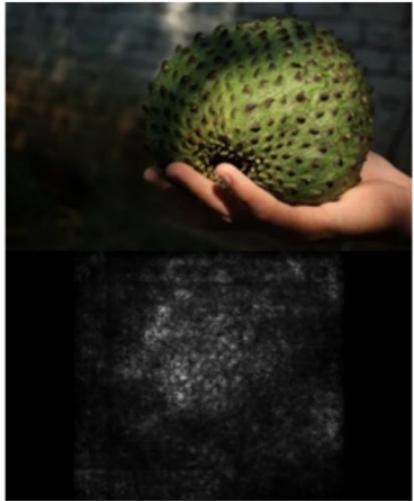
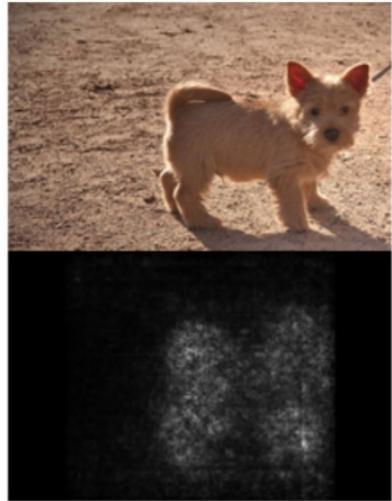
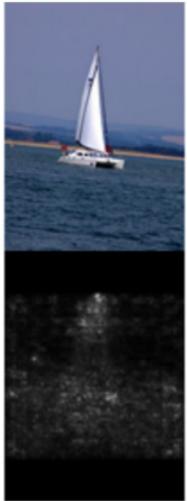
Compute gradient of (unnormalized) class score with respect to image pixels, take absolute value and max over RGB channels



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

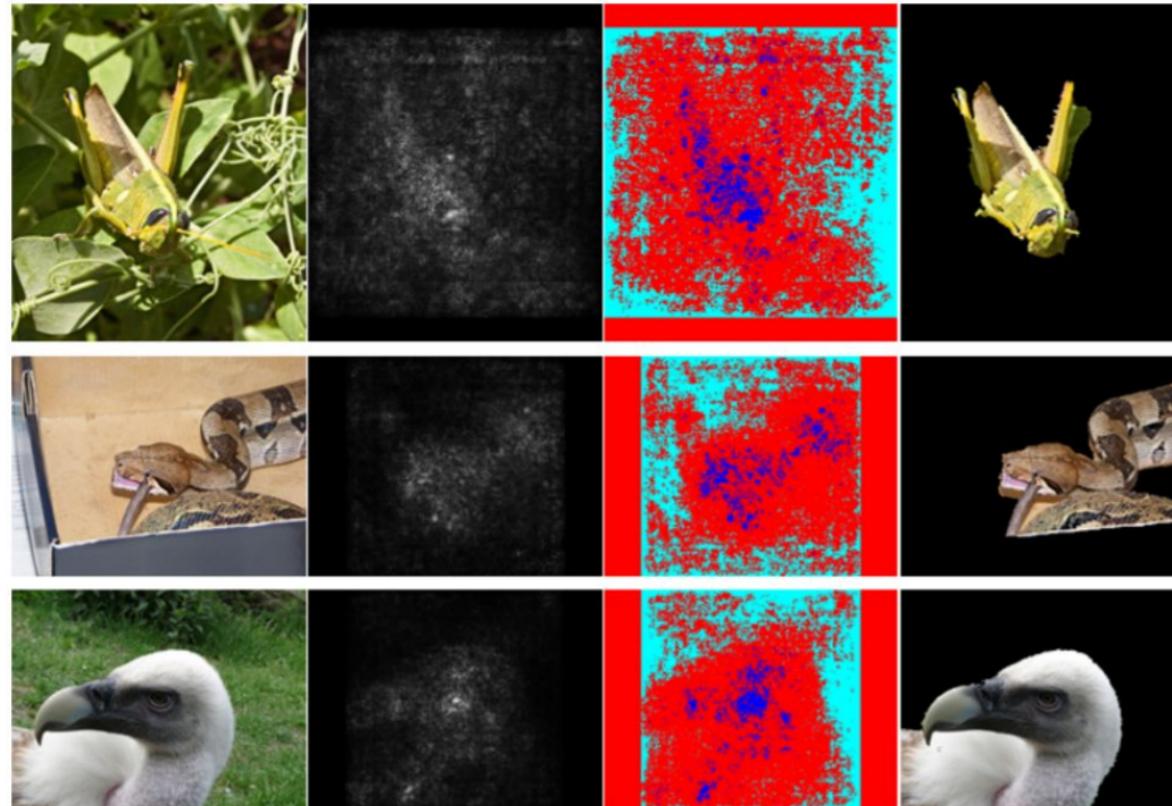
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

# Saliency Maps



Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.  
Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.

# Saliency Maps: Segmentation without supervision

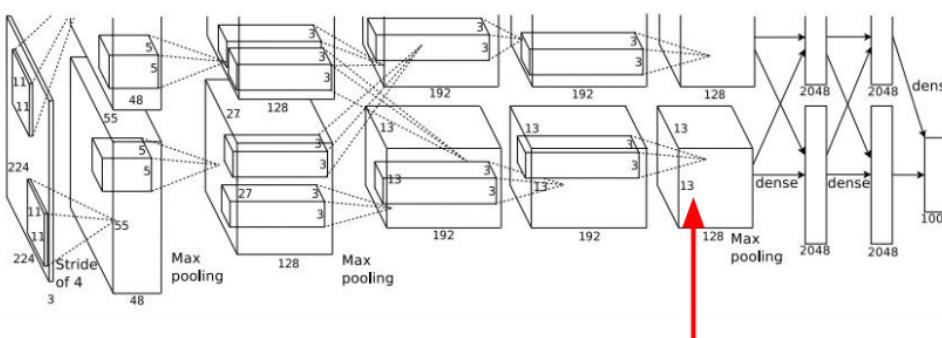


Use GrabCut on  
saliency map

Simonyan, Vedaldi, and Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Figures copyright Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman, 2014; reproduced with permission.  
Rother et al. "Grabcut: Interactive foreground extraction using iterated graph cuts". ACM TOG 2004

# Intermediate Features via (guided) backprop



Pick a single intermediate neuron, e.g. one value in  $128 \times 13 \times 13$  conv5 feature map

Compute gradient of neuron value with respect to image pixels

Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014  
Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015

# Intermediate features via (guided) backprop



Maximally activating patches  
(Each row is a different neuron)



