

# **Self Supervision**

through Context, Color,  
and Physical Interactions

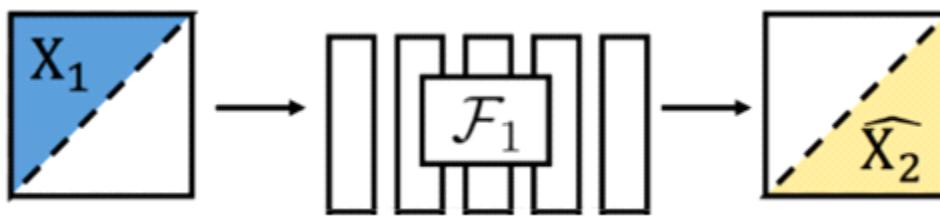
Lama Affara

# Basic Learning Types

- Supervised Learning
  - Classifying data into labels
- Unsupervised Learning
  - Clustering data, Anomaly detection
- Semi/Weak Supervised Learning
  - Missing/Incomplete labels
- Reinforcement Learning
  - Rewards from actions
- What about Self Supervision?

# Self Supervision

- No supervision from other agents
- Supervision can come from other modalities or from time



# Self Supervision Examples

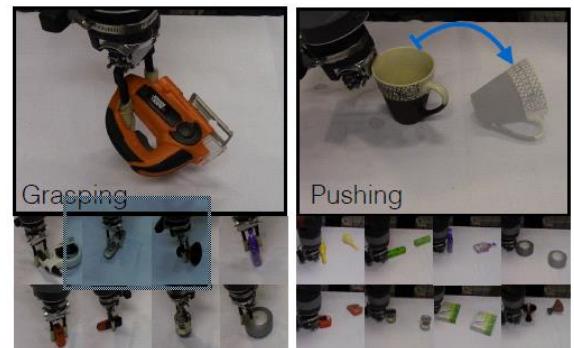
## Context



## Color



## Physical Interactions



# Self Supervision Examples

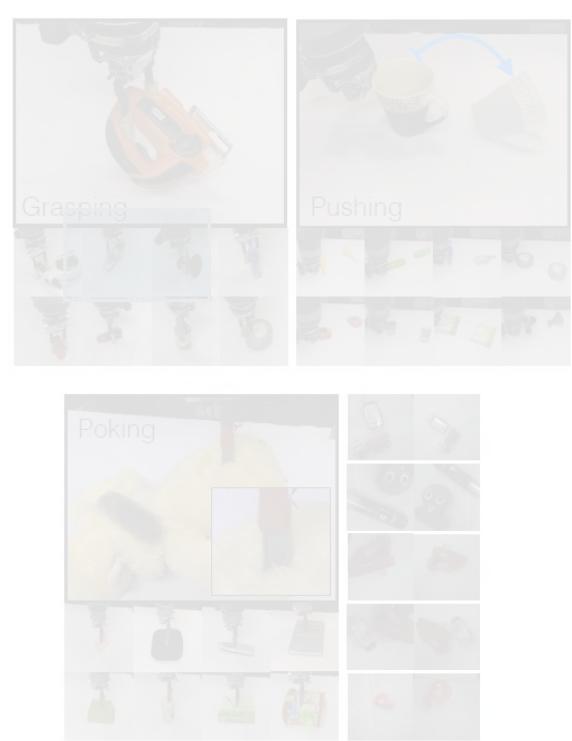
## Context



## Color



## Physical Interactions



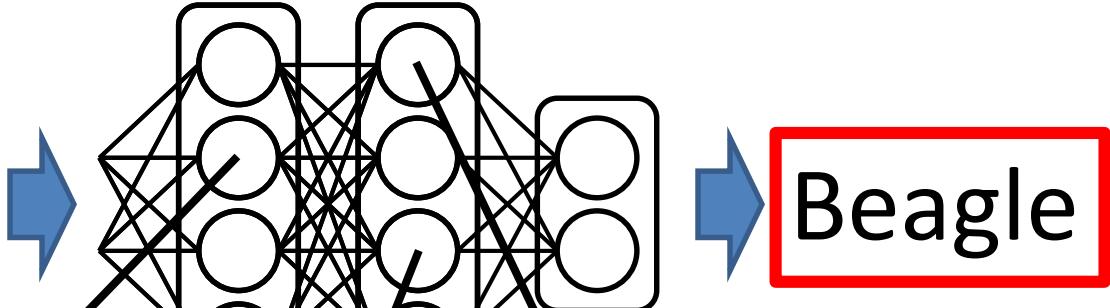
# Using Patch Context for Self Supervision

## UNSUPERVISED VISUAL REPRESENTATION LEARNING BY CONTEXT PREDICTION



Carl Doersch<sup>1,2</sup> Abhinav Gupta<sup>1</sup> Alexei A. Efros<sup>2</sup>

# ImageNet + Deep Learning



Materials?

Pose?

Parts?

Geometry?

Boundaries?

*Do we even need this task? Labels?*

# Context as Supervision

[Collobert & Weston 2008; Mikolov et al. 2013]

house, where the professor lived without his wife and child; or so he said jokingly sometimes: "Here's where I live. My house." His daughter often added, without resentment, for the visitor's information, "It started out to be for me, but it's really his." And she might reach in to bring forth an inch-high table lamp with fluted shade, or a blue dish the size of her little fingernail, marked "Kitty" and half full of eternal milk, but she was sure to replace these, after they had been admired, pretty near exactly where they had been. The little house was very orderly, and just big enough for all it contained, though to some tastes the bric-à-brac in the parlor might seem excessive. The daughter's preference was for the store-bought gimmicks and appliances, the toasters and carpet sweepers of Lilliput, but she knew that most adult visitors would

Deep  
Net

# Context Prediction for Images

?

?

?

?



?

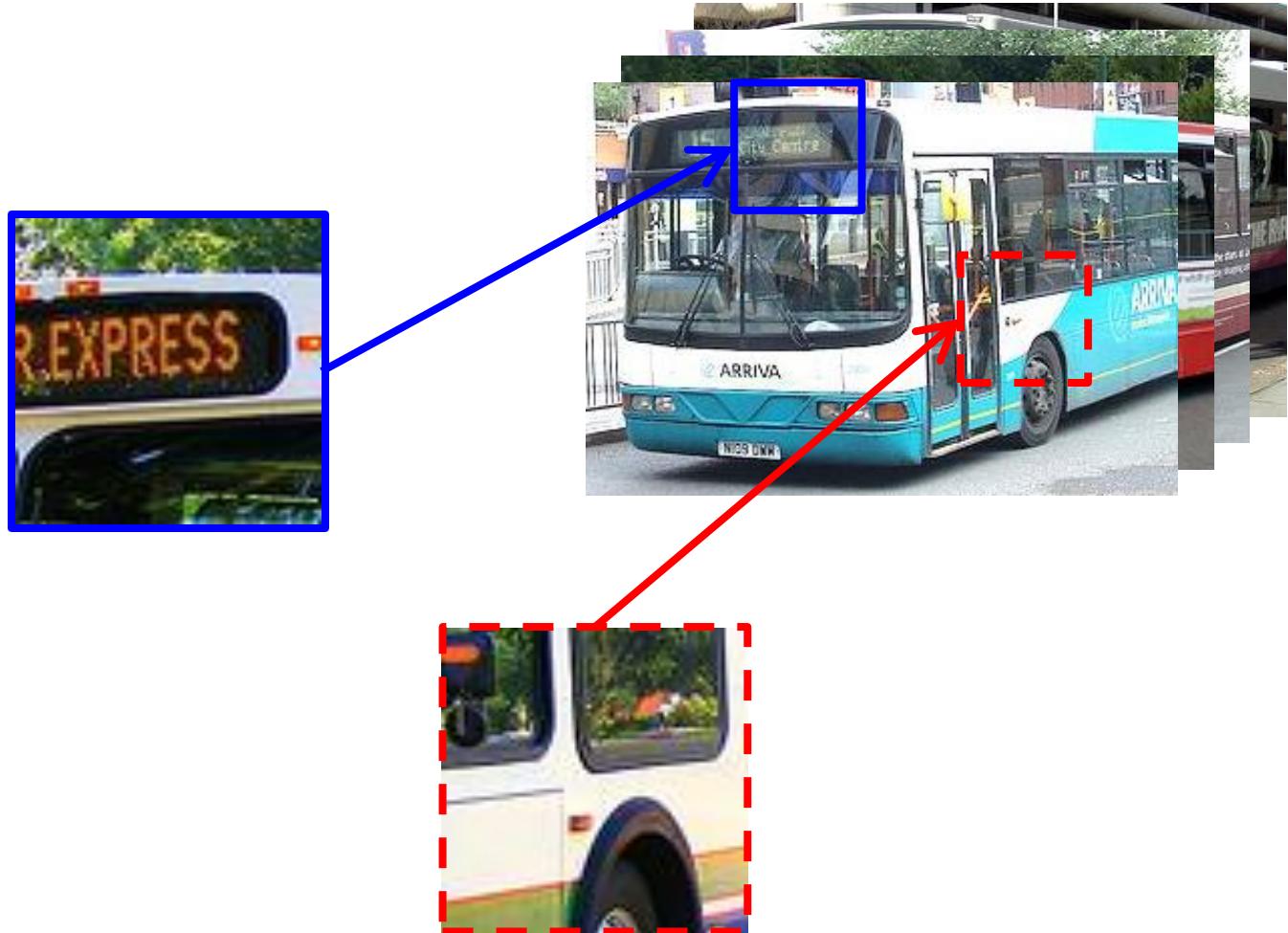
?

?

