

# Supervised Pretraining

Jia-Bin Huang

Virginia Tech

ECE 6554 Advanced Computer Vision

# Administrative stuffs

- Project proposal due March 2<sup>nd</sup>
  - 1-page summary of
- Feedback on paper summary
  - Explicit structure

# Discussion – Think-pair-share

- Discuss
  - strength,
  - weakness, and
  - potential extension
- Share with class

# Today's class

- Training tricks for CNN
- Transfer learning via supervised pretraining

# Training CNN with gradient descent

- A CNN as composition of functions

$$f_{\mathbf{w}}(\mathbf{x}) = f_L(\dots (f_2(f_1(\mathbf{x}; \mathbf{w}_1); \mathbf{w}_2) \dots); \mathbf{w}_L)$$

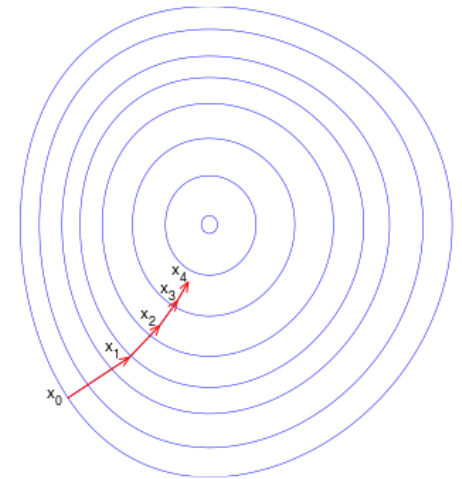
- Parameters

$$\mathbf{w} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_L)$$

- Empirical loss function

$$L(\mathbf{w}) = \frac{1}{n} \sum_i l(z_i, f_{\mathbf{w}}(\mathbf{x}_i))$$

- Gradient descent



New weight

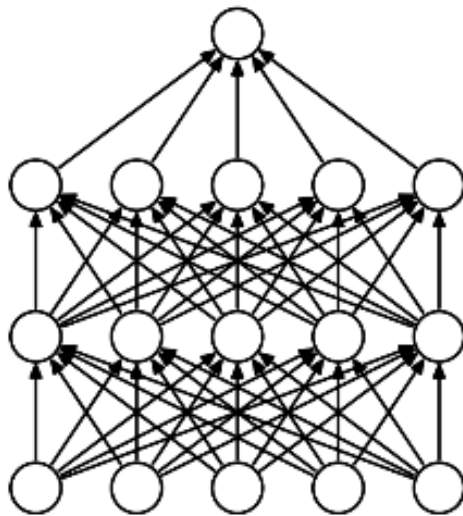
$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta_t \frac{\partial f}{\partial \mathbf{w}}(\mathbf{w}^t)$$

Old weight

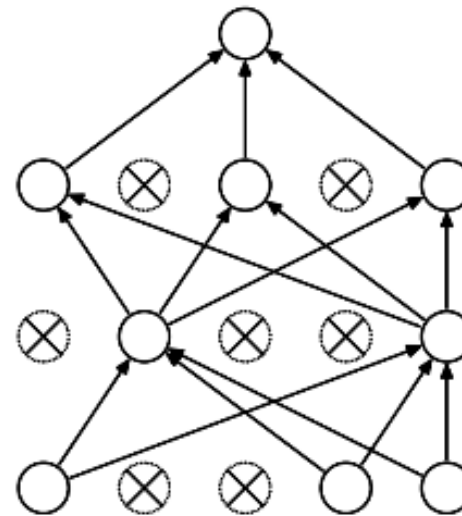
Learning rate

Gradient

# Dropout



(a) Standard Neural Net



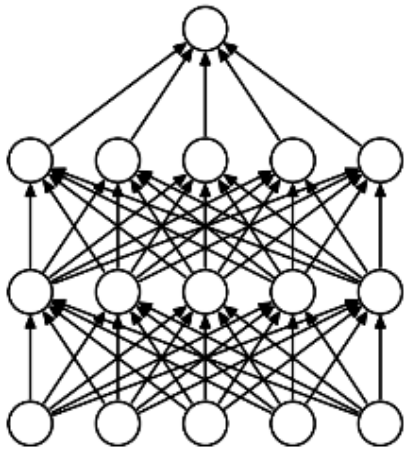
(b) After applying dropout.

Intuition: successful conspiracies

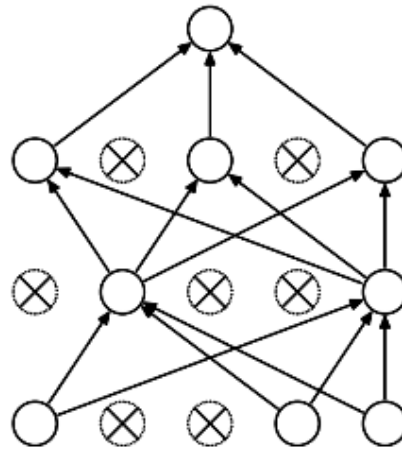
- 50 people planning a conspiracy
- Strategy A: plan a big conspiracy involving 50 people
  - Likely to fail. 50 people need to play their parts correctly.
- Strategy B: plan 10 conspiracies each involving 5 people
  - Likely to succeed!