Guided Backprop

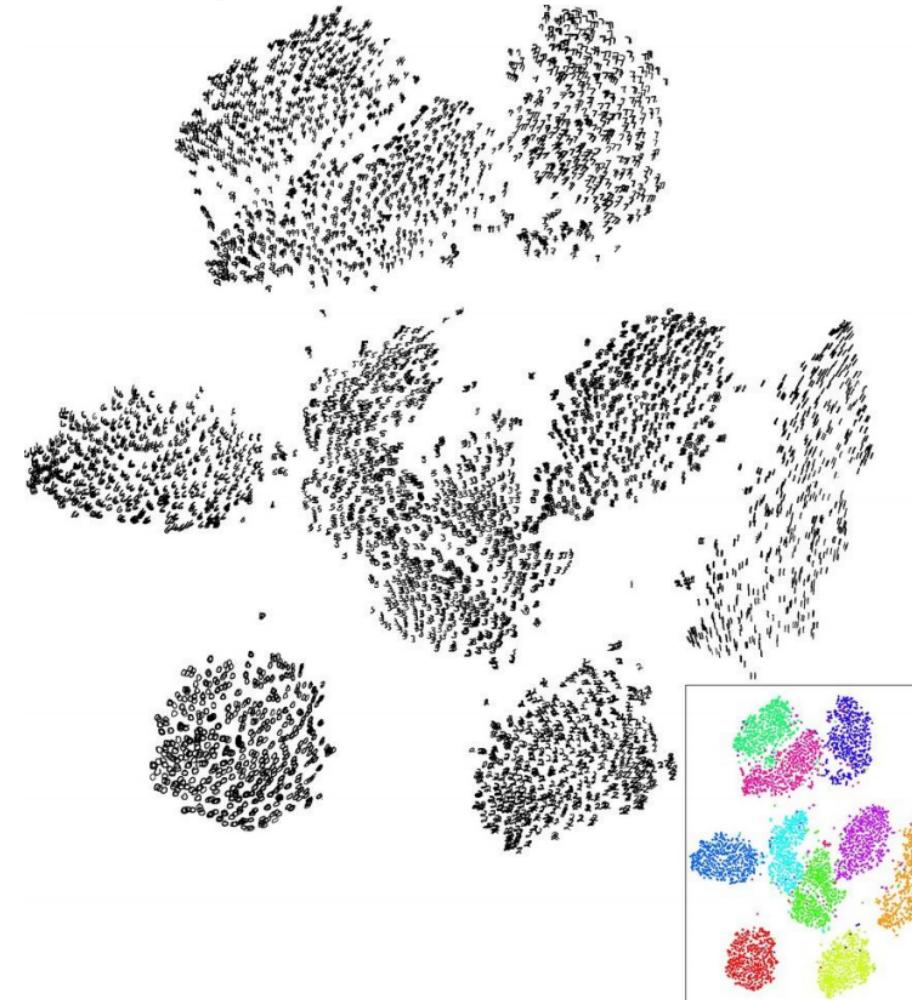
Zeiler and Fergus, "Visualizing and Understanding Convolutional Networks", ECCV 2014  
Springenberg et al, "Striving for Simplicity: The All Convolutional Net", ICLR Workshop 2015  
Figure copyright Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, Martin Riedmiller, 2015; reproduced with permission.

# Last Layer: Dimensionality Reduction

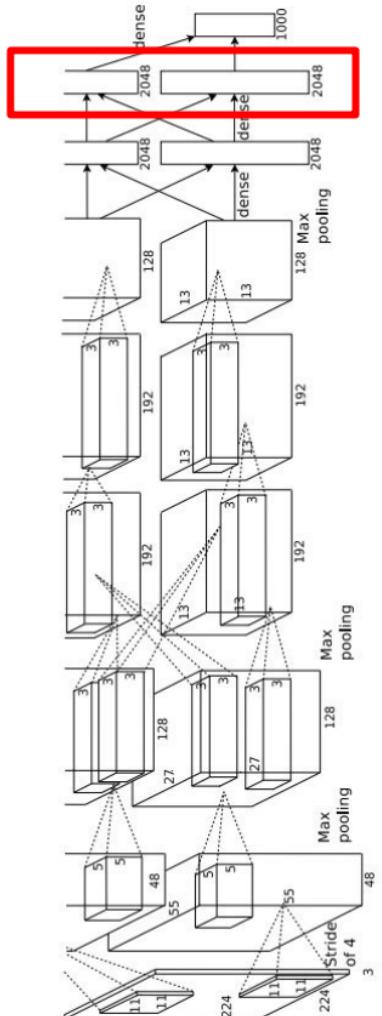
Visualize the “space” of FC7 feature vectors by reducing dimensionality of vectors from 4096 to 2 dimensions

Simple algorithm: Principal Component Analysis (PCA)

More complex: t-SNE



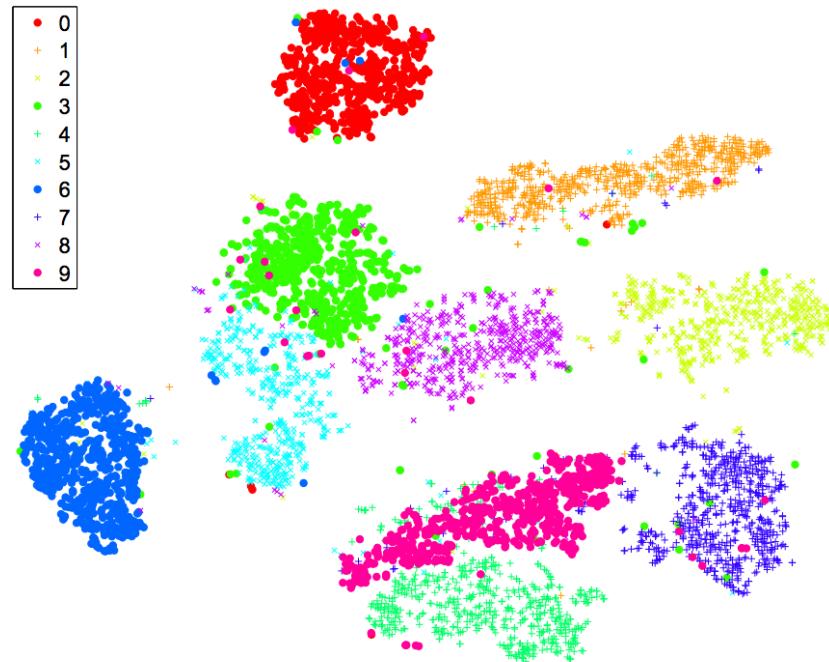
Van der Maaten and Hinton, “Visualizing Data using t-SNE”, JMLR 2008  
Figure copyright Laurens van der Maaten and Geoff Hinton, 2008. Reproduced with permission.