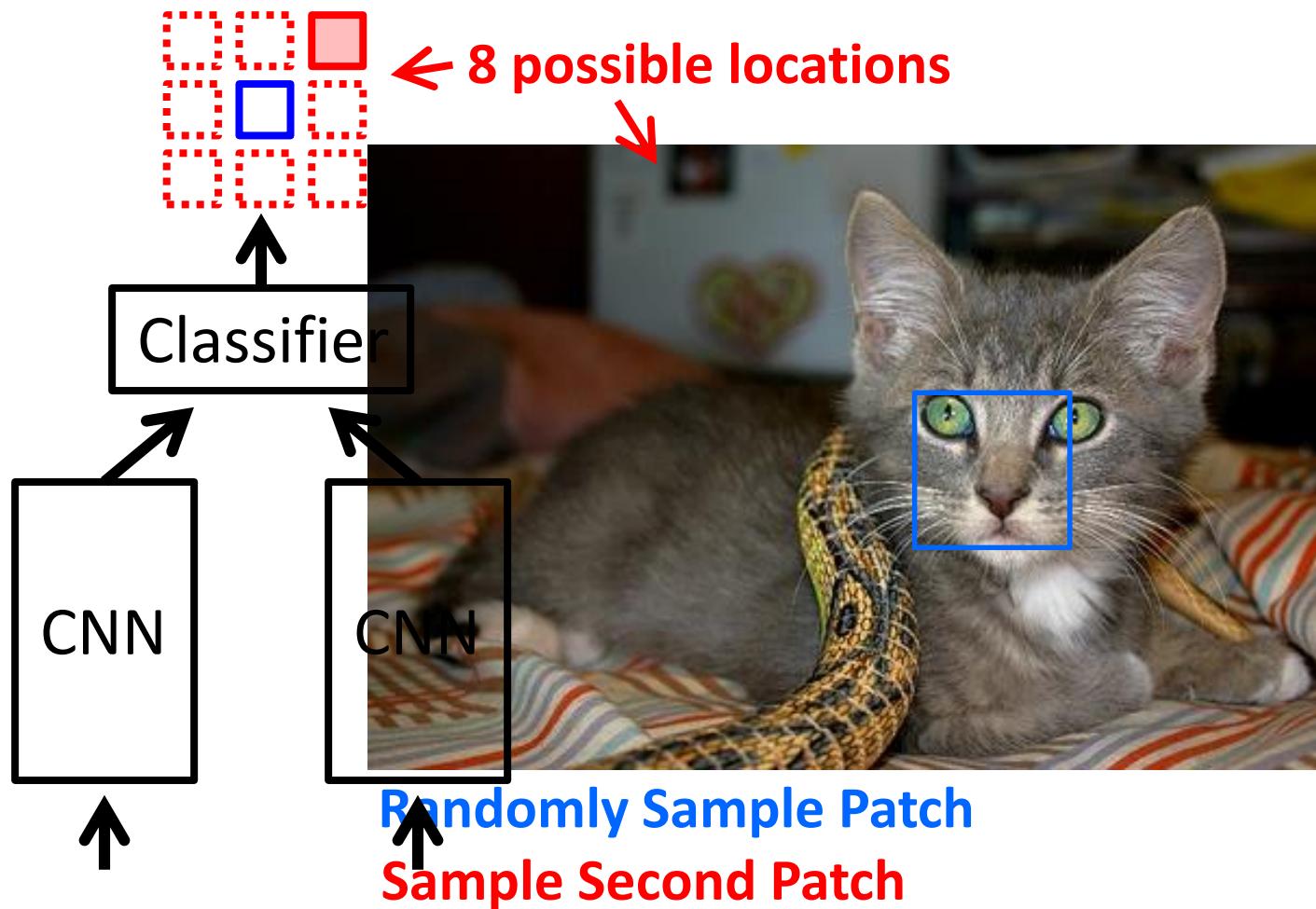
A

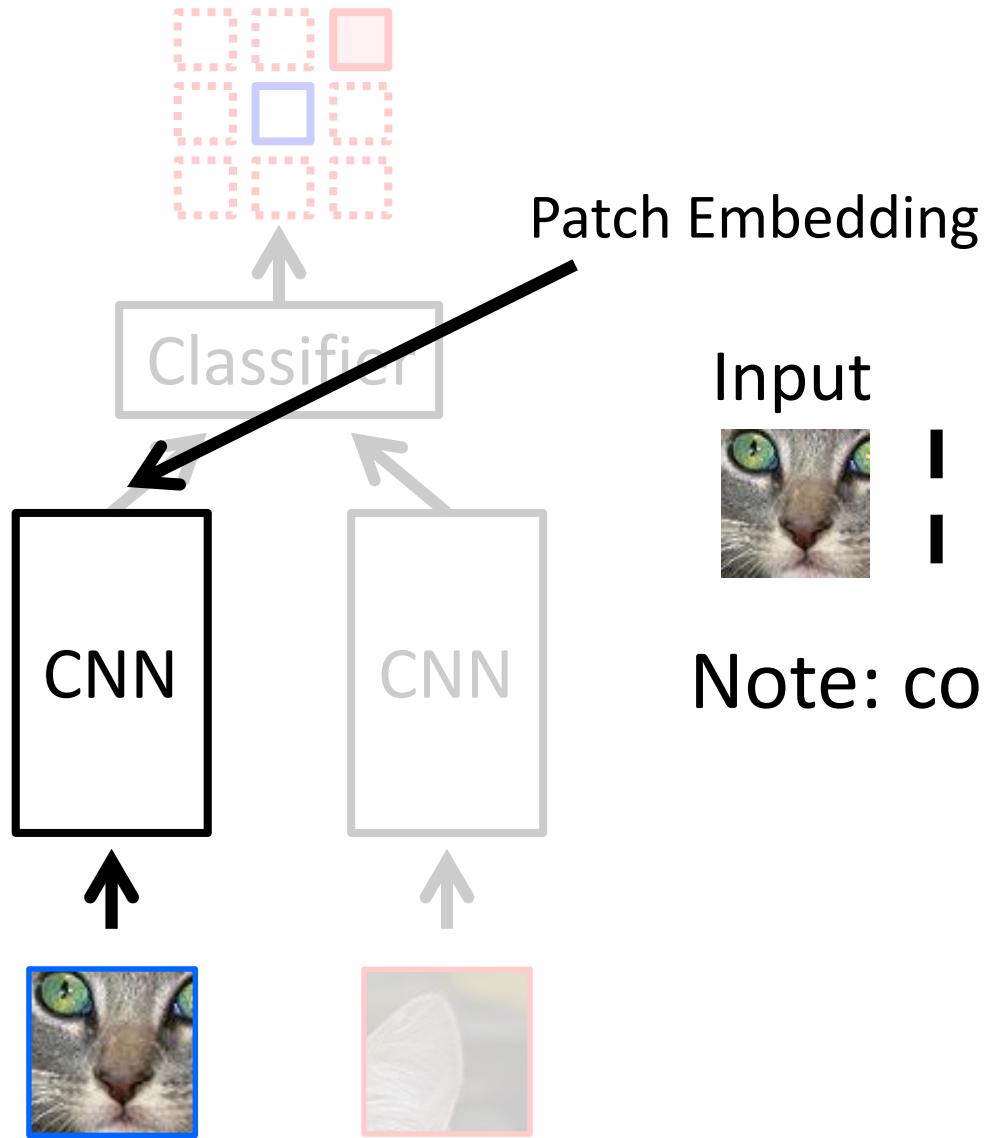
B

# Semantics from a non-semantic task



# Relative Position Task





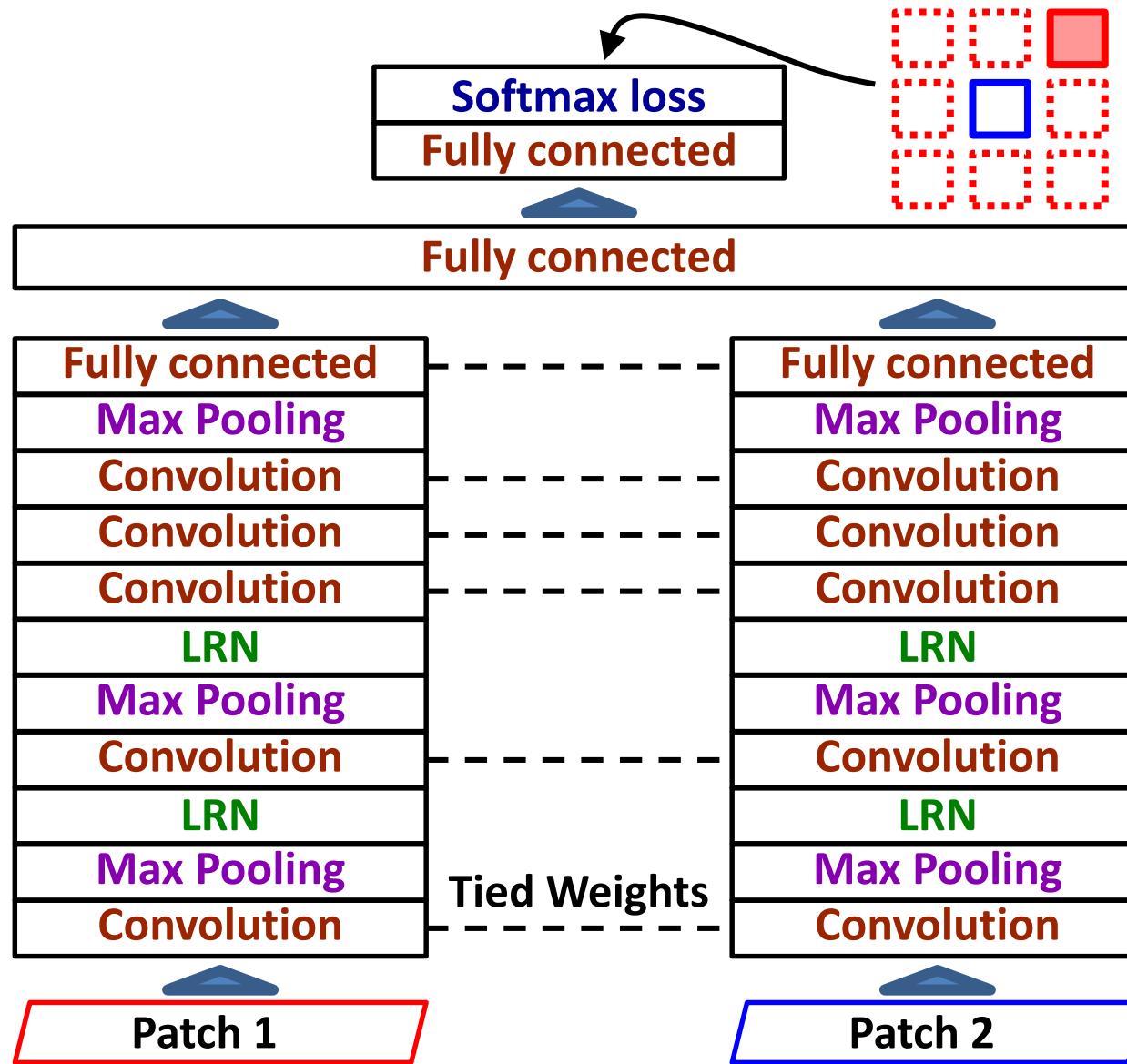
Patch Embedding

Input

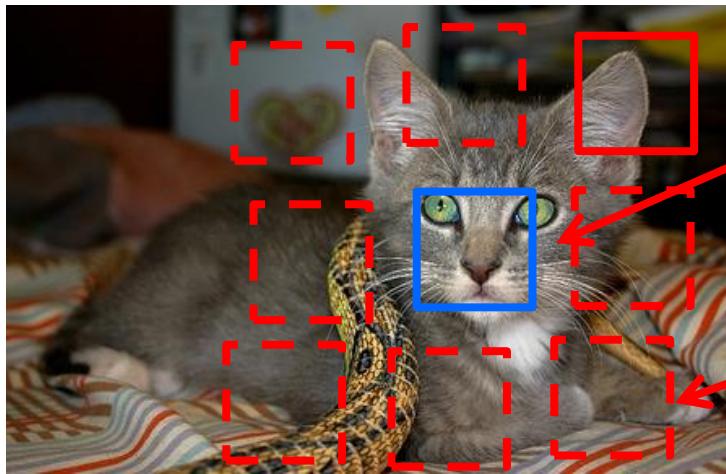
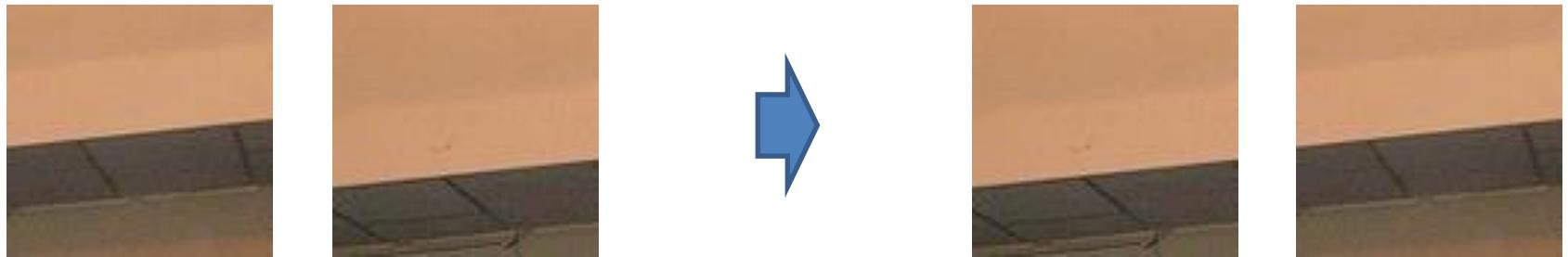


Note: connects ***across*** instances!