Dropout: A simple way to prevent neural networks from overfitting [[Srivastava JMLR 2014](#)]

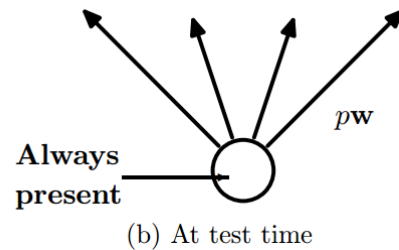
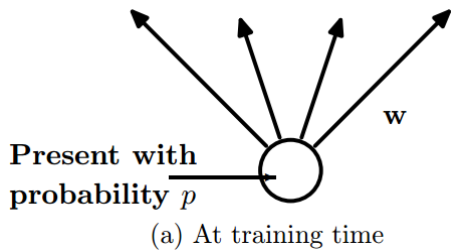
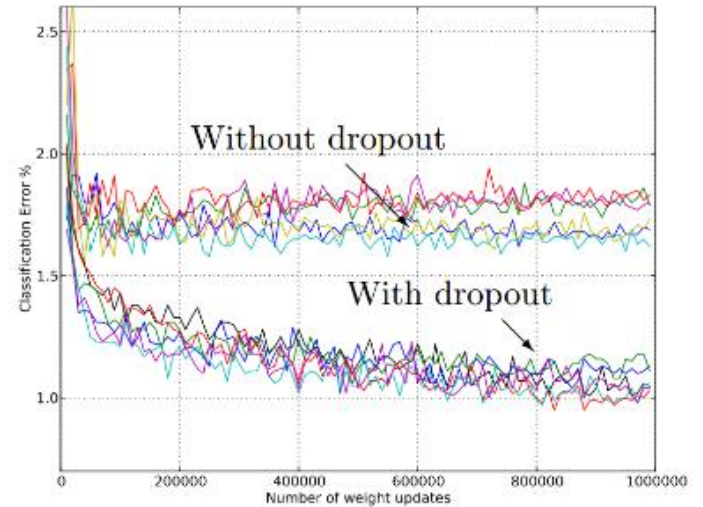
# Dropout



(a) Standard Neural Net



(b) After applying dropout.



**Main Idea:** approximately combining exponentially many different neural network architectures efficiently

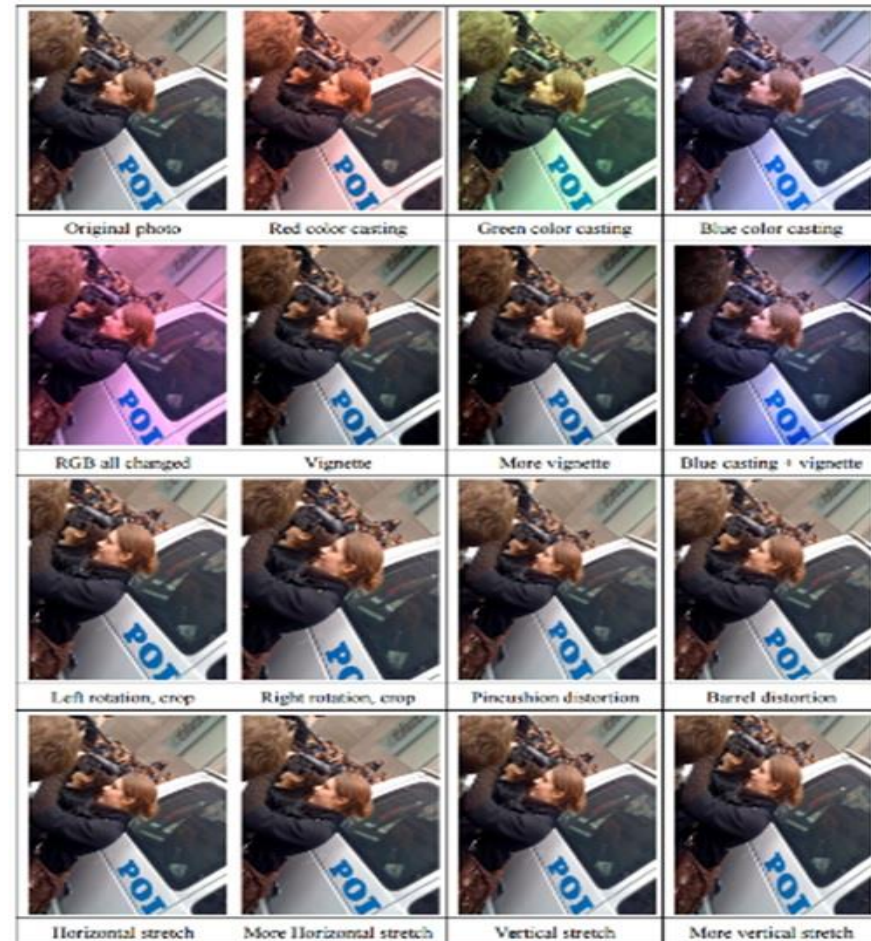
Model	Top-1 (val)	Top-5 (val)	Top-5 (test)
SVM on Fisher Vectors of Dense SIFT and Color Statistics	-	-	27.3
Avg of classifiers over FVs of SIFT, LBP, GIST and CSIFT	-	-	26.2
Conv Net + dropout (Krizhevsky et al., 2012)	40.7	18.2	-
Avg of 5 Conv Nets + dropout (Krizhevsky et al., 2012)	38.1	16.4	16.4

Table 6: Results on the ILSVRC-2012 validation/test set.