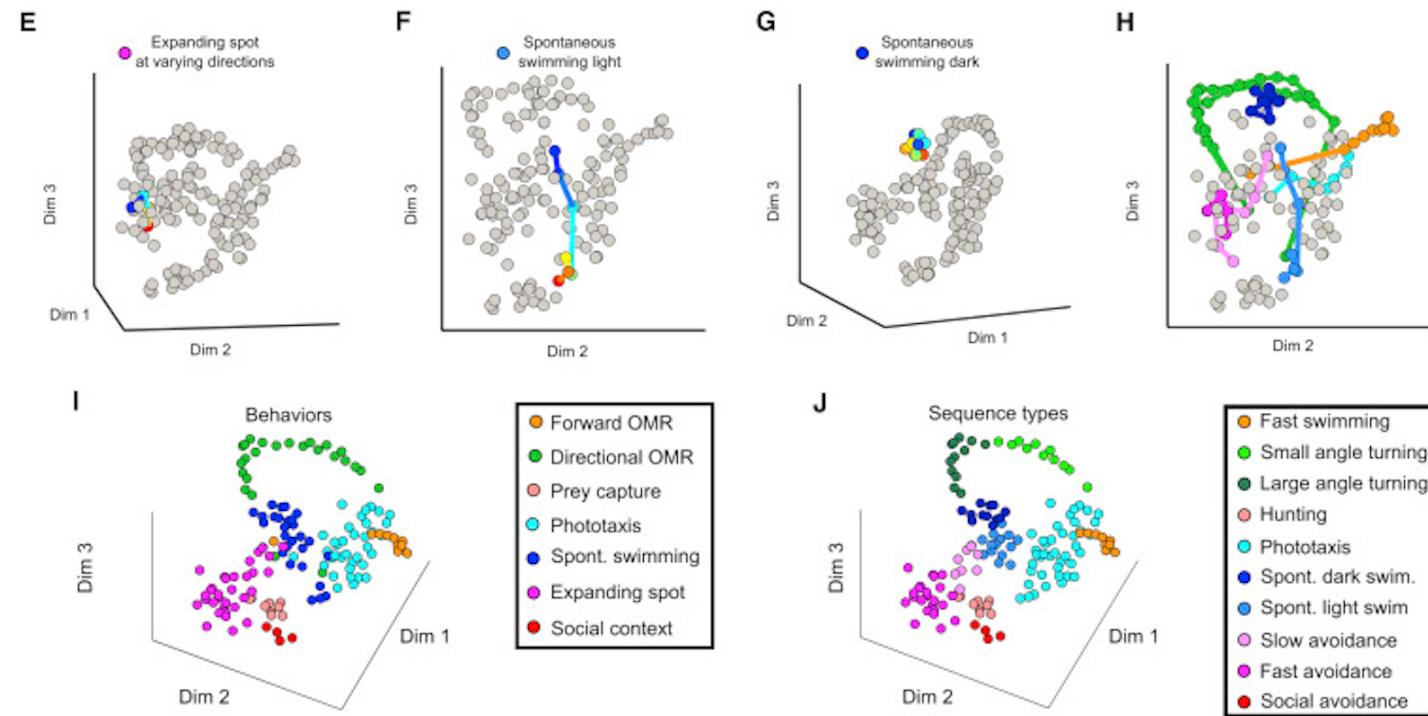
# TSNE for data visualization

- Reduce data dimensions to enable visualization in 2D or 3D
  - nD -> 2D or 3D
  - Preserve local structure of data to highlight groups
  - Unsupervised – clusters unlabeled data

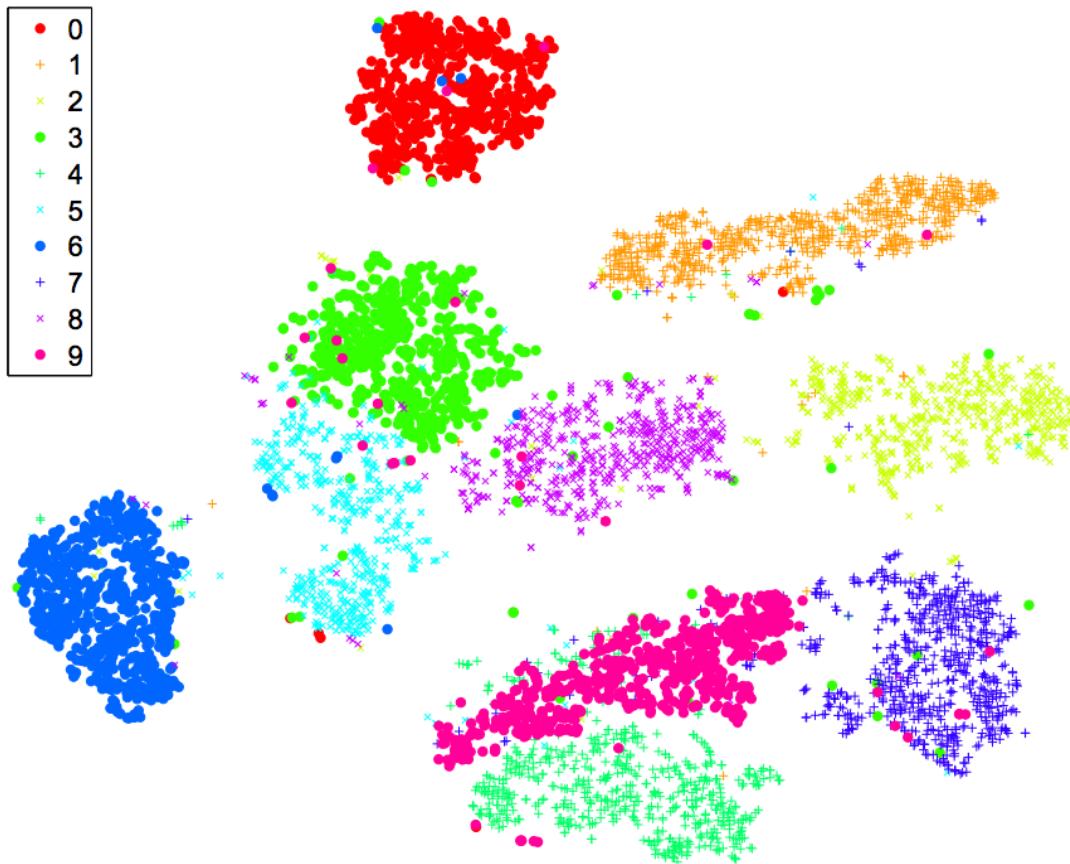
## Applied to MNIST digits



## Applied to movies of zebrafish behavior



## Aside about clustering data – why do we need deep learning at all?

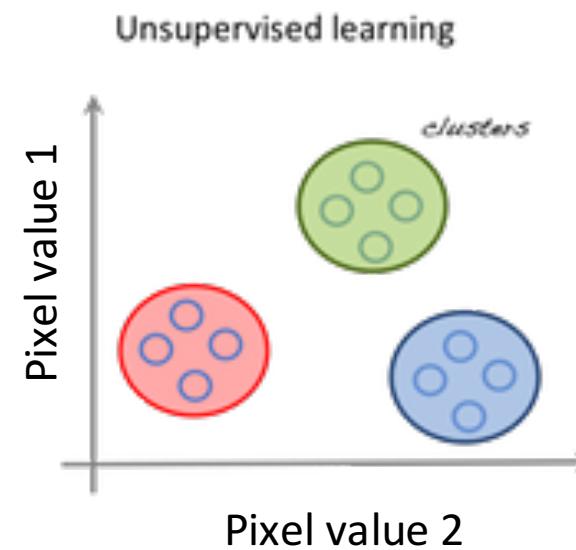
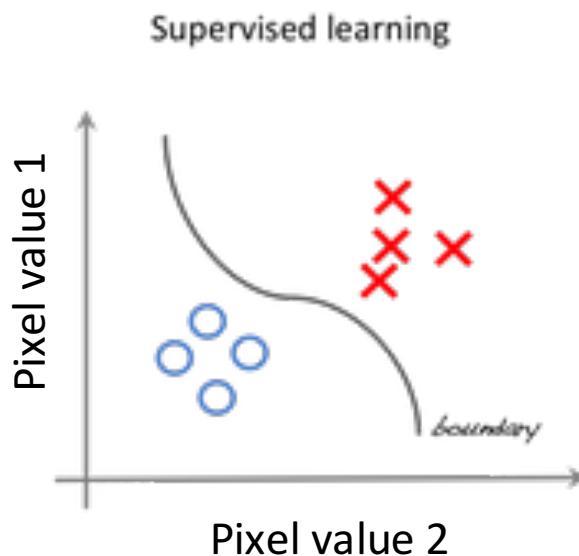


**Isn't this good enough?**

# Unsupervised learning in a nutshell

*Definition of Unsupervised Learning:*

Learning useful structure *without* labeled classes, optimization criterion, feedback signal, or any other information beyond the raw data

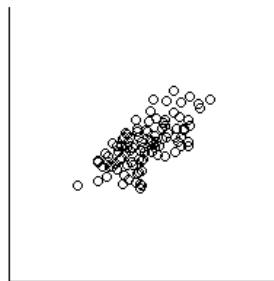


# Unsupervised learning in a nutshell

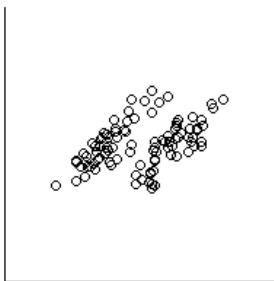
Mathematical tools for finding patterns in data:

- Eigenvector decomposition
- Principal component analysis
- Singular value decomposition