# Architecture



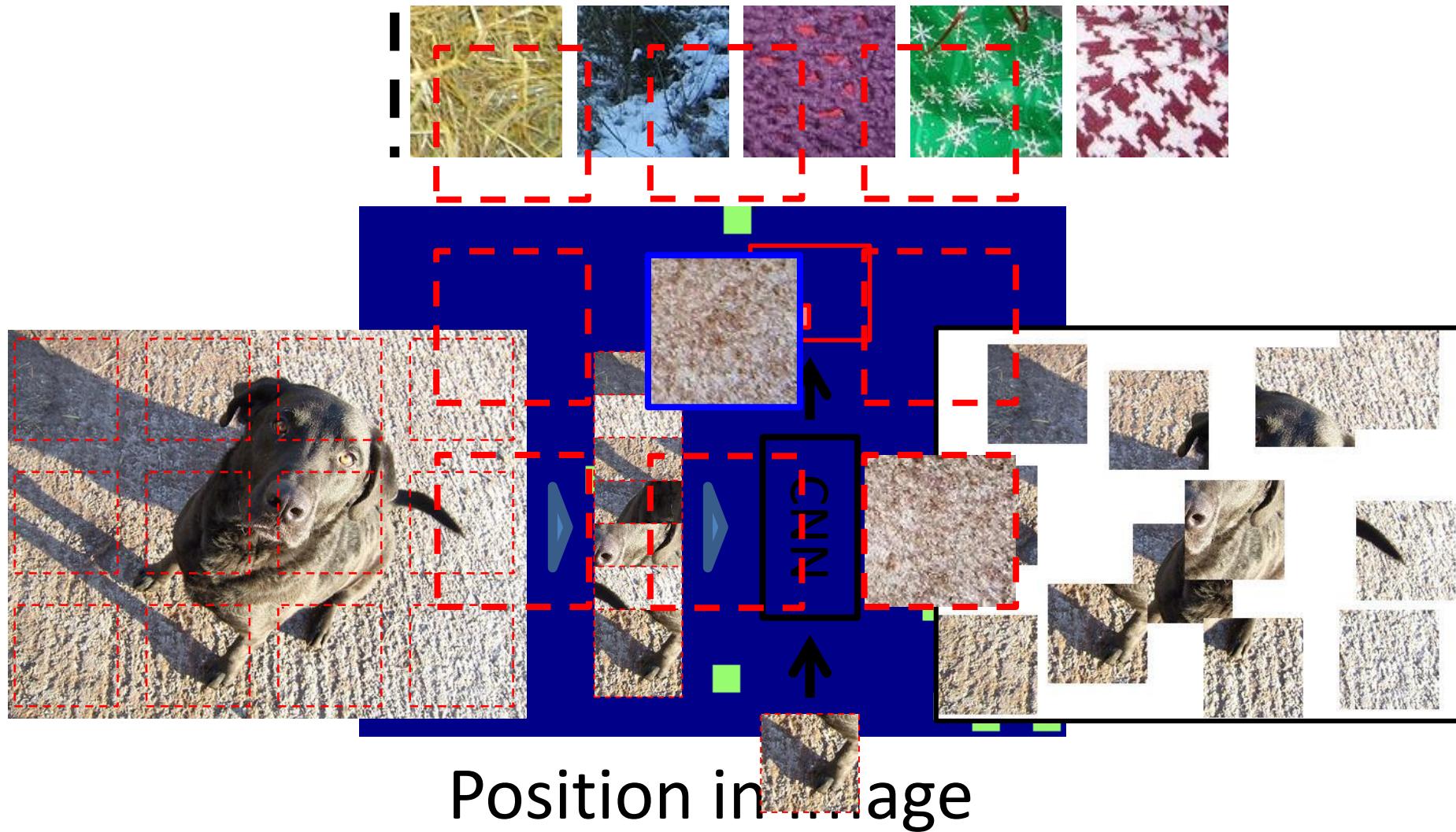
# Avoiding Trivial Shortcuts



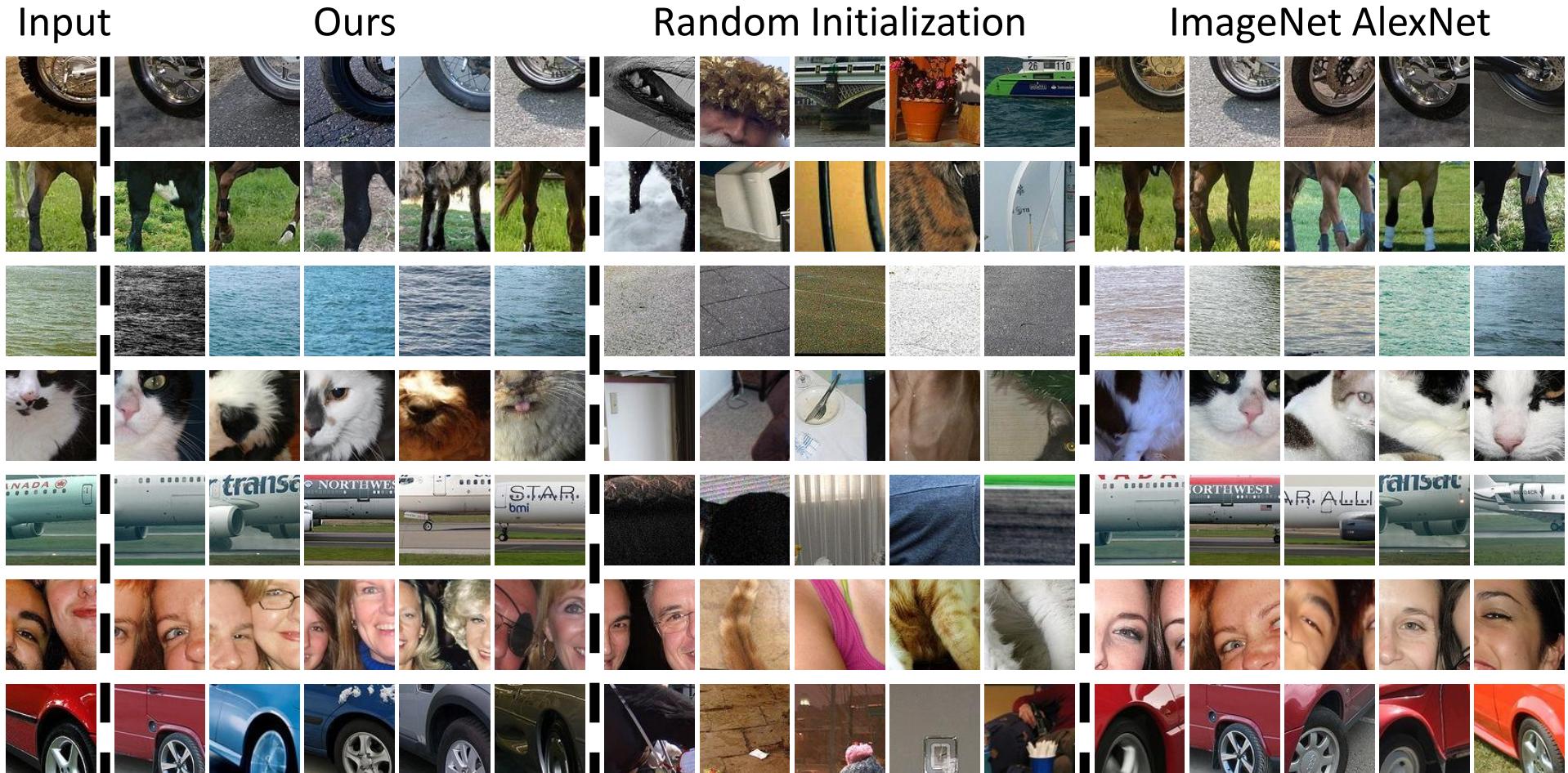
Include a gap

Jitter the patch locations

# A Not-So “Trivial” Shortcut



# What is learned?



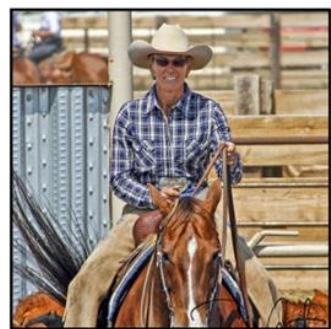
# Still don't capture everything



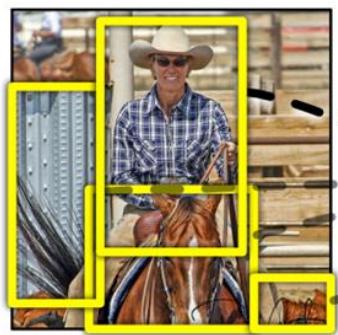
## You don't always need to learn!



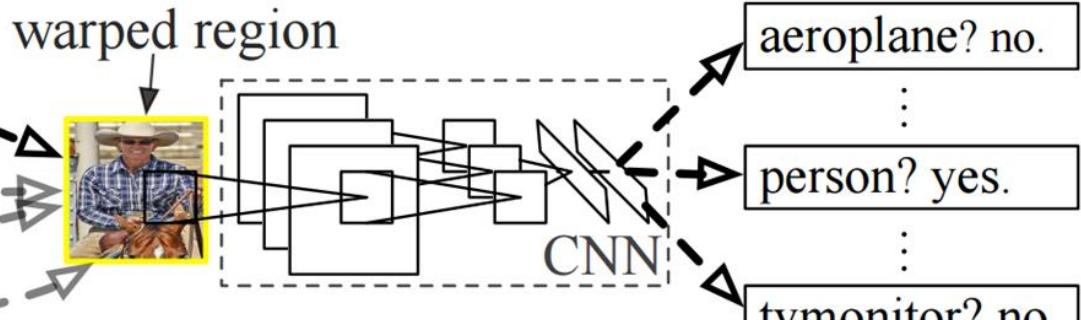
# Pre-Training for R-CNN



1. Input image



2. Extract region proposals (~2k)



3. Compute CNN features

4. Classify regions

Pre-train on relative-position task, w/o labels

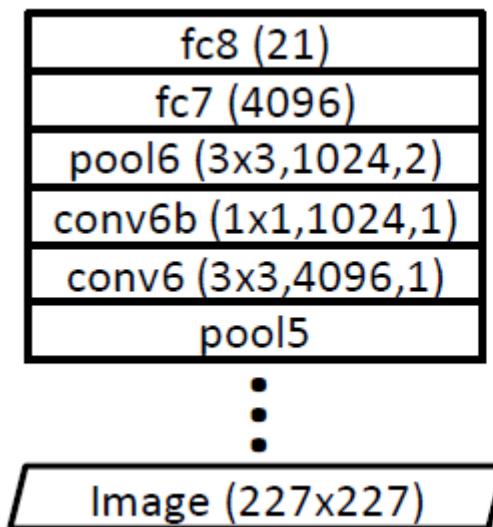
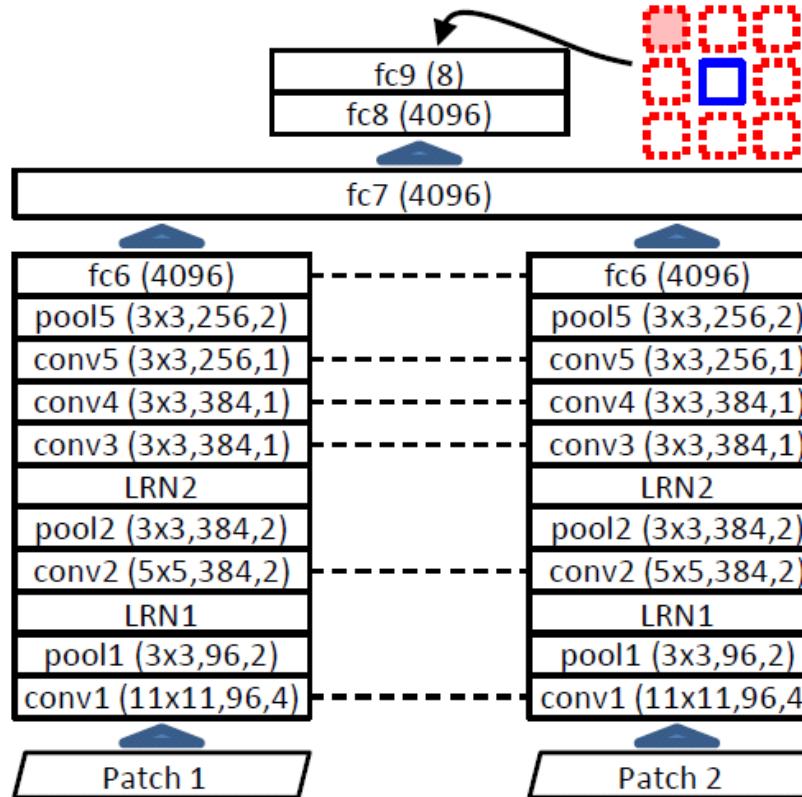
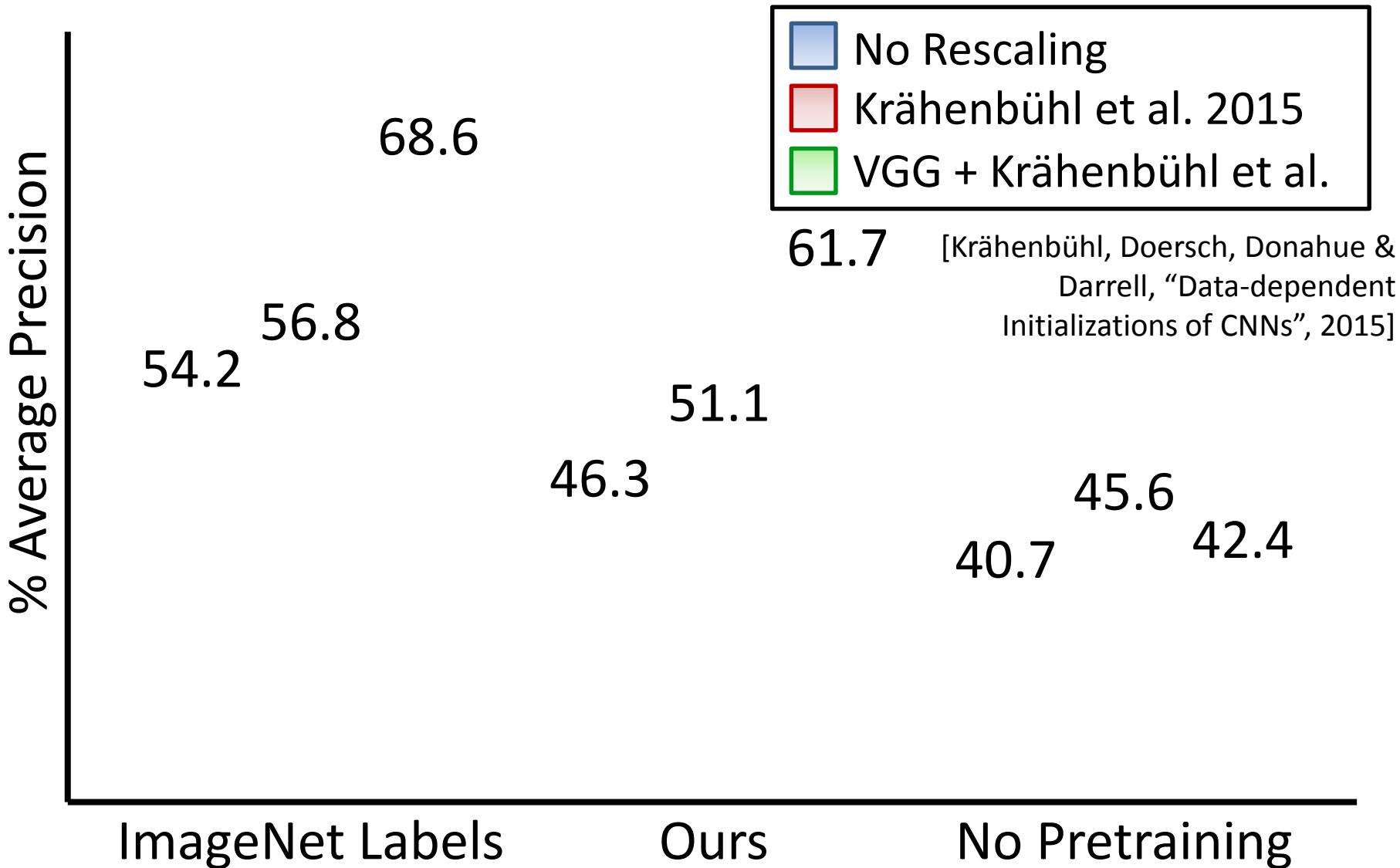


Figure 6. Our architecture for Pascal VOC detection. Layers from conv1 through pool5 are copied from our patch-based network (Figure 3). The new 'conv6' layer is created by converting the fc6 layer into a convolution layer. Kernel sizes, output units, and stride are given in parentheses, as in Figure 3.

# VOC 2007 Performance

(pretraining for R-CNN)