Dropout: A simple way to prevent neural networks from overfitting [[Srivastava JMLR 2014](#)]

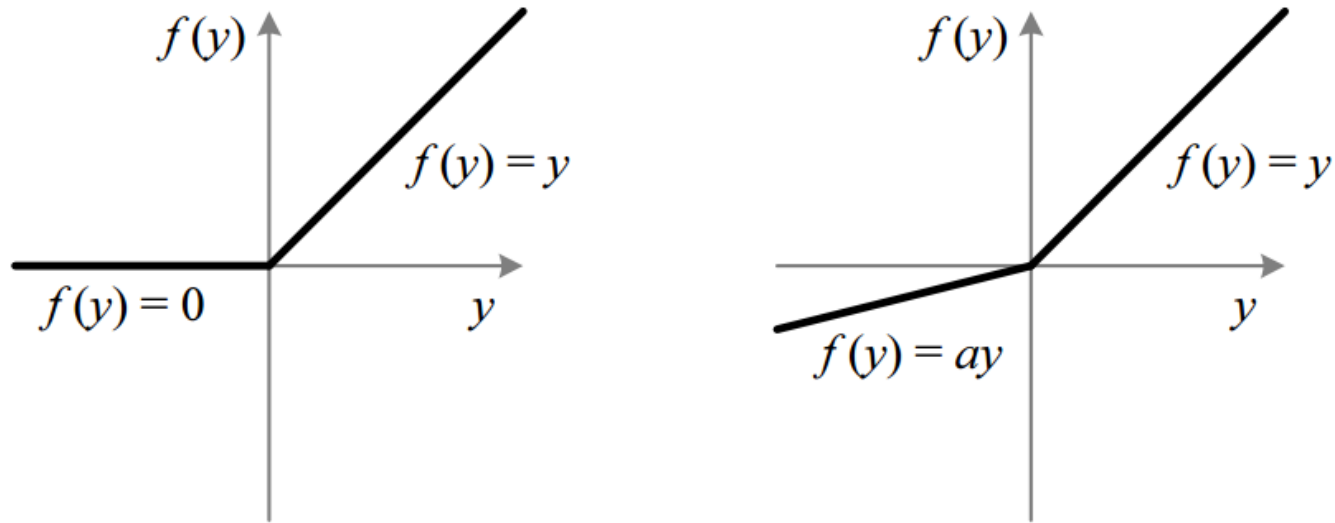
# Data Augmentation (Jittering)

- Create *virtual* training samples
  - Horizontal flip
  - Random crop
  - Color casting
  - Geometric distortion





# Parametric Rectified Linear Unit



	team	top-5 (test)
in competition ILSVRC 14	MSRA, SPP-nets [11]	8.06
	VGG [25]	7.32
	GoogLeNet [29]	6.66
post-competition	VGG [25] (arXiv v5)	6.8
	Baidu [32]	5.98
	<b>MSRA, PReLU-nets</b>	<b>4.94</b>

Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification [[He et al. 2015](#)]

# Batch Normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots x_m\}$ ;

Parameters to be learned:  $\gamma, \beta$

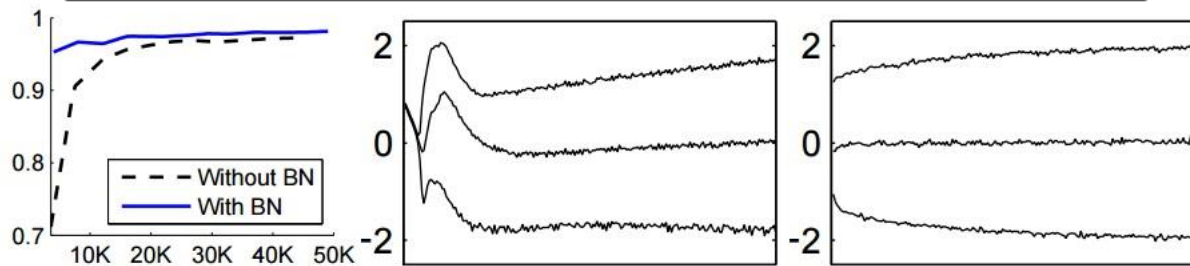
**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$



(a)

(b) Without BN

(c) With BN

Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift [[Ioffe and Szegedy 2015](#)]

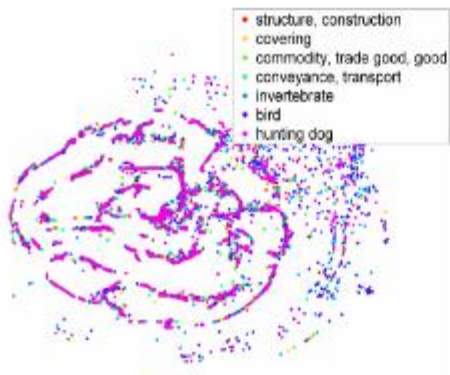
# Transfer Learning

- Improvement of learning in a **new** task through the *transfer of knowledge* from a **related** task that has already been learned.
- Weight initialization for CNN
- Two major strategies
  - ConvNet as fixed feature extractor
  - Fine-tuning the ConvNet

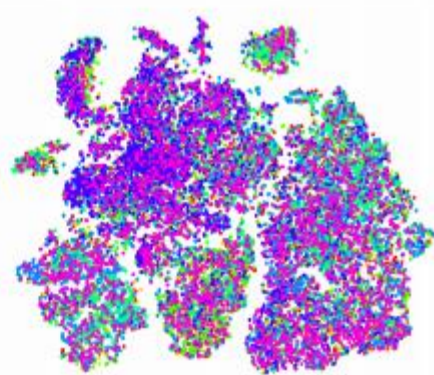
# When to finetune your model?