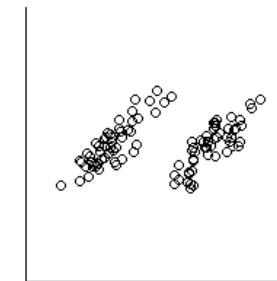
Dataset 1



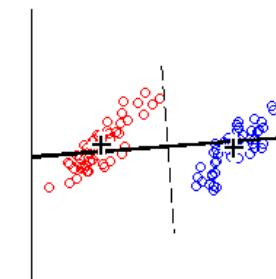
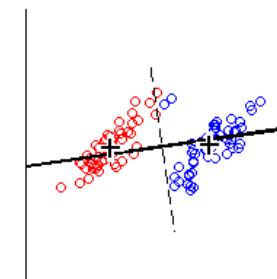
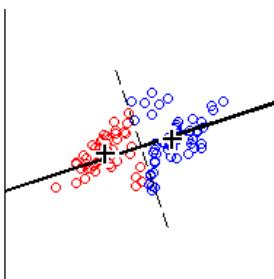
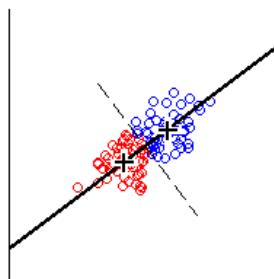
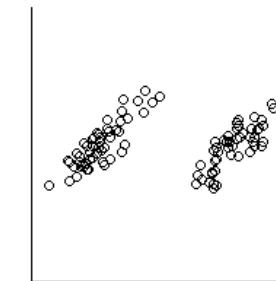
Dataset 2



Dataset 3

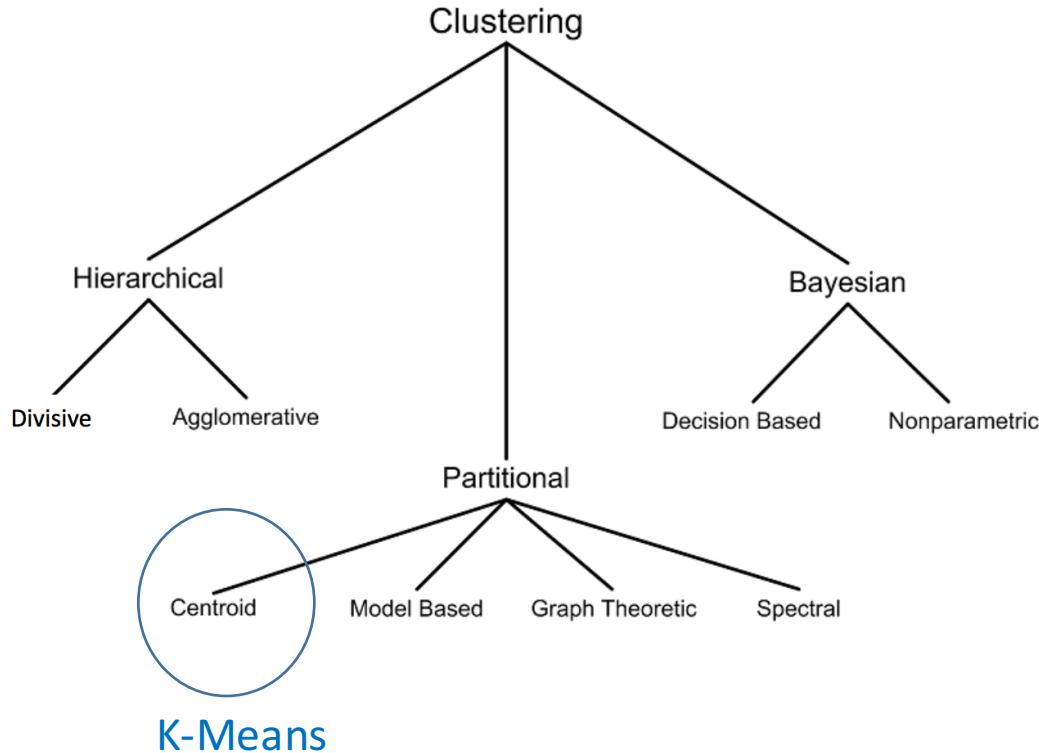


Dataset 4



# Iterative methods for unsupervised learning - Clustering

## Clustering techniques

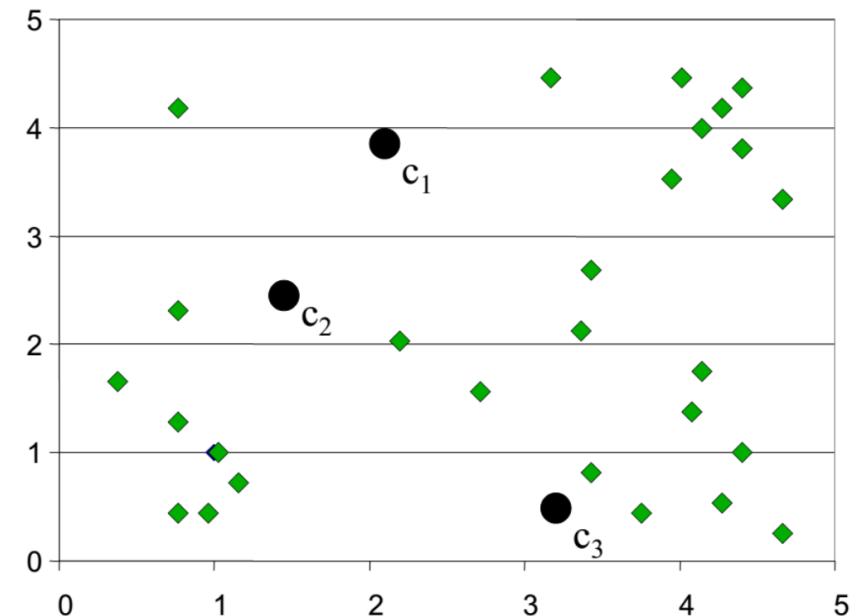


- **Hierarchical** algorithms find successive clusters using previously established clusters. These algorithms can be either **agglomerative** ("bottom-up") or **divisive** ("top-down"):
  - ➊ **Agglomerative algorithms** begin with each element as a separate cluster and merge them into successively larger clusters;
  - ➋ **Divisive algorithms** begin with the whole set and proceed to divide it into successively smaller clusters.
- **Partitional** algorithms typically determine all clusters at once, but can also be used as divisive algorithms in the hierarchical clustering.
- **Bayesian** algorithms try to generate a *posteriori distribution* over the collection of all partitions of the data.

# K-Means Clustering

- Given  $k$ , the  $k$ -means algorithm works as follows:
  - Choose  $k$  (random) data points (**seeds**) to be the initial **centroids**, cluster centers
  - Assign each data point to the closest **centroid**
  - Re-compute the **centroids** using the current cluster memberships
  - If a convergence criterion is not met, repeat steps 2 and 3