# Self Supervision Examples

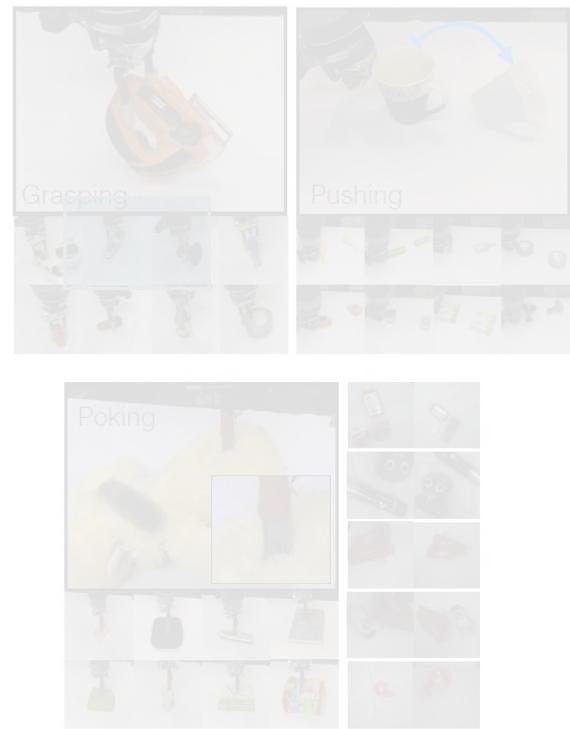
Context



Color



Physical Interactions



# Using Colorization for Self Supervision

# Colorful Image Colorization

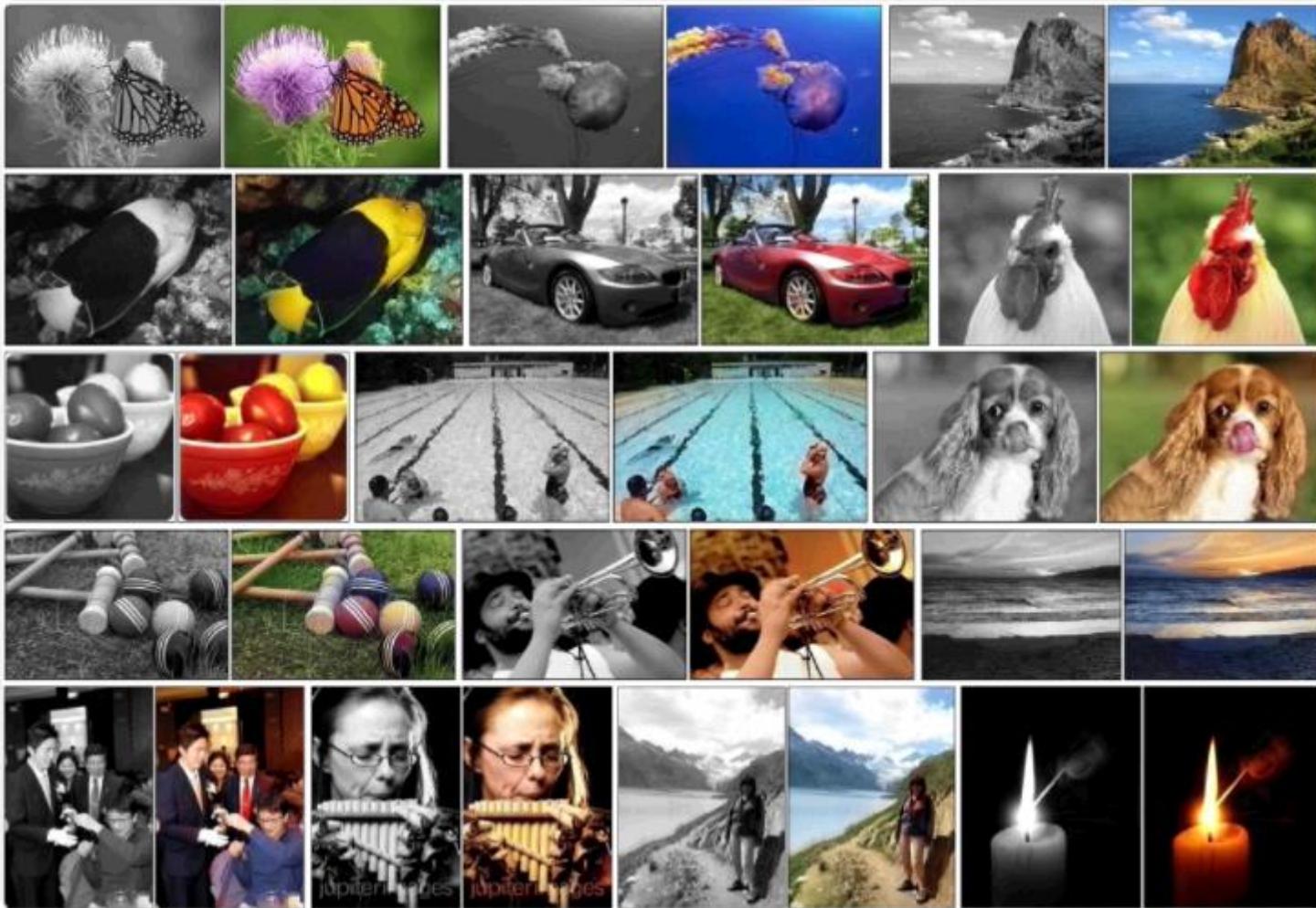
Richard Zhang

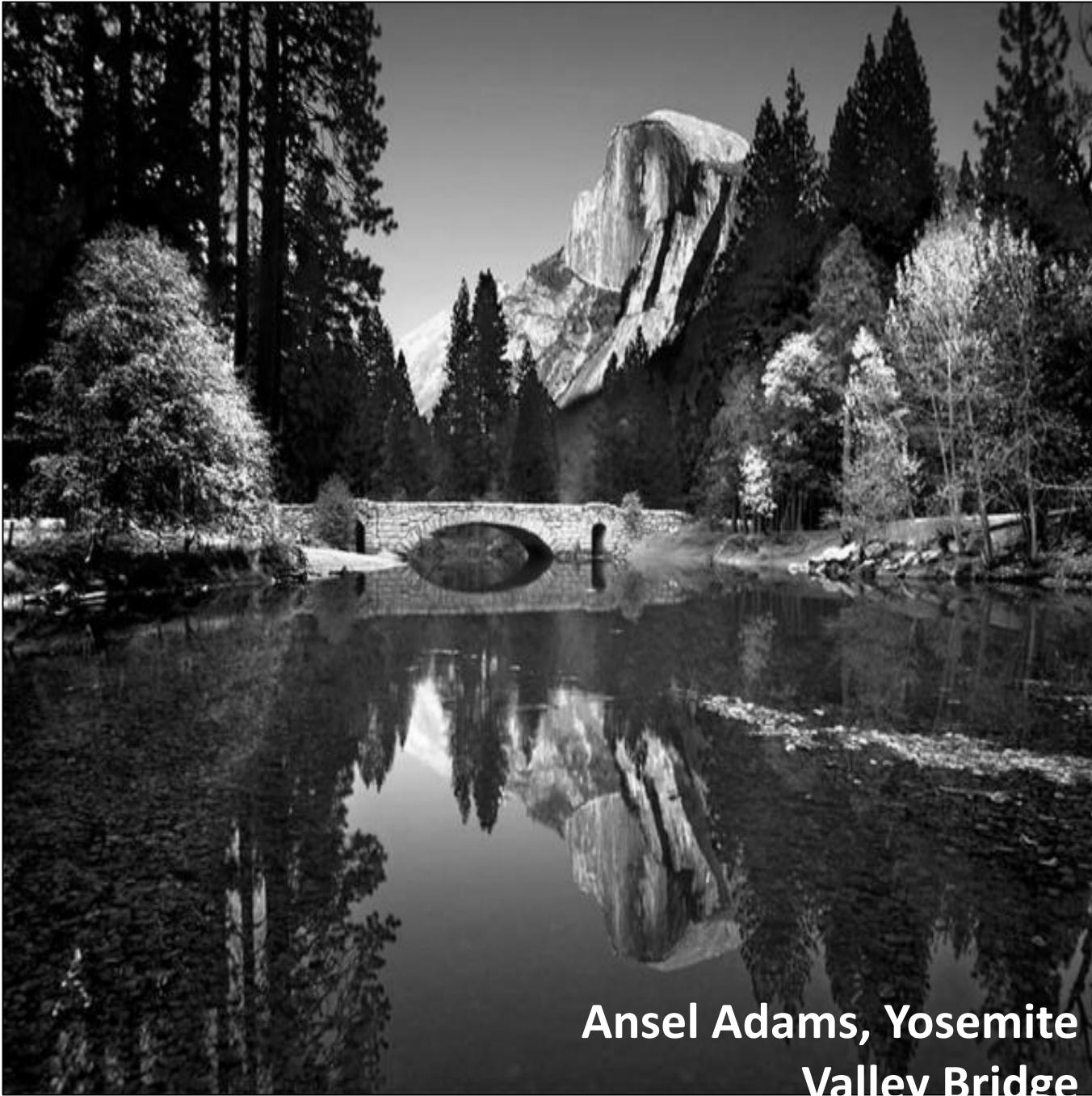
Phillip Isola

Alexei A. Efros

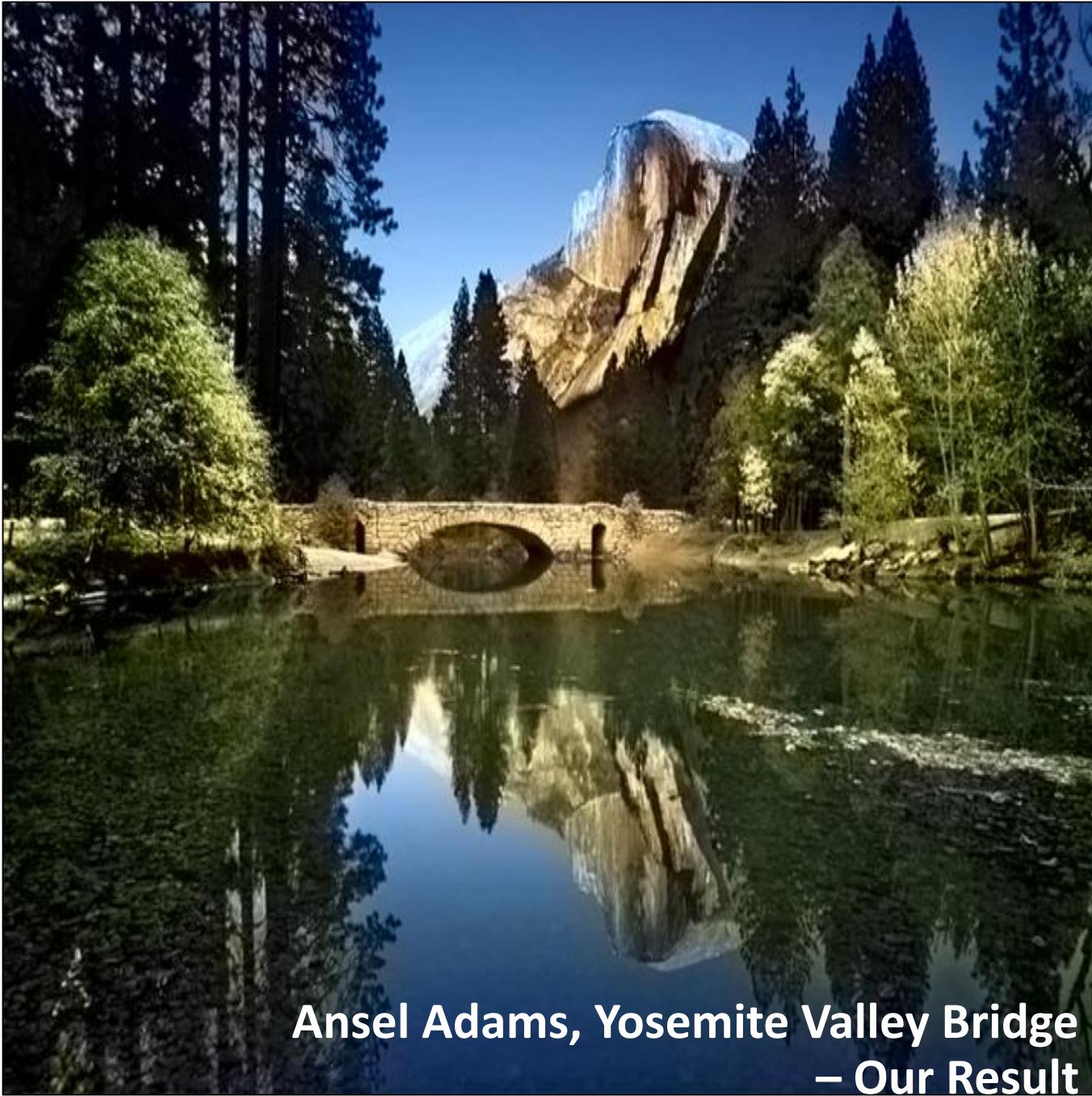
[Demo] [GitHub] [Talk] [Slides] [Paper]

Also check out our new work on [Interactive Deep Colorization!](#)

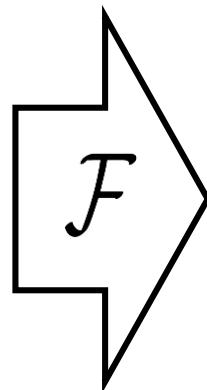




Ansel Adams, Yosemite  
Valley Bridge



**Ansel Adams, Yosemite Valley Bridge  
– Our Result**

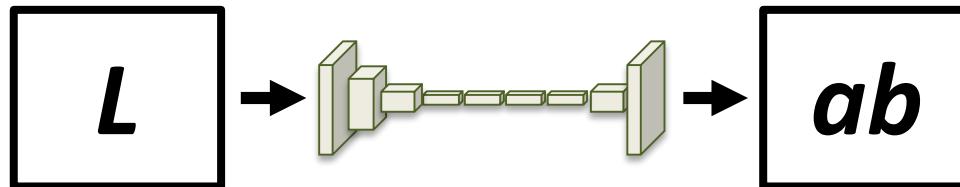


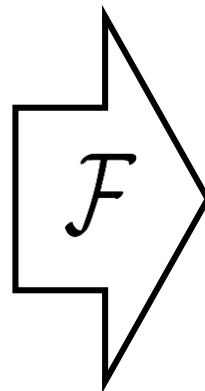
Grayscale image:  $L$  channel

$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

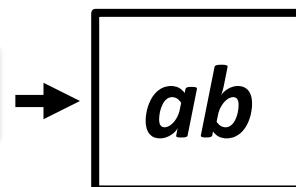
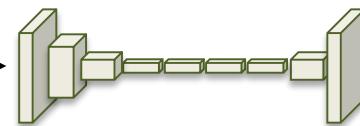
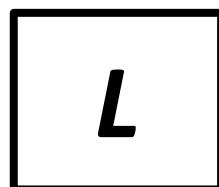
Color information:  $ab$  ch

$$\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$$





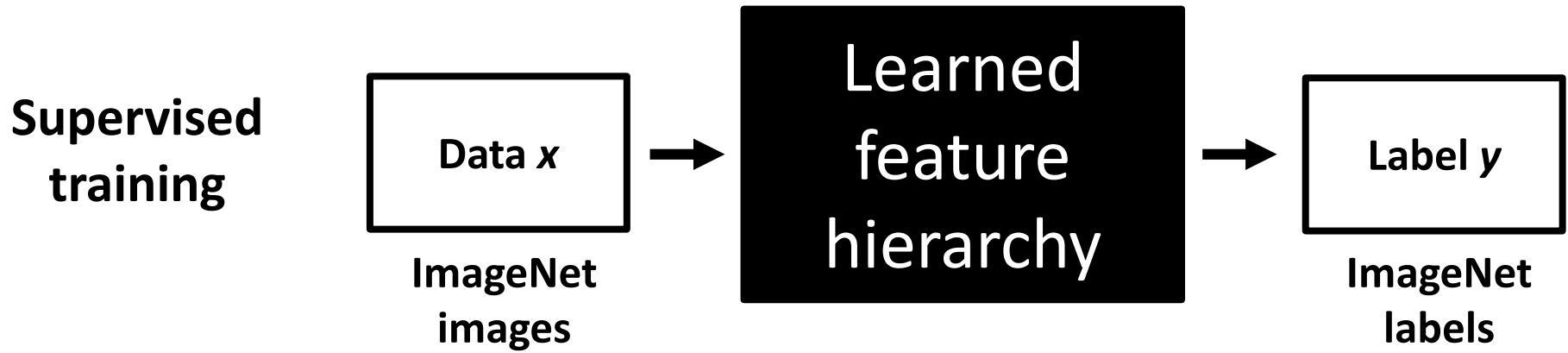
Grayscale image:  $L$  channel  
 $X \in \mathbb{R}^{H \times W \times 1}$