- New dataset is small and similar to original dataset.
  - train a linear classifier on the CNN codes
- New dataset is large and similar to the original dataset
  - fine-tune through the full network
- New dataset is small but very different from the original dataset
  - SVM classifier from activations somewhere earlier in the network
- New dataset is large and very different from the original dataset
  - fine-tune through the entire network

# Convolutional activation features



(a) LLC



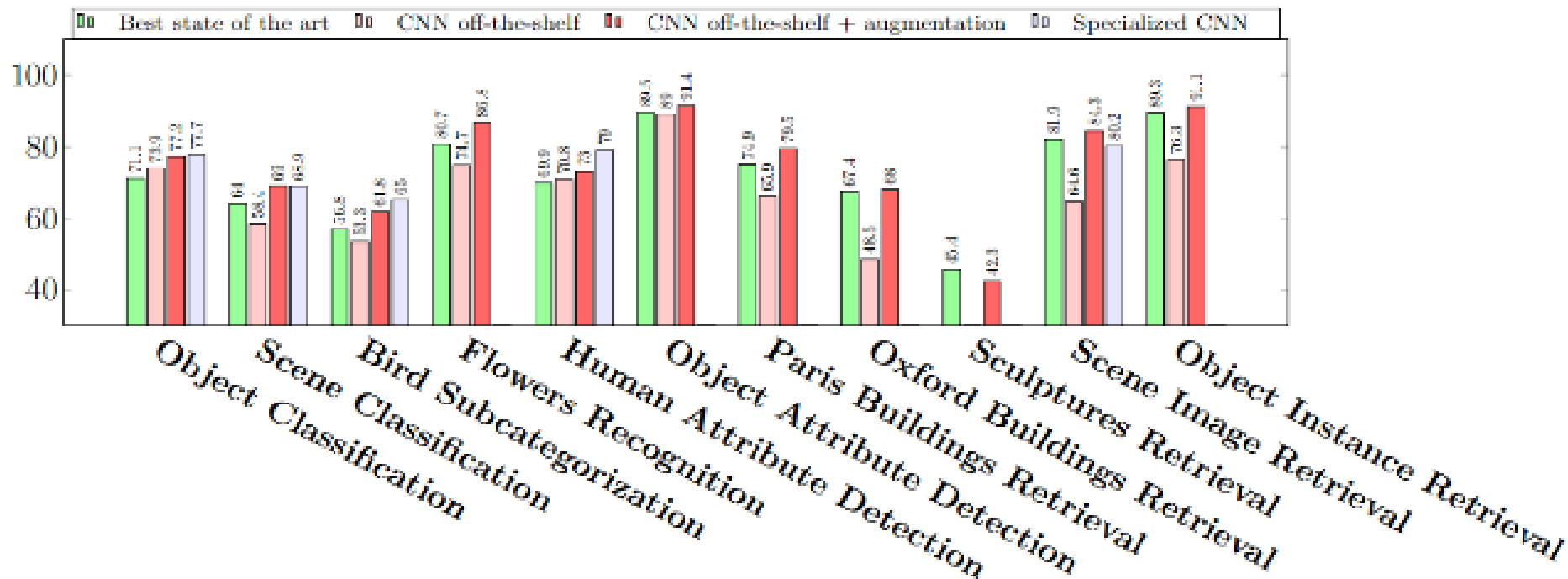
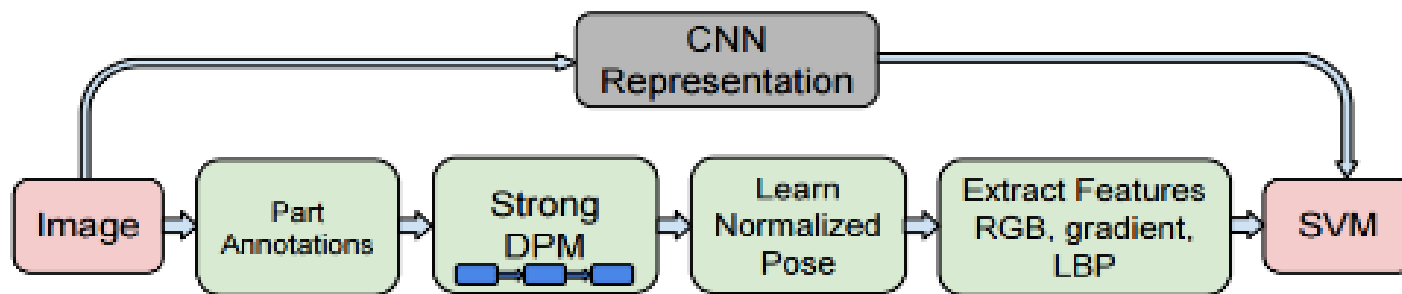
(b) GIST