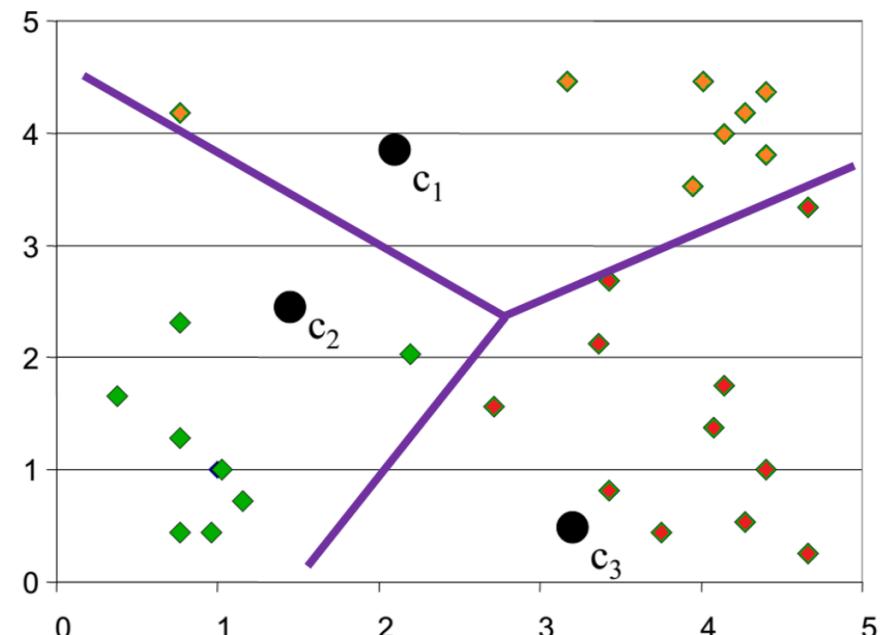
Randomly initialize seeds



# K-Means Clustering

- Given  $k$ , the  $k$ -means algorithm works as follows:
  - Choose  $k$  (random) data points (**seeds**) to be the initial **centroids**, cluster centers
  - Assign each data point to the closest **centroid**
  - Re-compute the **centroids** using the current cluster memberships
  - If a convergence criterion is not met, repeat steps 2 and 3

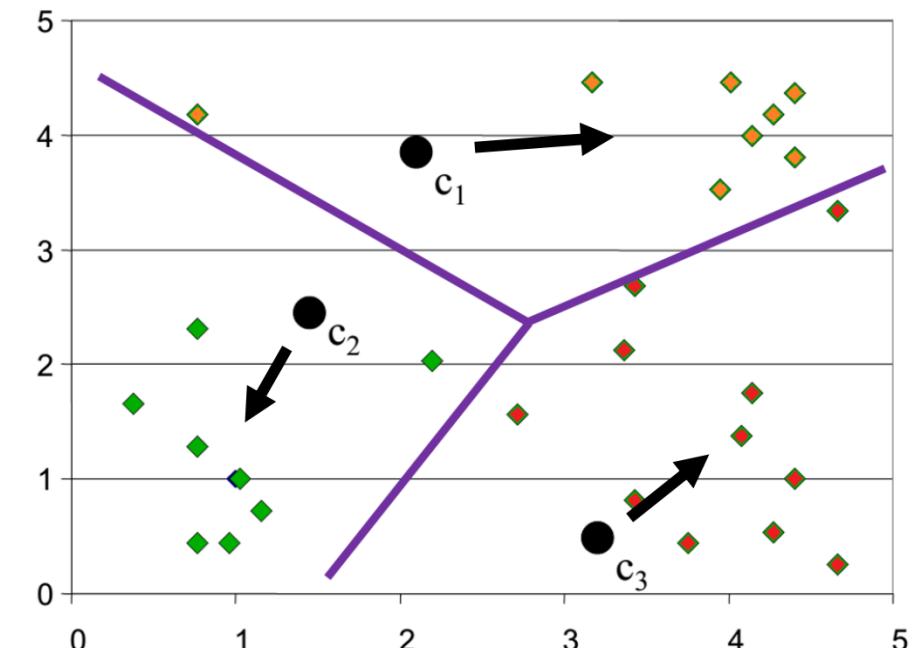
Determine cluster membership for each data point



# K-Means Clustering

- Given  $k$ , the  $k$ -means algorithm works as follows:
  - Choose  $k$  (random) data points (**seeds**) to be the initial **centroids**, cluster centers
  - Assign each data point to the closest **centroid**
  - Re-compute the **centroids** using the current cluster memberships
  - If a convergence criterion is not met, repeat steps 2 and 3