Concatenate (L,ab)  
(X,  $\hat{Y}$ )

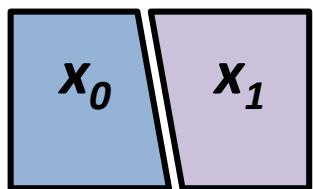


# Predicting Labels from Data

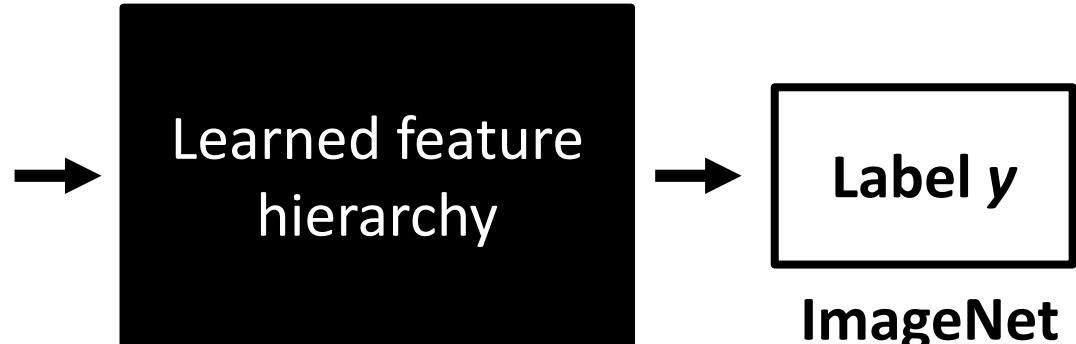


# Predicting Data from Data

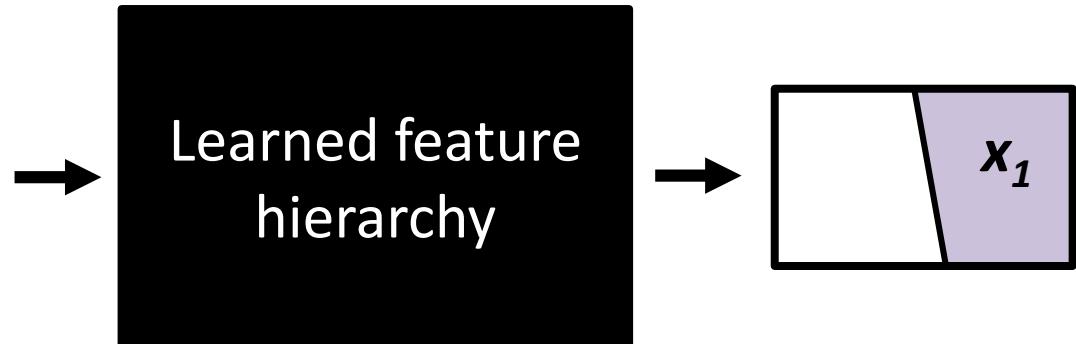
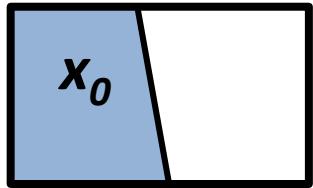
Supervised  
training



ImageNet  
images



Unsupervised/  
Self-supervised  
training



# Inherent Ambiguity



Grayscale

# Inherent Ambiguity



Our Output



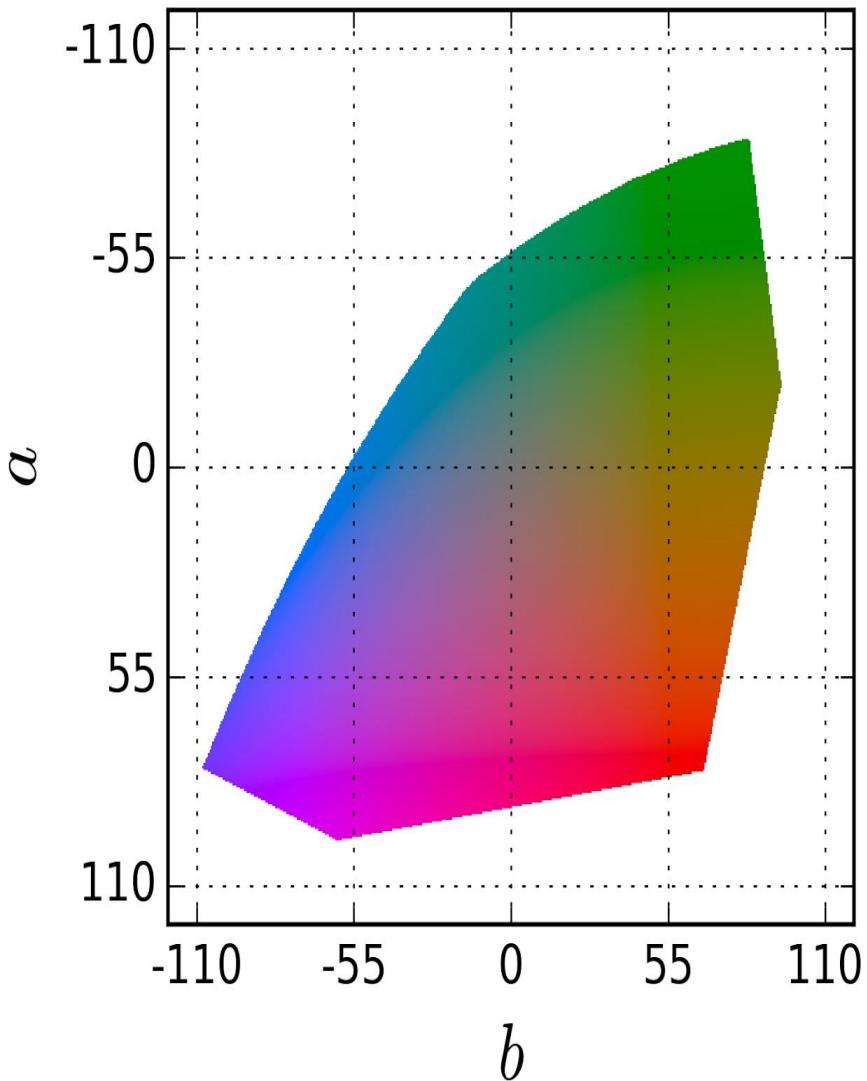
Ground Truth

# Better Loss Function

Colors in *ab* space

- Regression with L2 loss inadequate

$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$



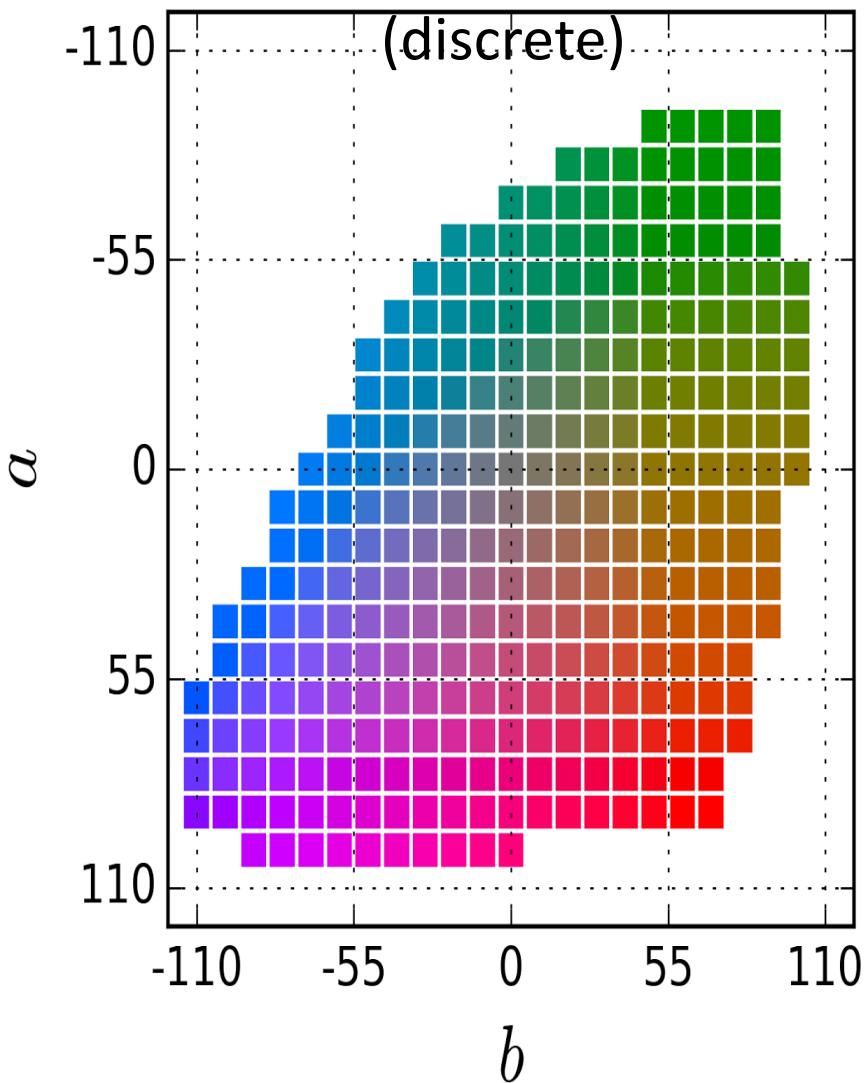
# Better Loss Function

Colors in *ab* space

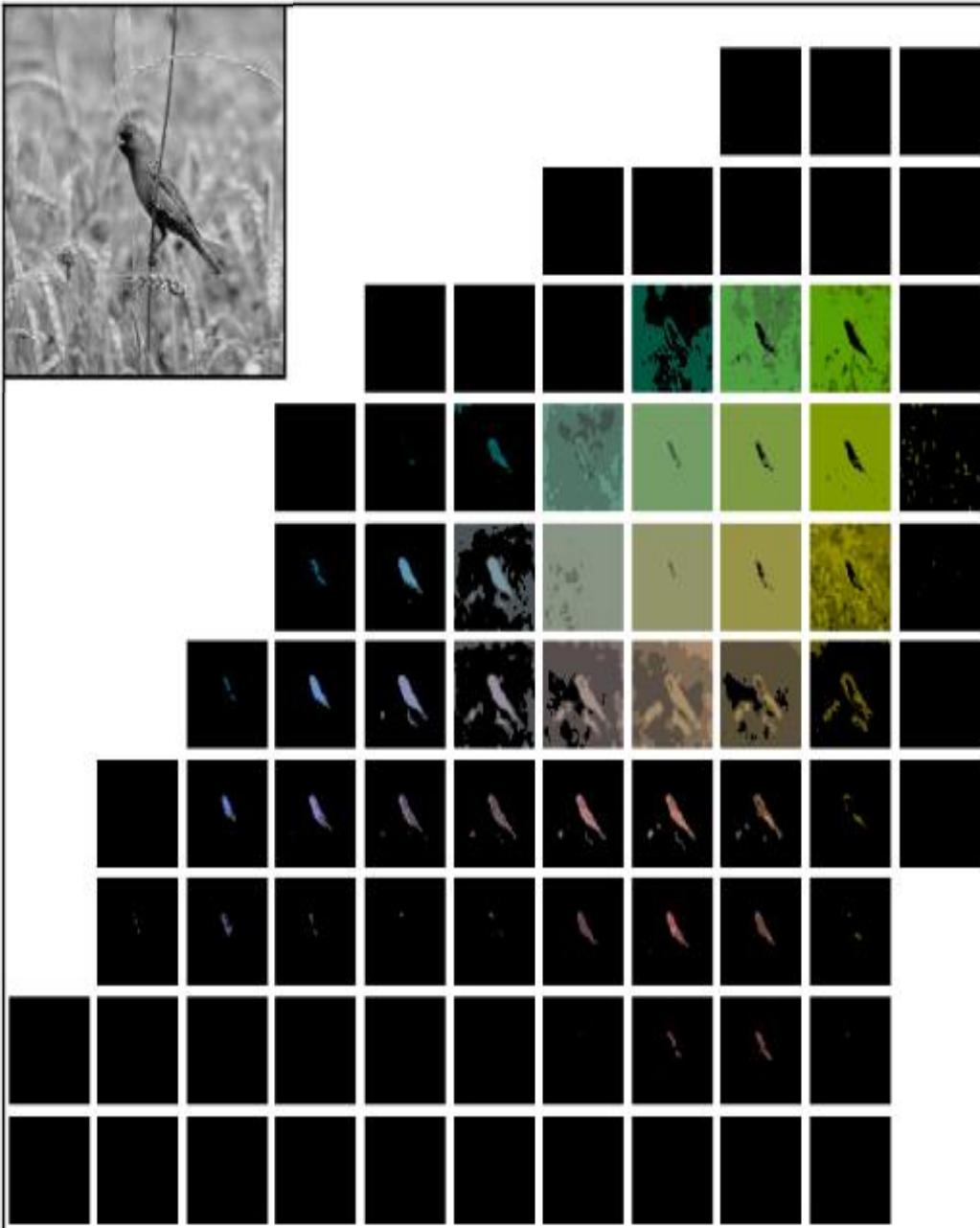
- Regression with L2 loss inadequate

$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$

$$L(\hat{\mathbf{Z}}, \mathbf{Z}) = -\frac{1}{HW} \sum_{h,w} \sum_q \mathbf{Z}_{h,w,q} \log(\hat{\mathbf{Z}}_{h,w,q})$$