(c) DeCAF<sub>1</sub>

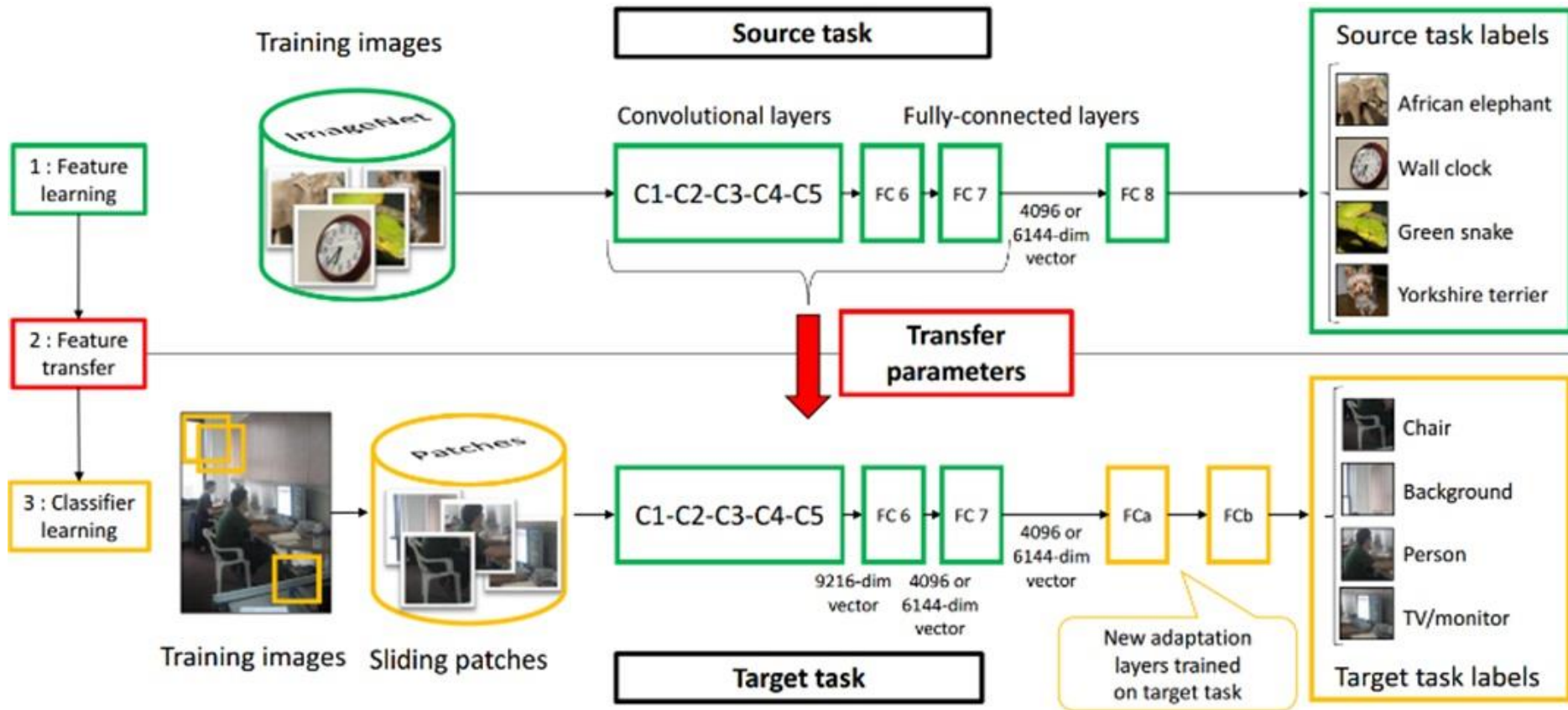


(d) DeCAF<sub>6</sub>



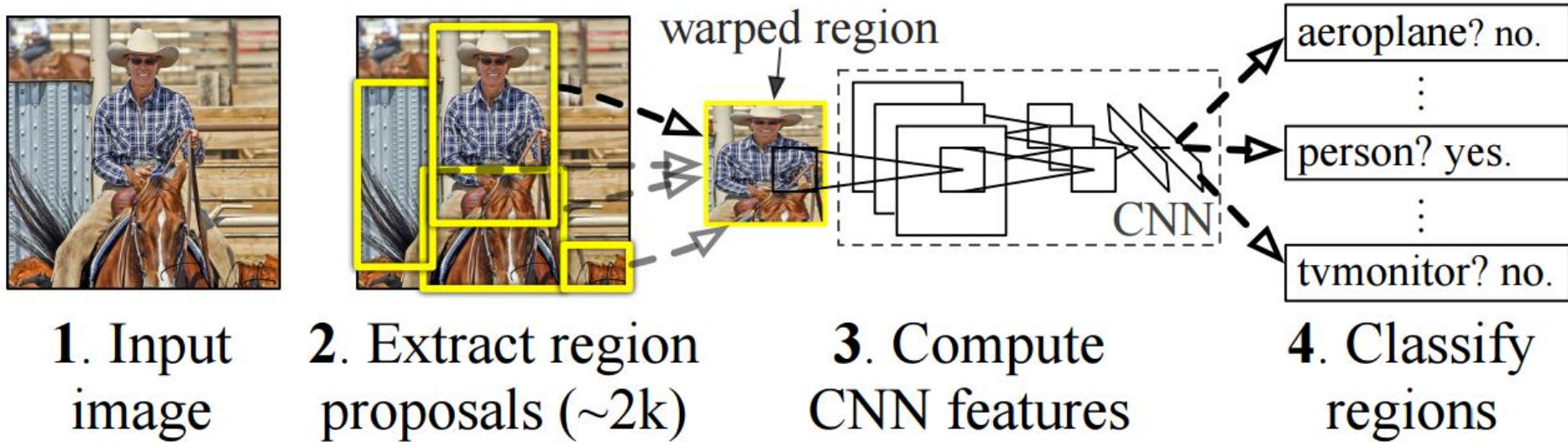
CNN Features off-the-shelf: an Astounding Baseline for Recognition

[[Razavian et al. 2014](#)]



Learning and Transferring Mid-Level Image Representations using Convolutional Neural Networks [[Oquab et al. CVPR 2014](#)]