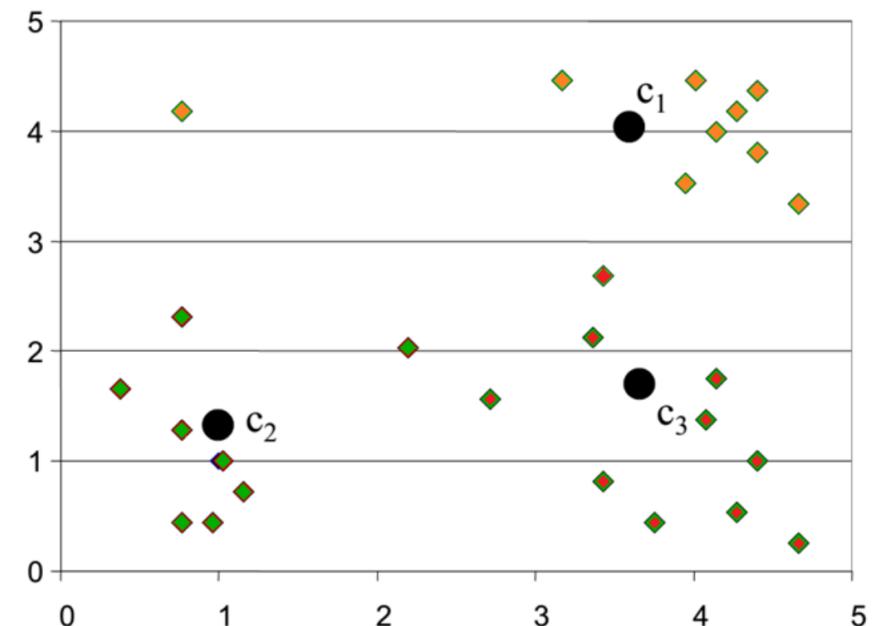
Compute and update new cluster center



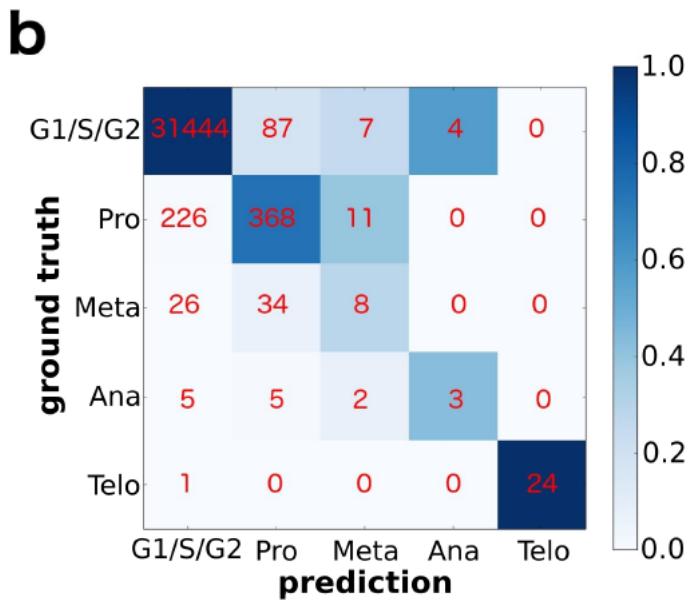
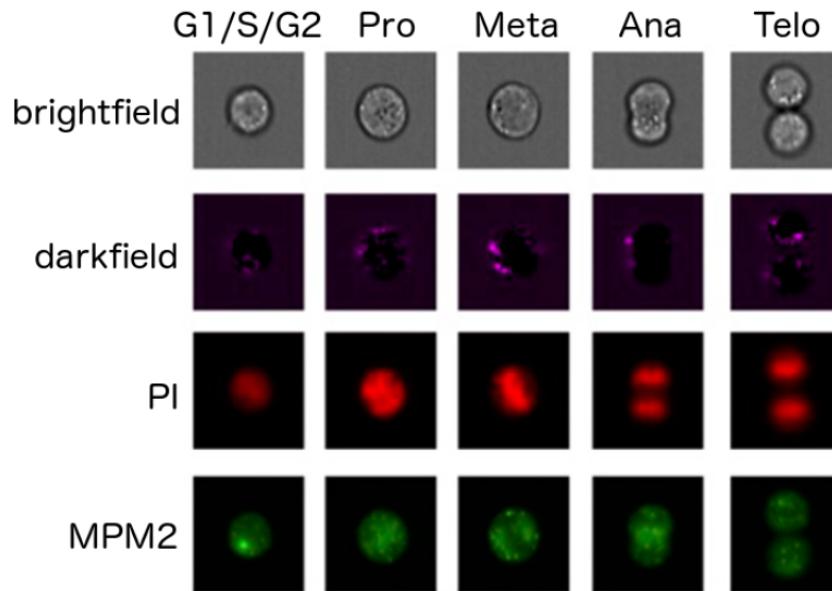
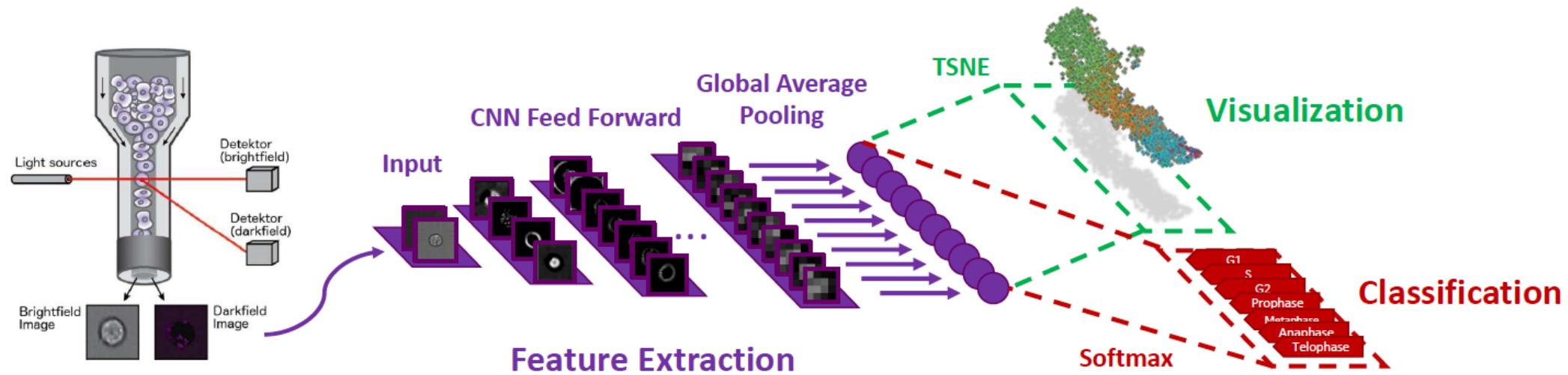
# K-Means Clustering

- Given  $k$ , the  $k$ -means algorithm works as follows:
  - Choose  $k$  (random) data points (**seeds**) to be the initial **centroids**, cluster centers
  - Assign each data point to the closest **centroid**
  - Re-compute the **centroids** using the current cluster memberships
  - If a convergence criterion is not met, repeat steps 2 and 3