*a*



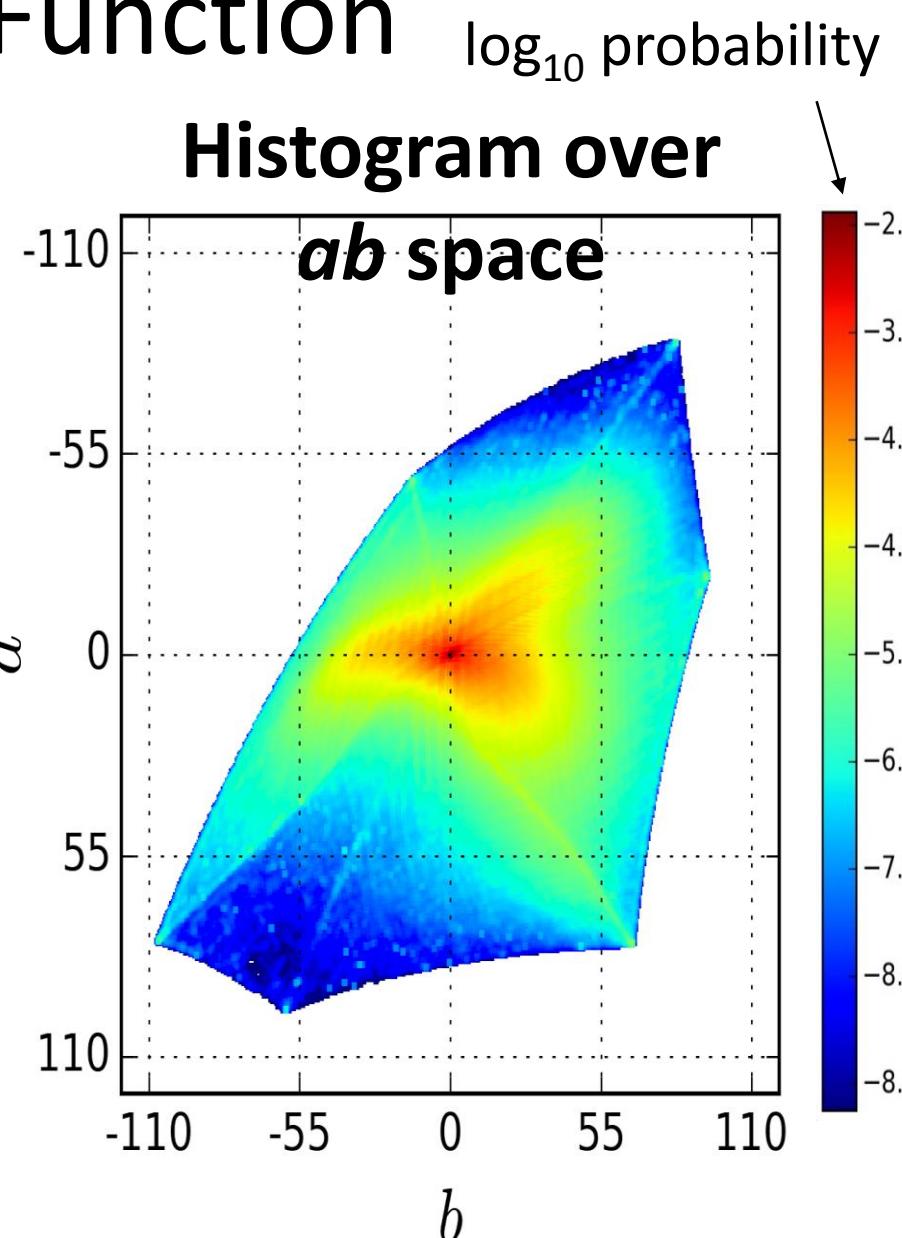
*b*

# Better Loss Function

- Regression with L2 loss inadequate

$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$

$$L(\hat{\mathbf{Z}}, \mathbf{Z}) = -\frac{1}{HW} \sum_{h,w} \sum_q \mathbf{Z}_{h,w,q} \log(\hat{\mathbf{Z}}_{h,w,q})$$



# Better Loss Function

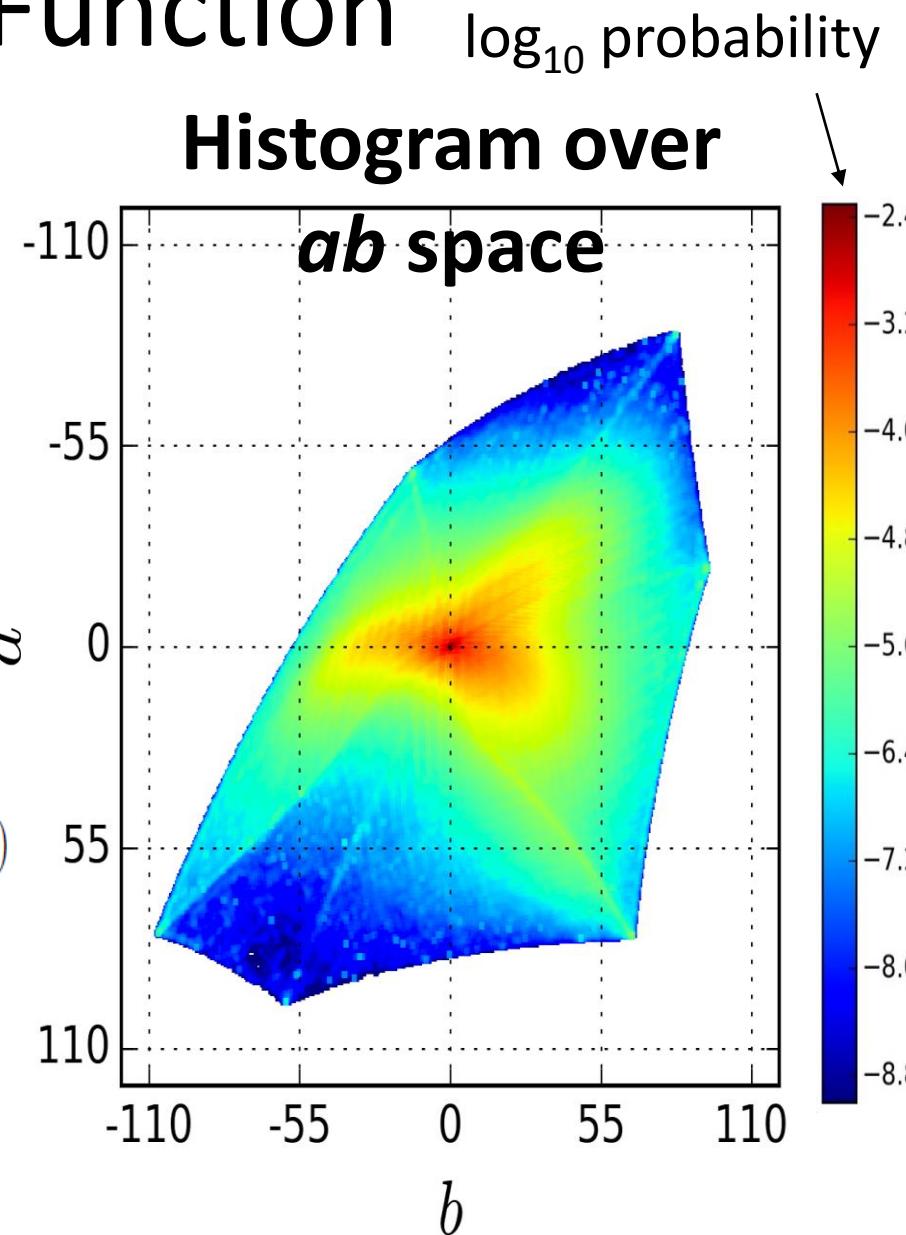
- Regression with L2 loss inadequate

$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$

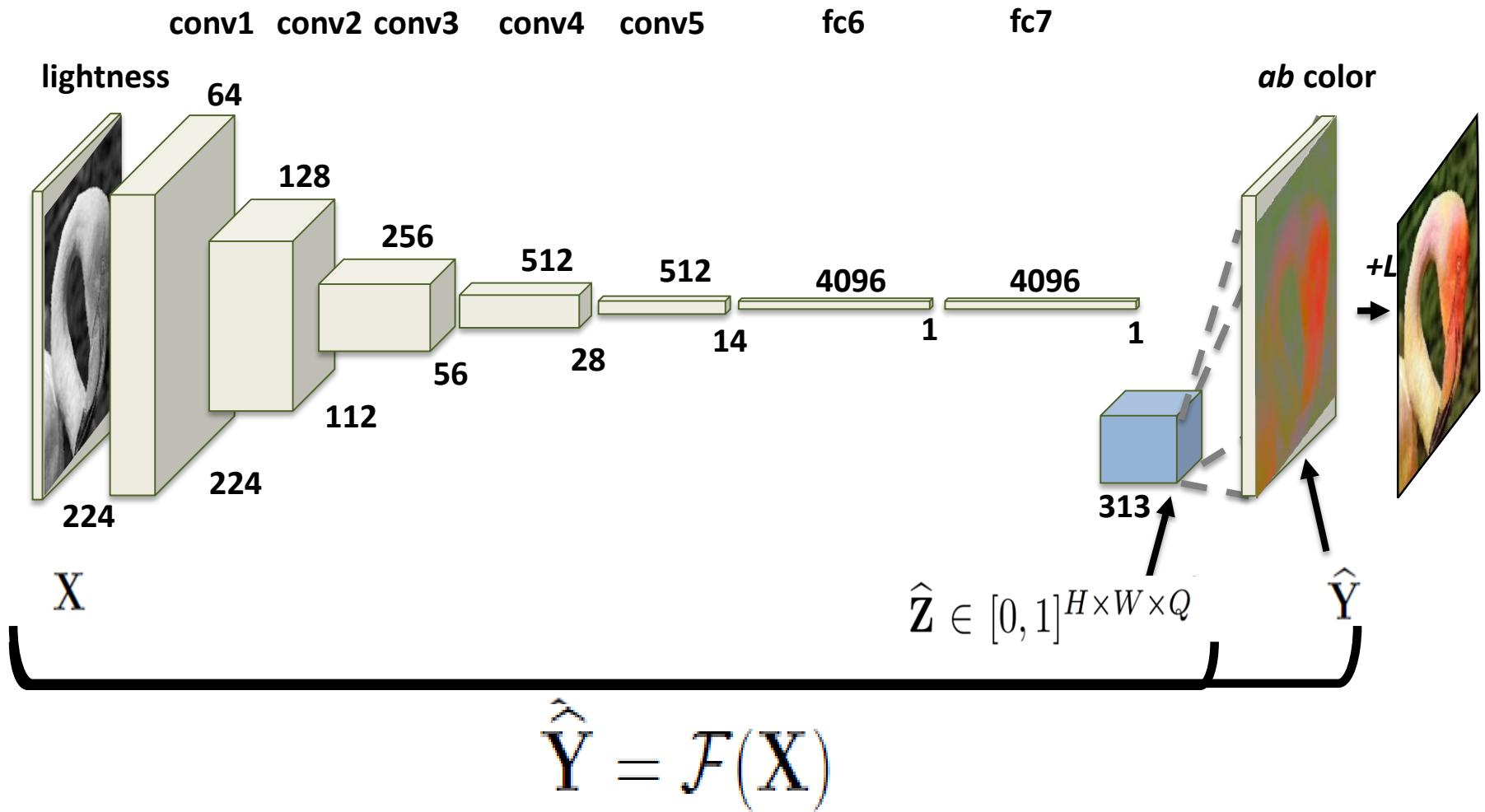
$$L(\hat{\mathbf{Z}}, \mathbf{Z}) = -\frac{1}{HW} \sum_{h,w} \sum_q \mathbf{Z}_{h,w,q} \log(\hat{\mathbf{Z}}_{h,w,q})$$

$$L(\hat{\mathbf{Z}}, \mathbf{Z}) = -\frac{1}{HW} \sum_{h,w} v(\mathbf{Z}_{h,w}) \sum_q \mathbf{Z}_{h,w,q} \log(\hat{\mathbf{Z}}_{h,w,q})$$

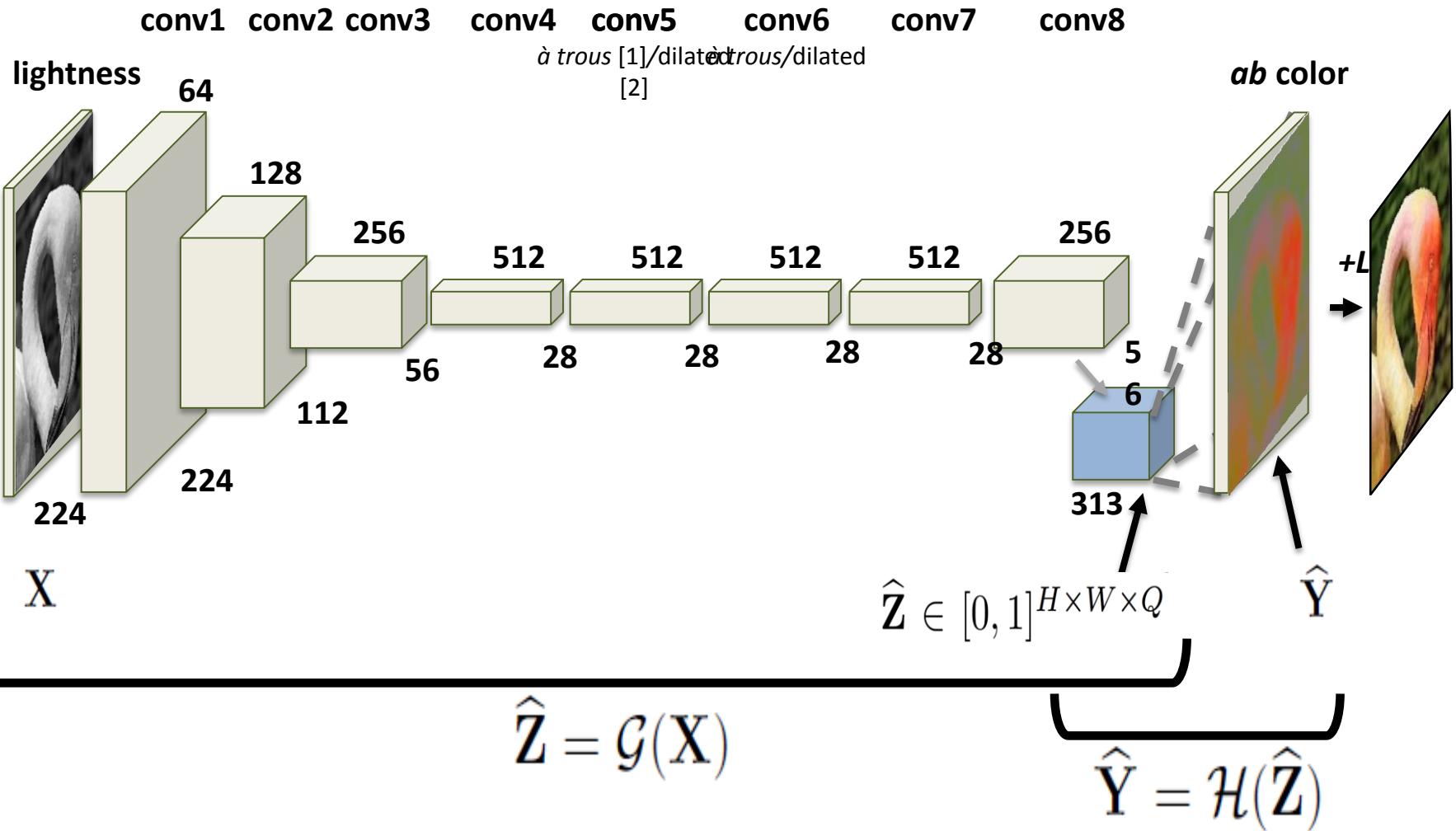
- **Class rebalancing** to encourage learning of *rare* colors