# R-CNN: *Regions with CNN features*

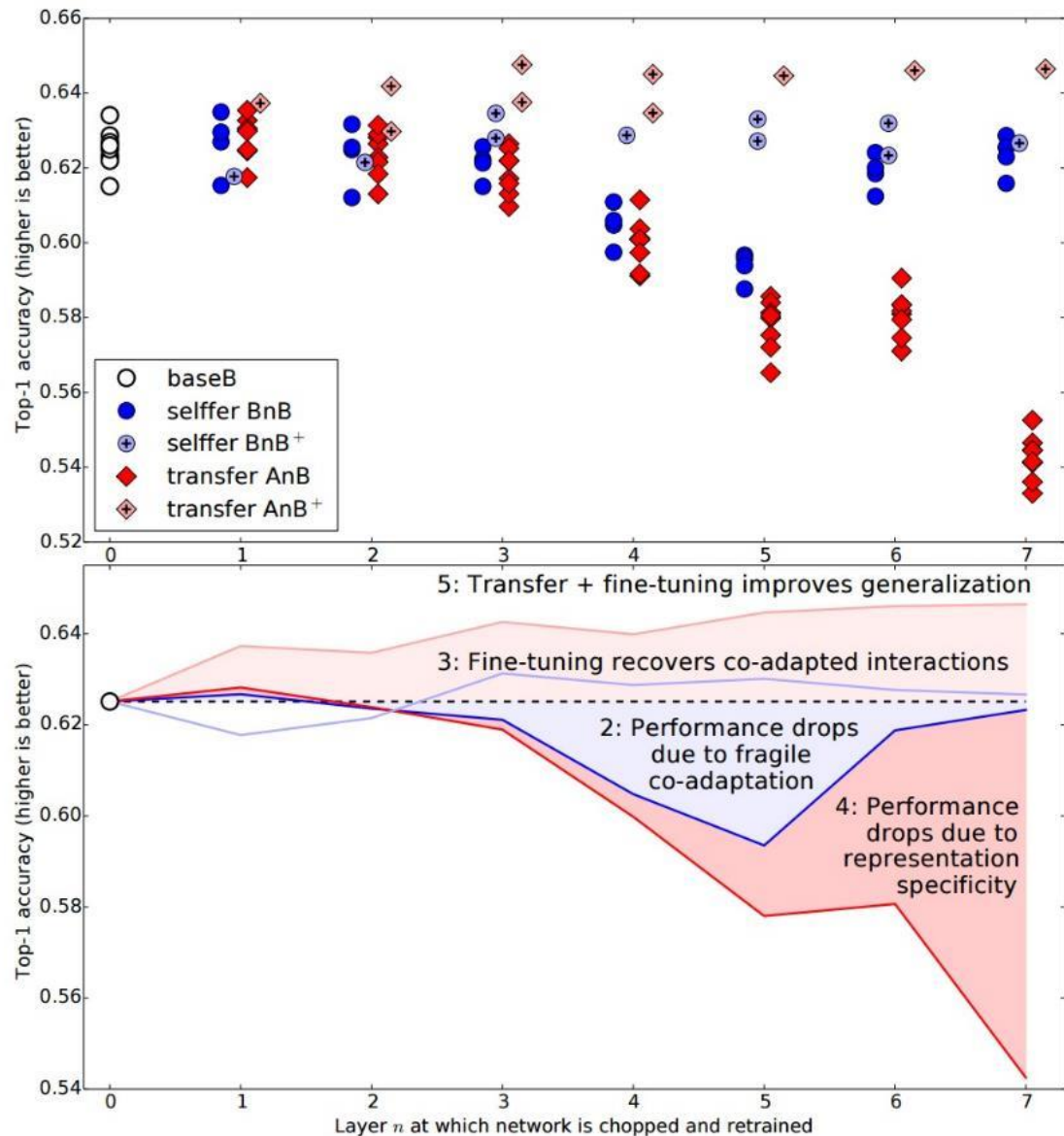
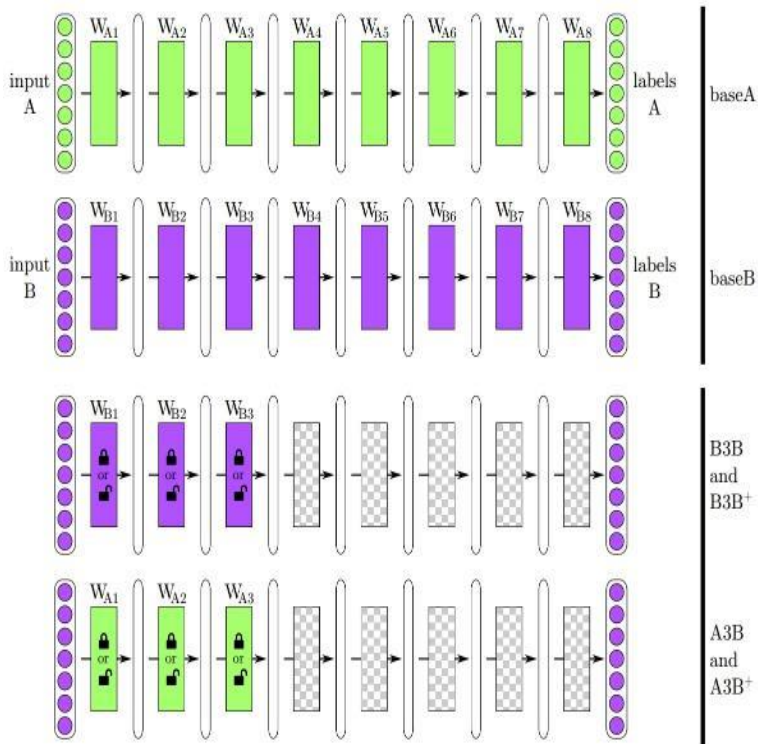