Result of first iteration

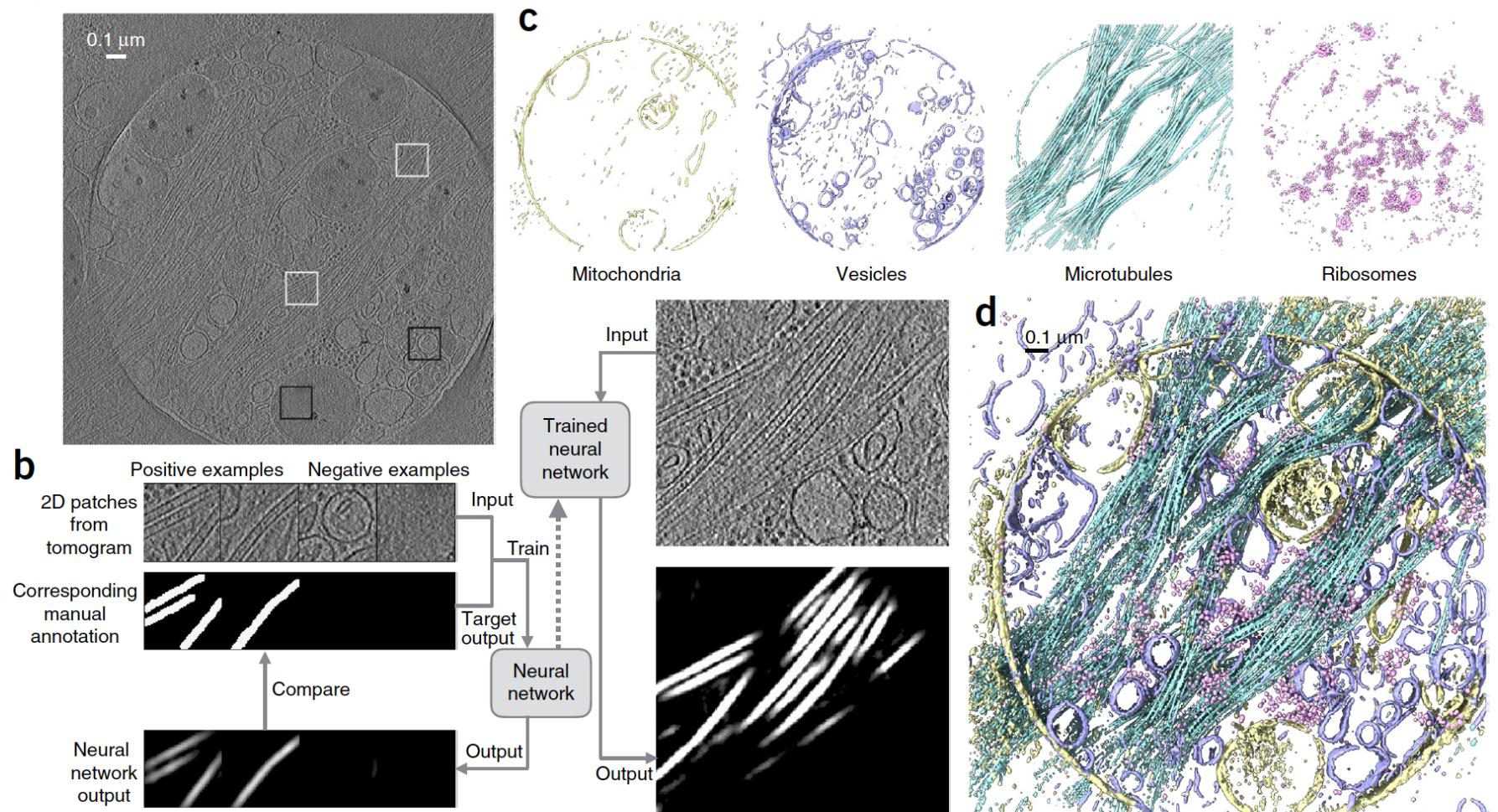


## **Examples of CNN's for biomedical image classification**



# Convolutional neural networks for automated annotation of cellular cryo-electron tomograms

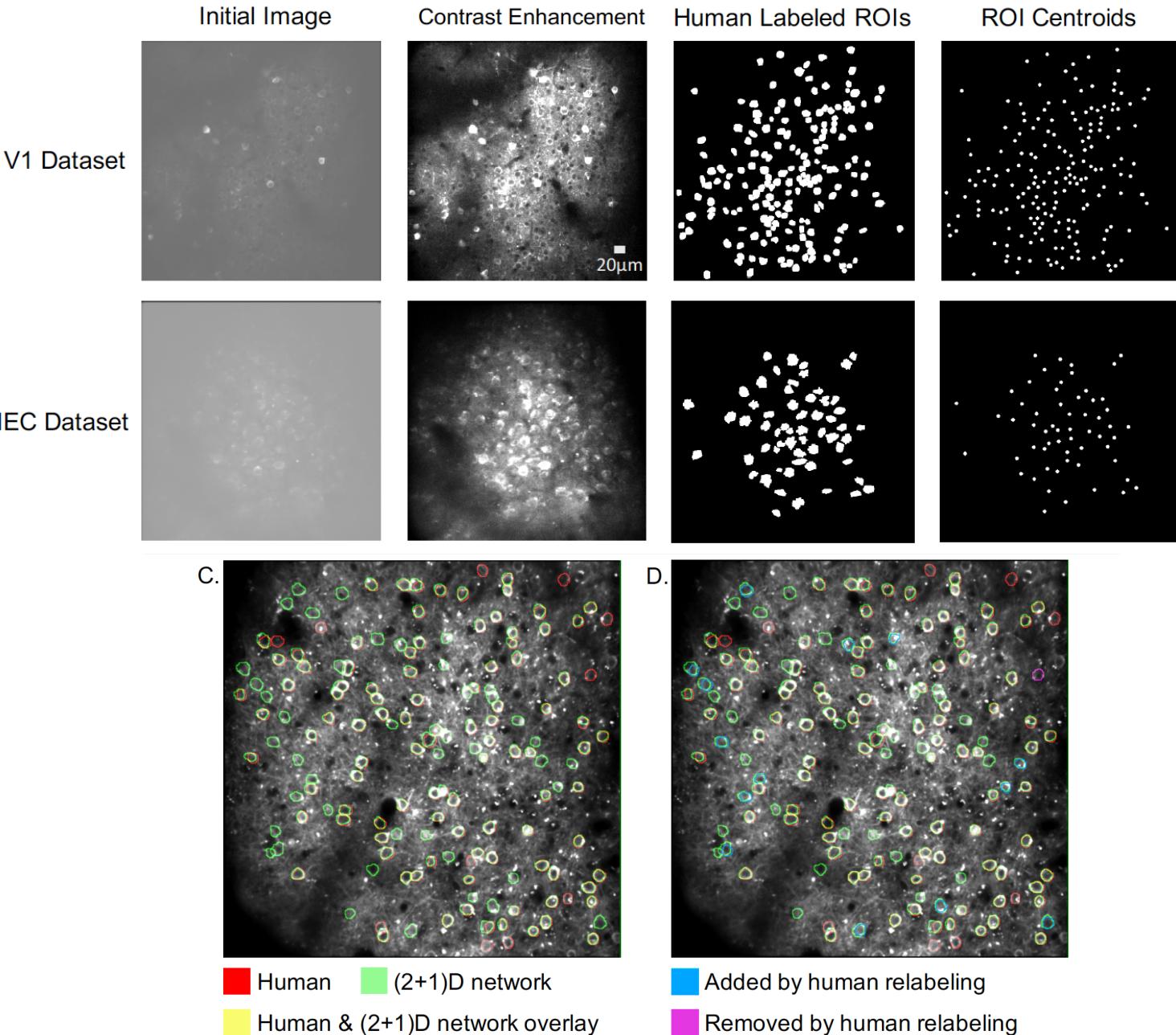
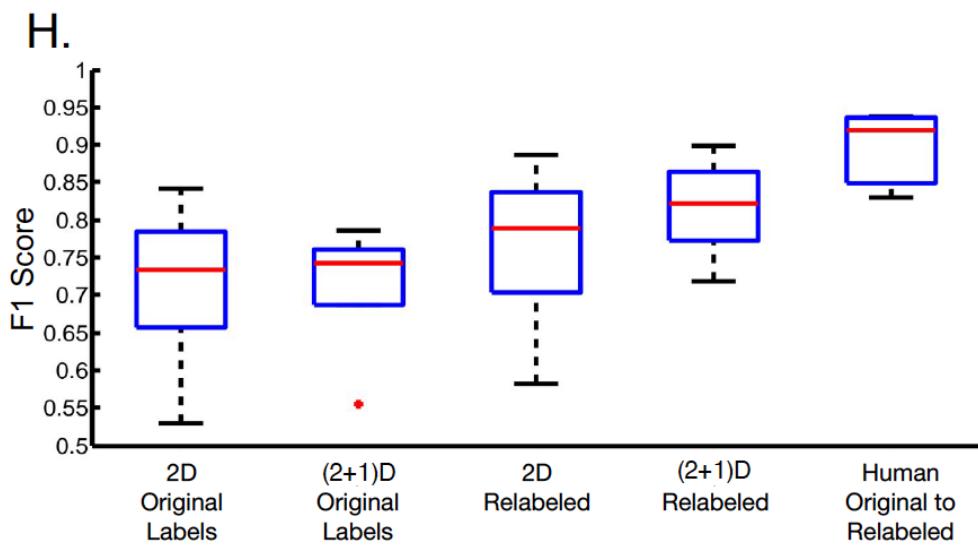
Muyuan Chen<sup>1,2</sup>, Wei Dai<sup>2,4</sup>, Stella Y Sun<sup>2</sup>,  
Darius Jonasch<sup>2</sup>, Cynthia Y He<sup>3</sup>, Michael F Schmid<sup>2</sup>,  
Wah Chiu<sup>2</sup> & Steven J Ludtke<sup>2</sup> 



# Automatic Neuron Detection in Calcium Imaging Data Using Convolutional Networks

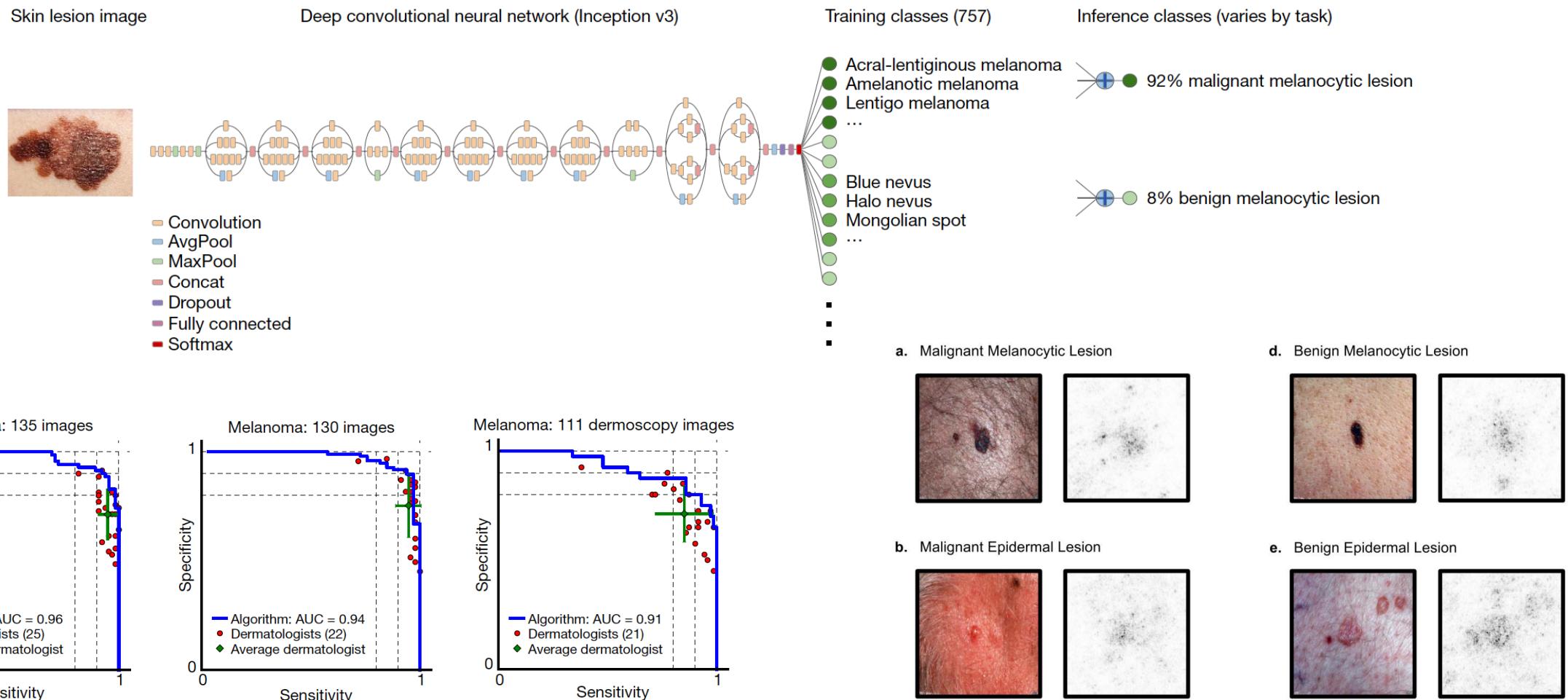
Noah J. Apthorpe<sup>1\*</sup> Alexander J. Riordan<sup>2\*</sup> Rob E. Aguilar,<sup>1</sup> Jan Homann<sup>2</sup>  
Yi Gu<sup>2</sup> David W. Tank<sup>2</sup> H. Sebastian Seung<sup>1,2</sup>