# Network Architecture



# Network Architecture



GrdaptTruth

L2 Regression

Class w/ Rebalancing



# Failure Cases



# Biases





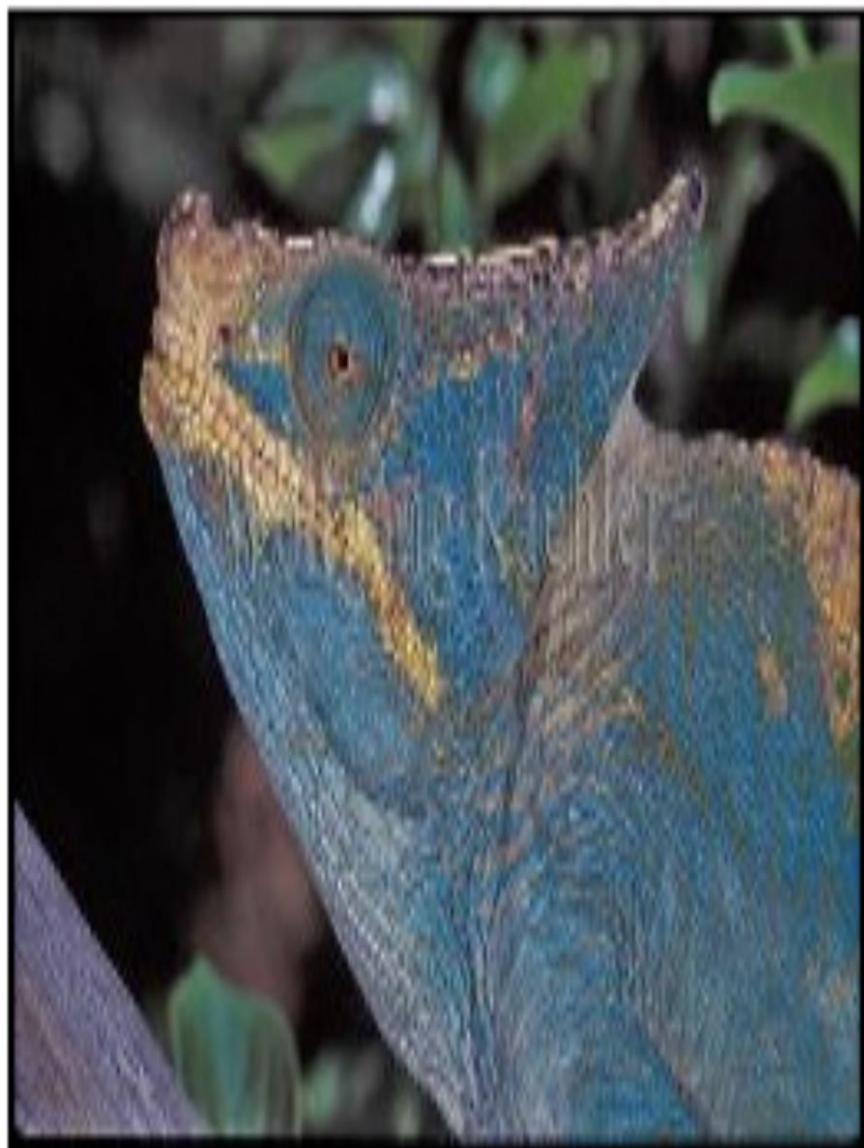
**Fake, 0% fooled**



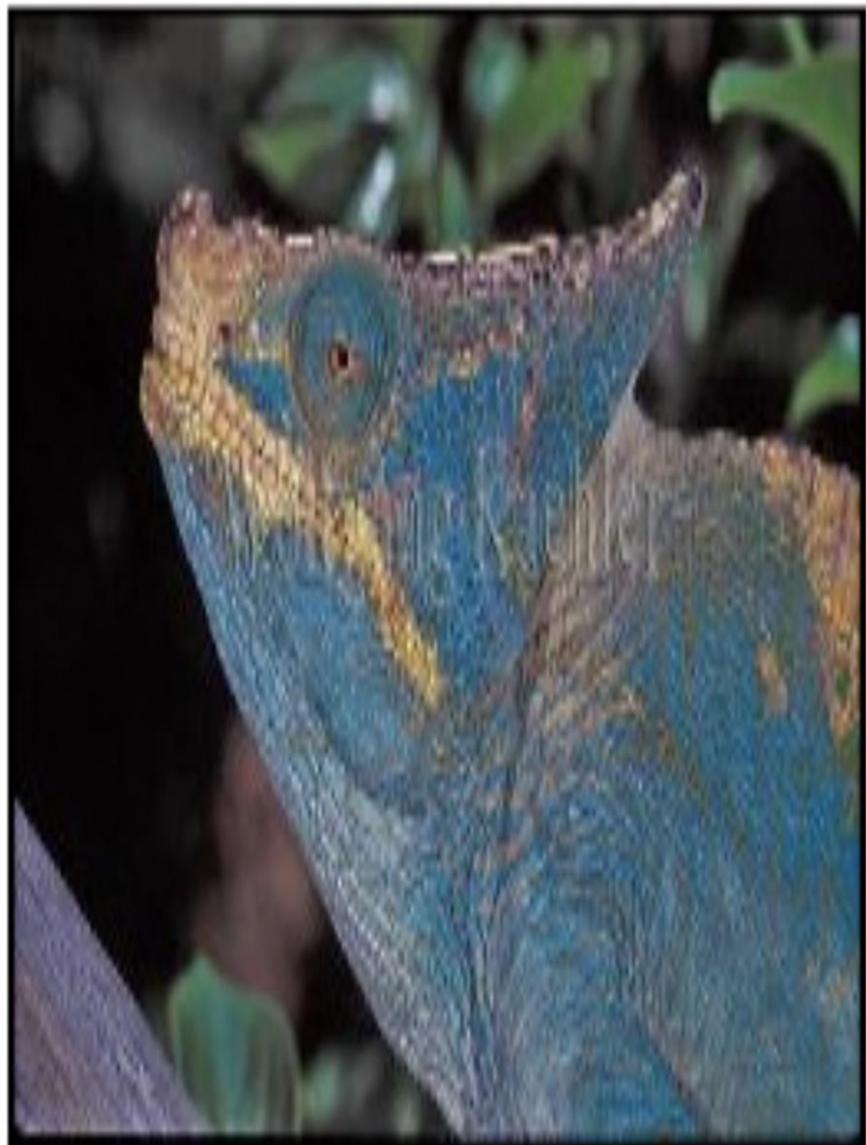
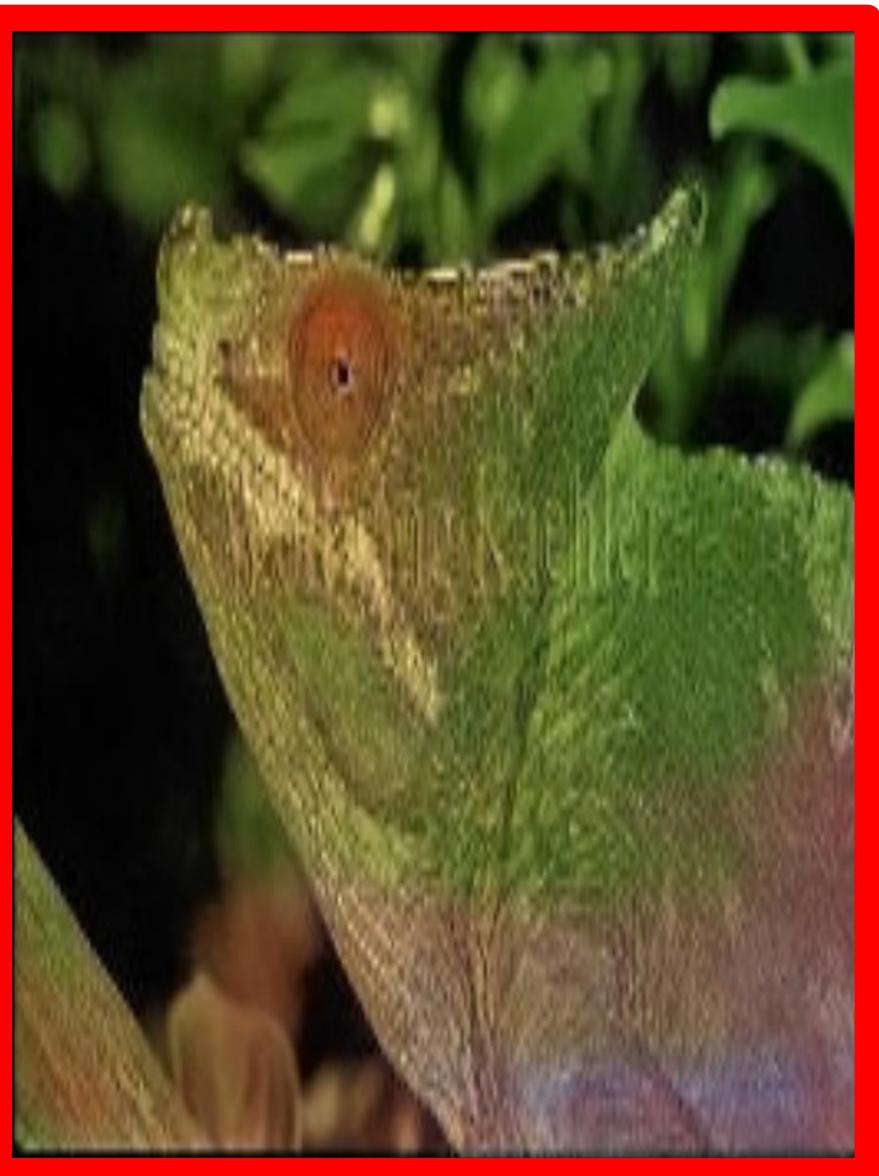