VOC 2007 test	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv	mAP
R-CNN pool <sub>5</sub>	51.8	60.2	36.4	27.8	23.2	52.8	60.6	49.2	18.3	47.8	44.3	40.8	56.6	58.7	42.4	23.4	46.1	36.7	51.3	55.7	44.2
R-CNN fc <sub>6</sub>	59.3	61.8	43.1	34.0	25.1	53.1	60.6	52.8	21.7	47.8	42.7	47.8	52.5	58.5	44.6	25.6	48.3	34.0	53.1	58.0	46.2
R-CNN fc <sub>7</sub>	57.6	57.9	38.5	31.8	23.7	51.2	58.9	51.4	20.0	50.5	40.9	46.0	51.6	55.9	43.3	23.3	48.1	35.3	51.0	57.4	44.7
R-CNN FT pool <sub>5</sub>	58.2	63.3	37.9	27.6	26.1	54.1	66.9	51.4	26.7	55.5	43.4	43.1	57.7	59.0	45.8	28.1	50.8	40.6	53.1	56.4	47.3
R-CNN FT fc <sub>6</sub>	63.5	66.0	47.9	37.7	29.9	62.5	70.2	60.2	32.0	57.9	47.0	53.5	60.1	64.2	52.2	31.3	55.0	50.0	57.7	63.0	53.1
R-CNN FT fc <sub>7</sub>	64.2	69.7	50.0	41.9	32.0	62.6	71.0	60.7	32.7	58.5	46.5	56.1	60.6	66.8	54.2	31.5	52.8	48.9	57.9	64.7	54.2
R-CNN FT fc <sub>7</sub> BB	<b>68.1</b>	<b>72.8</b>	<b>56.8</b>	<b>43.0</b>	<b>36.8</b>	<b>66.3</b>	<b>74.2</b>	<b>67.6</b>	<b>34.4</b>	<b>63.5</b>	<b>54.5</b>	<b>61.2</b>	<b>69.1</b>	<b>68.6</b>	<b>58.7</b>	<b>33.4</b>	<b>62.9</b>	<b>51.1</b>	<b>62.5</b>	<b>64.8</b>	<b>58.5</b>
DPM v5 [20]	33.2	60.3	10.2	16.1	27.3	54.3	58.2	23.0	20.0	24.1	26.7	12.7	58.1	48.2	43.2	12.0	21.1	36.1	46.0	43.5	33.7
DPM ST [28]	23.8	58.2	10.5	8.5	27.1	50.4	52.0	7.3	19.2	22.8	18.1	8.0	55.9	44.8	32.4	13.3	15.9	22.8	46.2	44.9	29.1
DPM HSC [31]	32.2	58.3	11.5	16.3	30.6	49.9	54.8	23.5	21.5	27.7	34.0	13.7	58.1	51.6	39.9	12.4	23.5	34.4	47.4	45.2	34.3

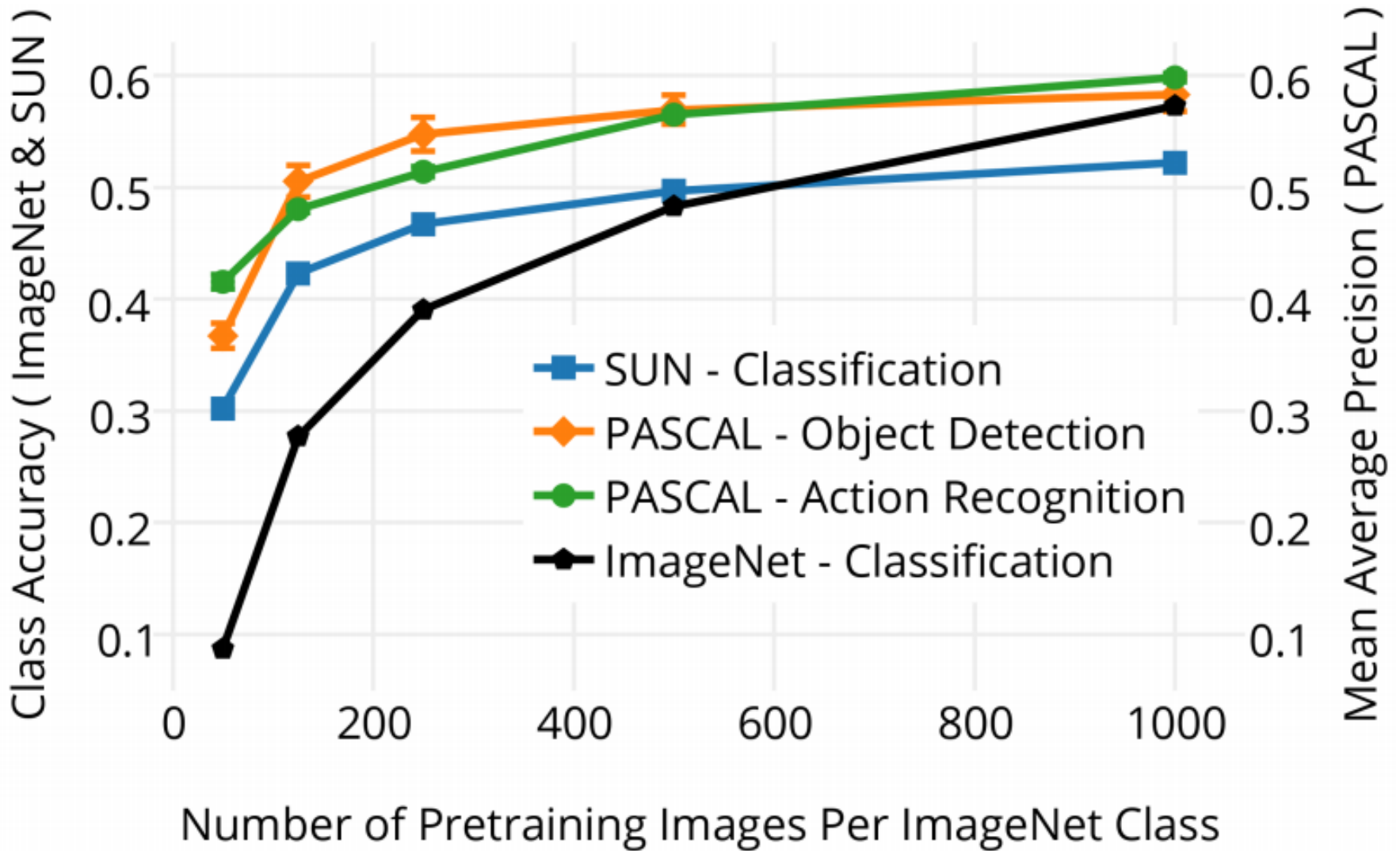
Rich feature hierarchies for accurate object detection and semantic segmentation, [[Girshick et al. CVPR 2014](#)]