<sup>1</sup>Computer Science Department <sup>2</sup>Princeton Neuroscience Institute  
Princeton University



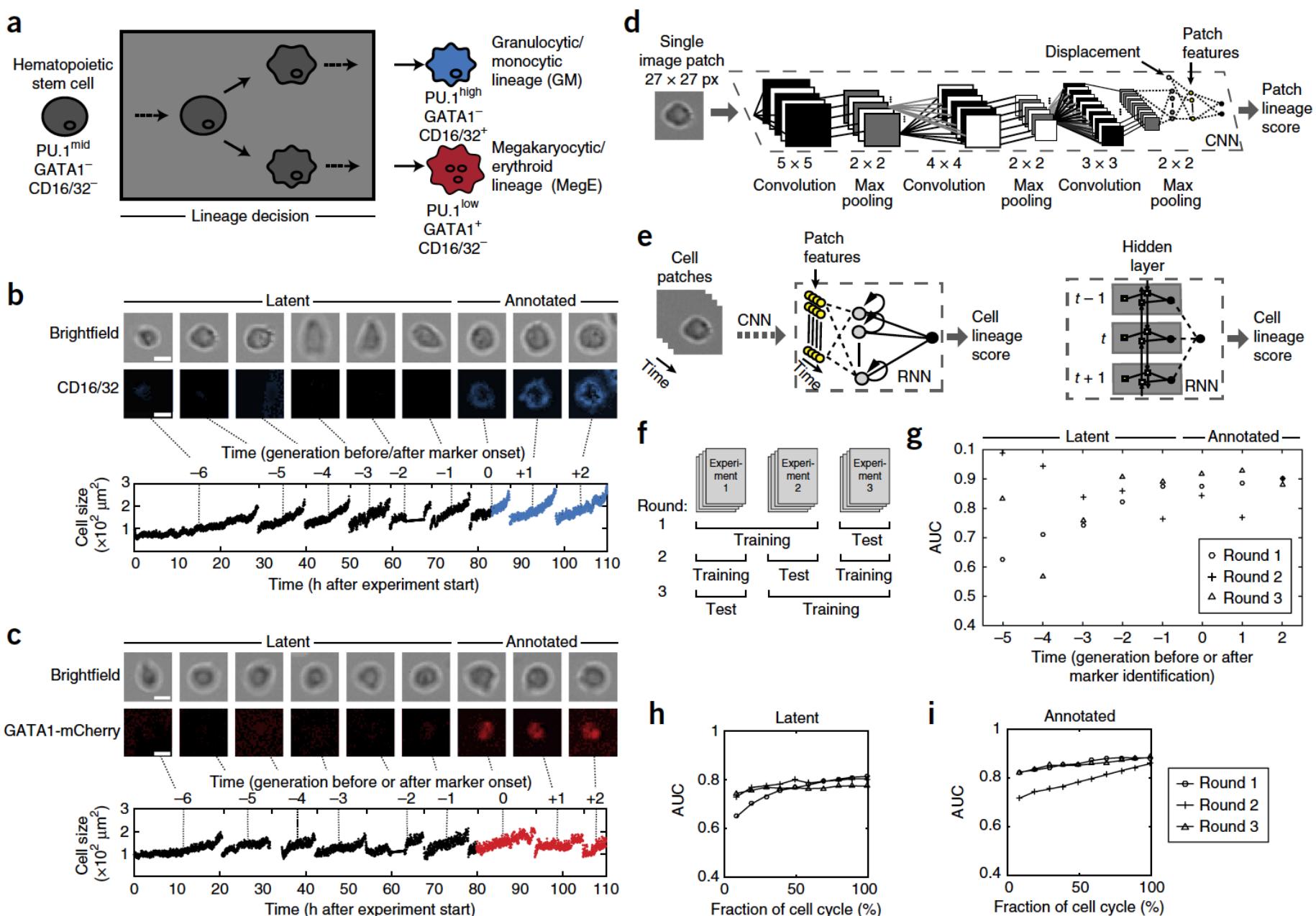
# Dermatologist-level classification of skin cancer with deep neural networks

Andre Esteva<sup>1\*</sup>, Brett Kuprel<sup>1\*</sup>, Roberto A. Novoa<sup>2,3</sup>, Justin Ko<sup>2</sup>, Susan M. Swetter<sup>2,4</sup>, Helen M. Blau<sup>5</sup> & Sebastian Thrun<sup>6</sup>



# Prospective identification of hematopoietic lineage choice by deep learning

Felix Buggenthin<sup>1,6</sup>, Florian Buettner<sup>1,2,6</sup>, Philipp S Hoppe<sup>3,4</sup>, Max Endele<sup>3</sup>, Manuel Kroiss<sup>1,5</sup>, Michael Strasser<sup>1</sup>, Michael Schwarzfischer<sup>1</sup>, Dirk Loeffler<sup>3,4</sup>, Konstantinos D Kokkaliaris<sup>3,4</sup>, Oliver Hilsenbeck<sup>3,4</sup>, Timm Schroeder<sup>3,4</sup>, Fabian J Theis<sup>1,5</sup> & Carsten Marr<sup>1</sup>



## **A good time to talk about the class project!**

- Select a “base” dataset
- Simulate parameters of a physical (imaging) system with base dataset
- Train deep neural net with simulated dataset
- Report results

What you'll need to submit:

- 1) The project's source code
- 2) A short research-style paper (3 pages minimum, 5 pages maximum) that includes an introduction, results, a discussion section, references and at least 2 figures
- 3) A completed web template containing the main results from the research paper
- 4) An 8-minute presentation that each student will deliver to the class

## Example project topics:

Can we design a new lens/transducer/antenna shape to improve classification of X?

What is the tradeoff between image resolution and classification accuracy for X?

Can we determine an optimal set of colors to improve fluorophore distinguishability?

If we capture 2 images that are overlapped on one sensor, what is the best way to pre-blur them to then be able to tell them apart? Or to be able to classify them together?

If we just had a few sensors, how should we arrange them e.g. a mask to be able to predict the position of X?

Is there some optimal shift-variant blur that we can use for a particular task?

Or, given a shift-variant PSF image, can we establish a good deconvolution using locally connected layers?

What is the optimal way to layout filters on a sensor to capture a color image for classification? Or an HDR image?

HDR image generation with filters over pixels – what is optimal design?

What if we could make a sensor with different sized pixels – how should they be laid out to achieve the best X?