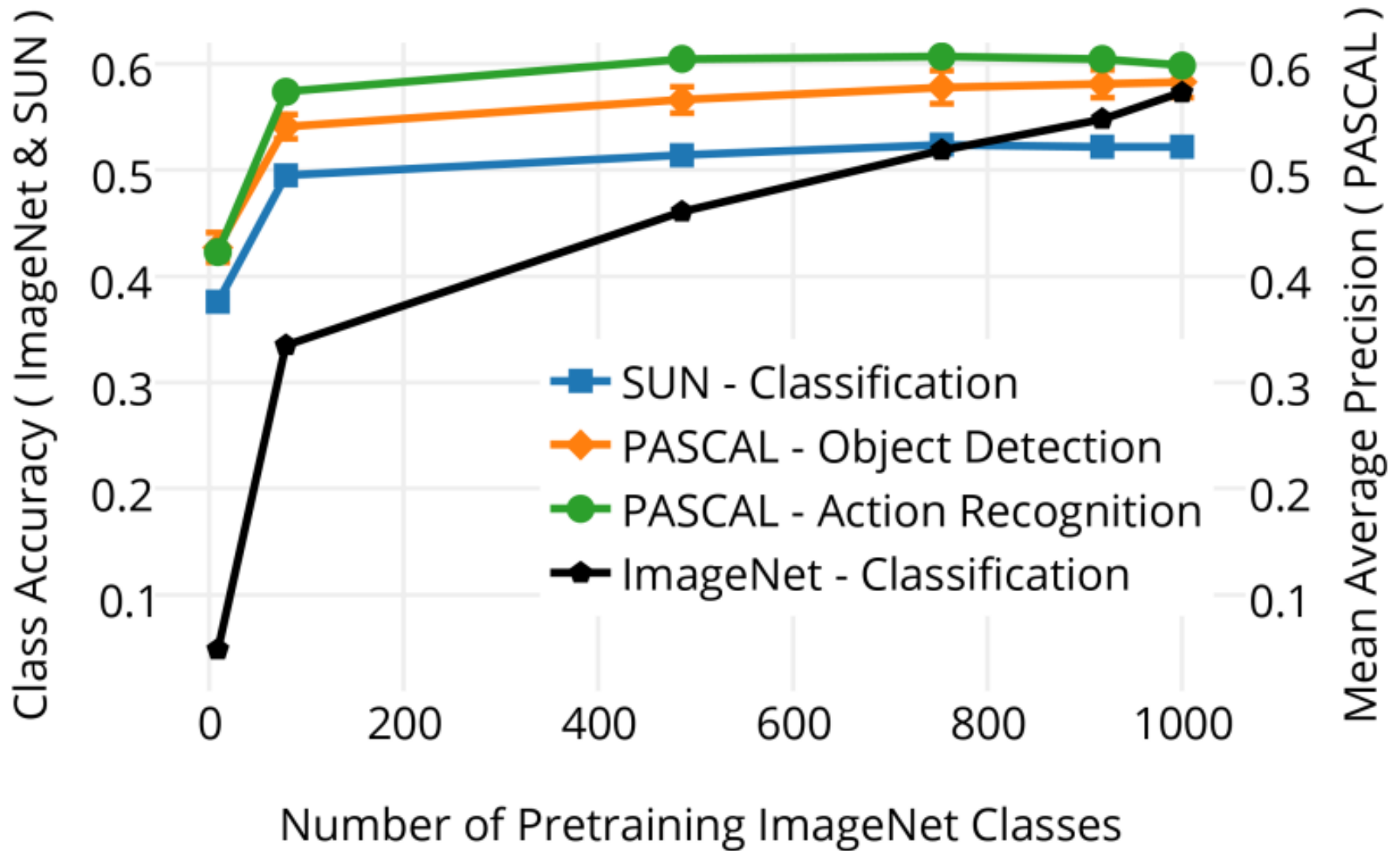
# How transferable are features in CNN?



How transferable are features in deep neural networks [[Yosinski NIPS 2014](#)]



What makes ImageNet good for transfer learning? [[Huh arXiv 2016](#)]



What makes ImageNet good for transfer learning? [[Huh arXiv 2016](#)]

Orig  
918  
753  
486  
79  
9  
Rnd



What makes ImageNet good for transfer learning? [[Huh arXiv 2016](#)]

# Things to remember

- Training CNN
  - Dropout
  - Data augmentation
  - Activation
  - Batch normalization
- Transfer learning
  - Two strategies
    - CNN code
    - Finetuning
  - When and how to transfer
  - Characteristics of transfer learning