**Fake, 55% fooled**





**Fake, 58% fooled**





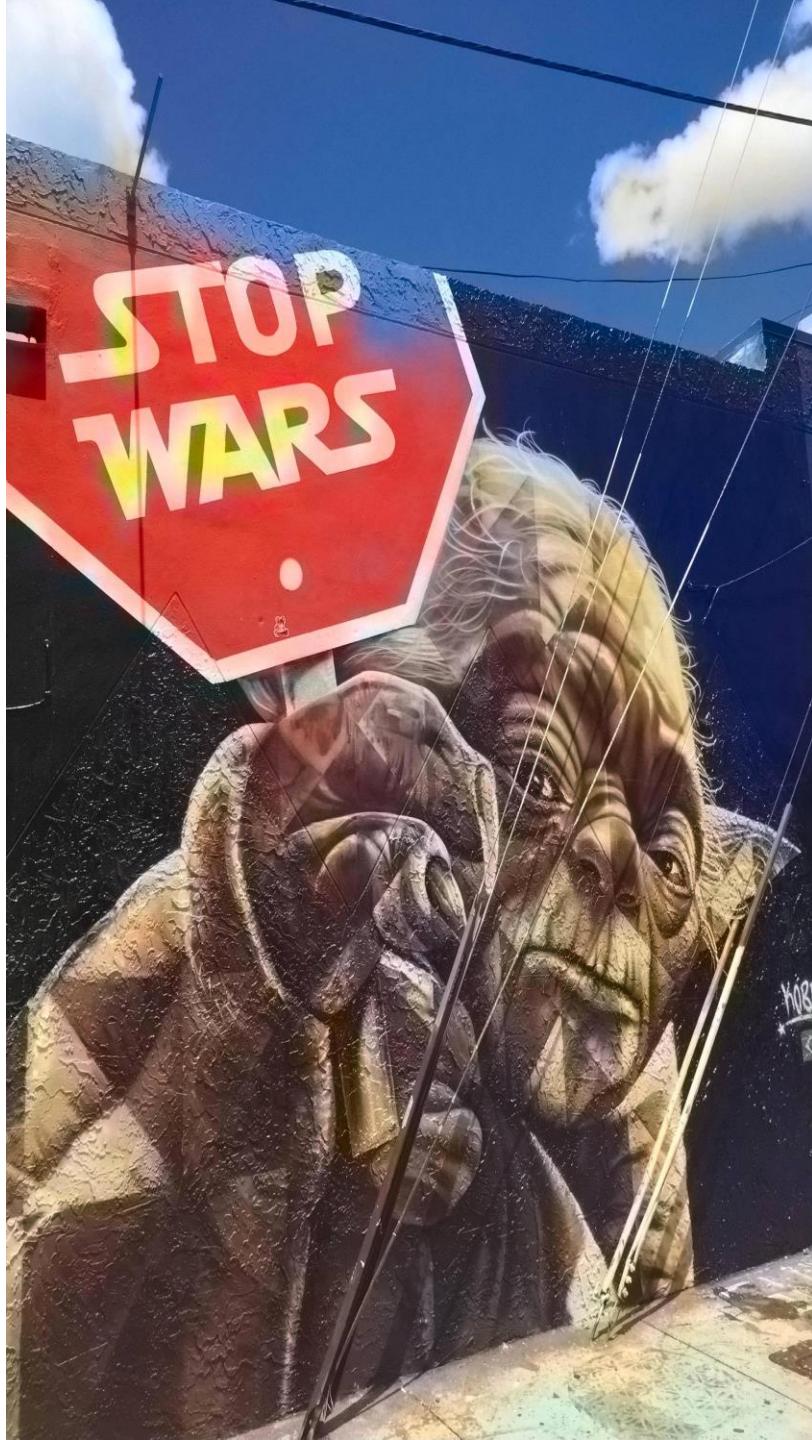
**from Reddit /u/SherySantucci**



**Recolorized by Reddit ColorizeBot**

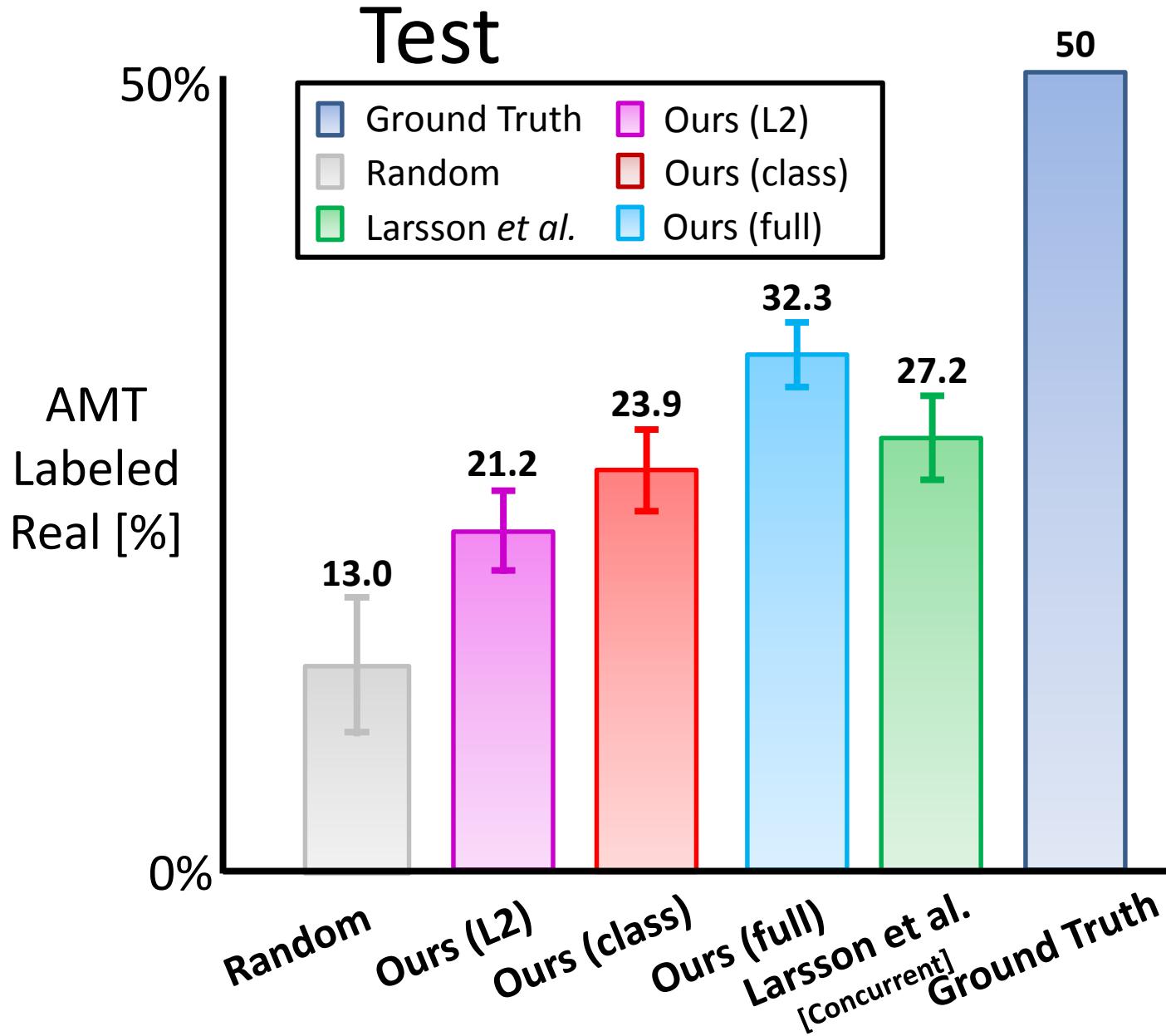


**Photo taken  
by Reddit  
/u/Timteroo,  
Mural from  
street artist  
Eduardo Kobra**



Recolorized by Reddit  
ColorizeBot

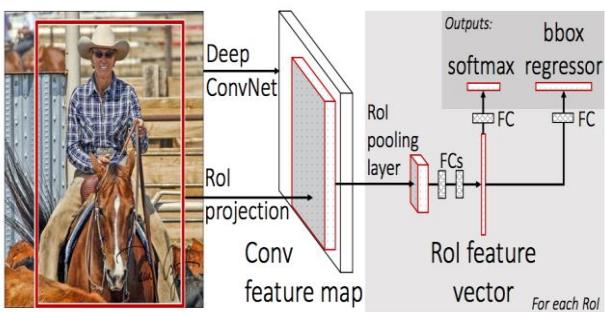
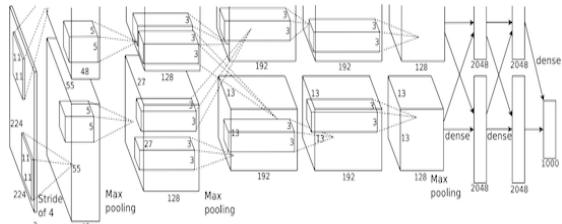
# Perceptual Realism



1600  
images  
tested per  
algorithm

# Dataset & Task Generalization on PASCAL VOC

Does the feature representation  
*transfer* to other datasets and tasks?

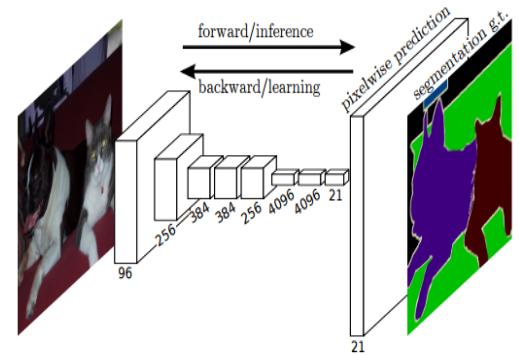


## Classification

Krähenbühl et al. In ICLR,  
2016.

## Detection

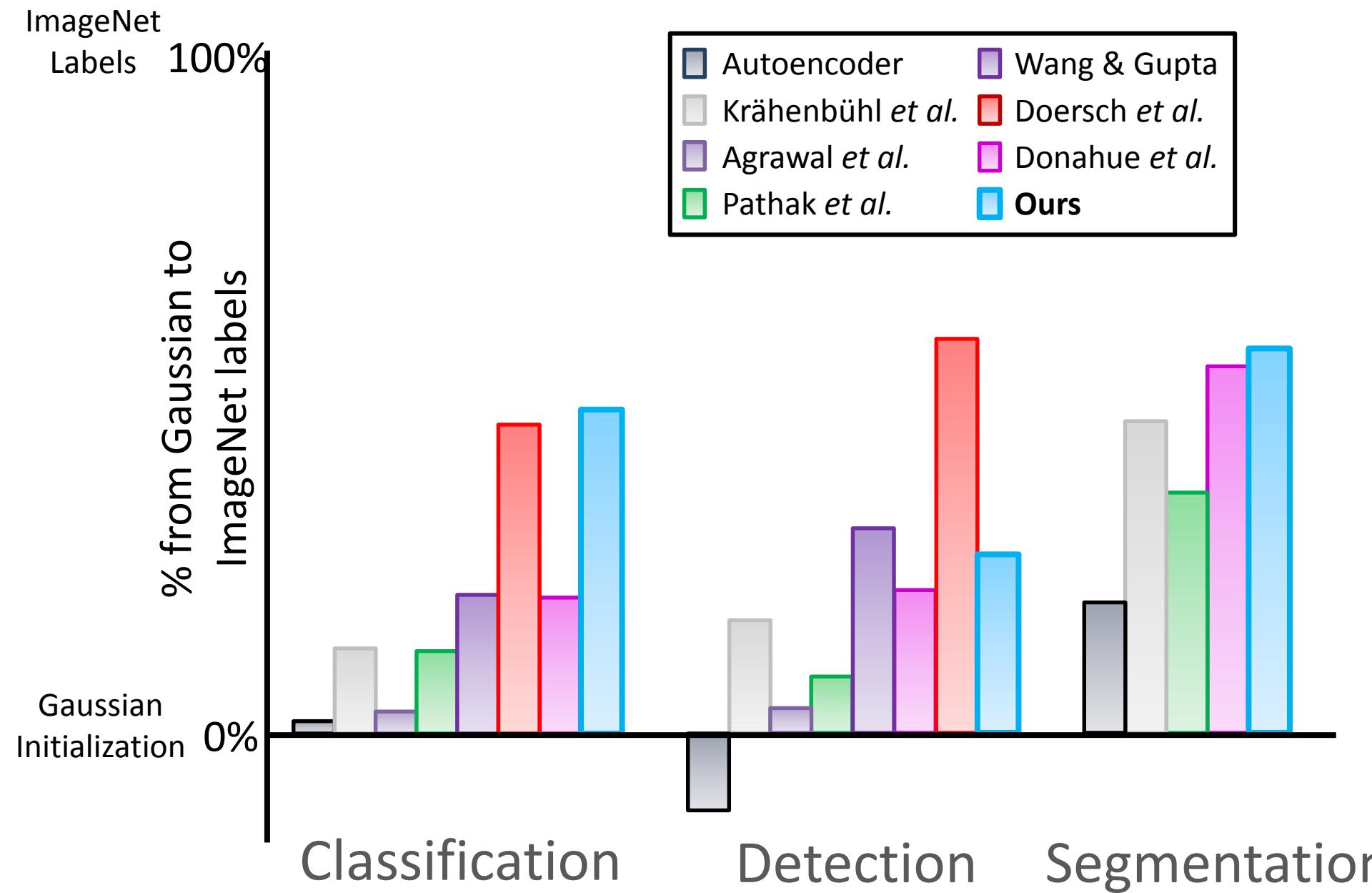
Fast R-CNN. Girshick. In  
ICCV, 2015.



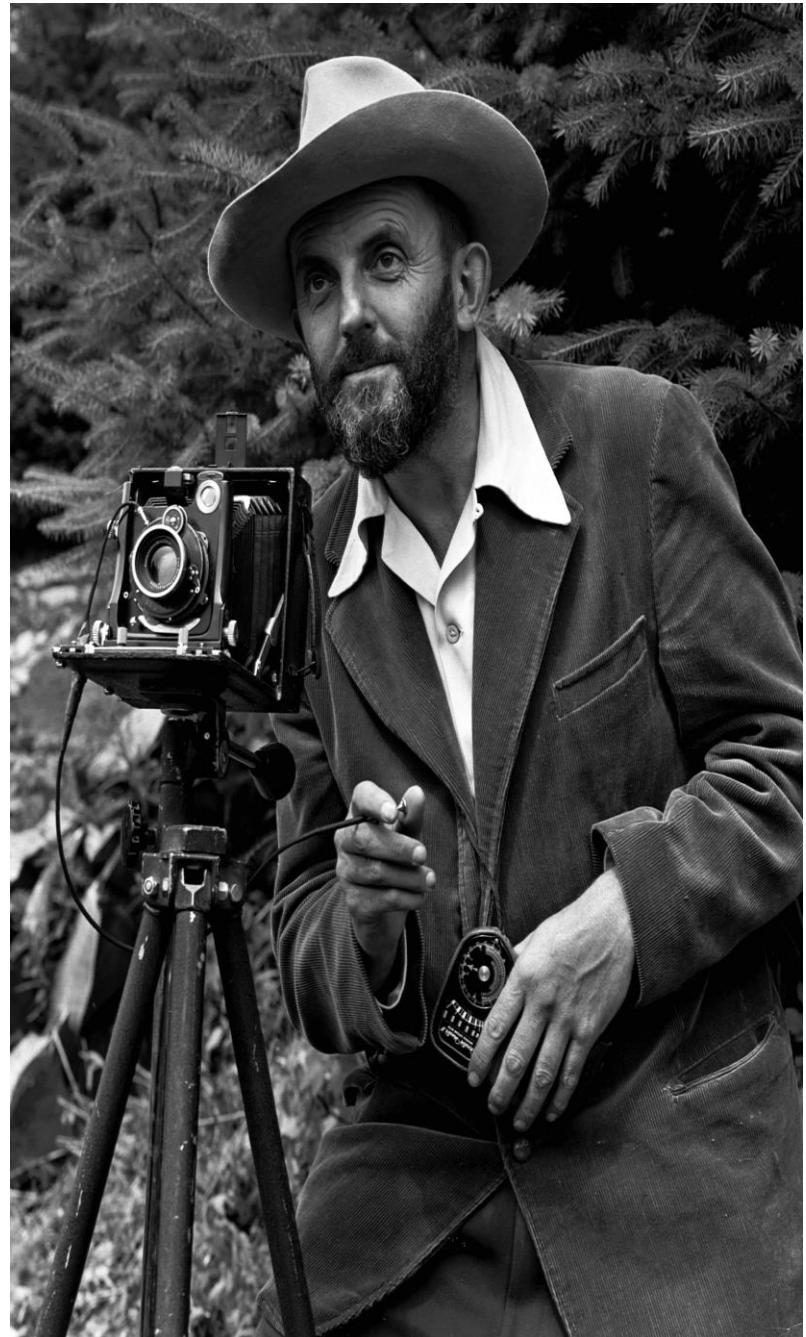
## Segmentation

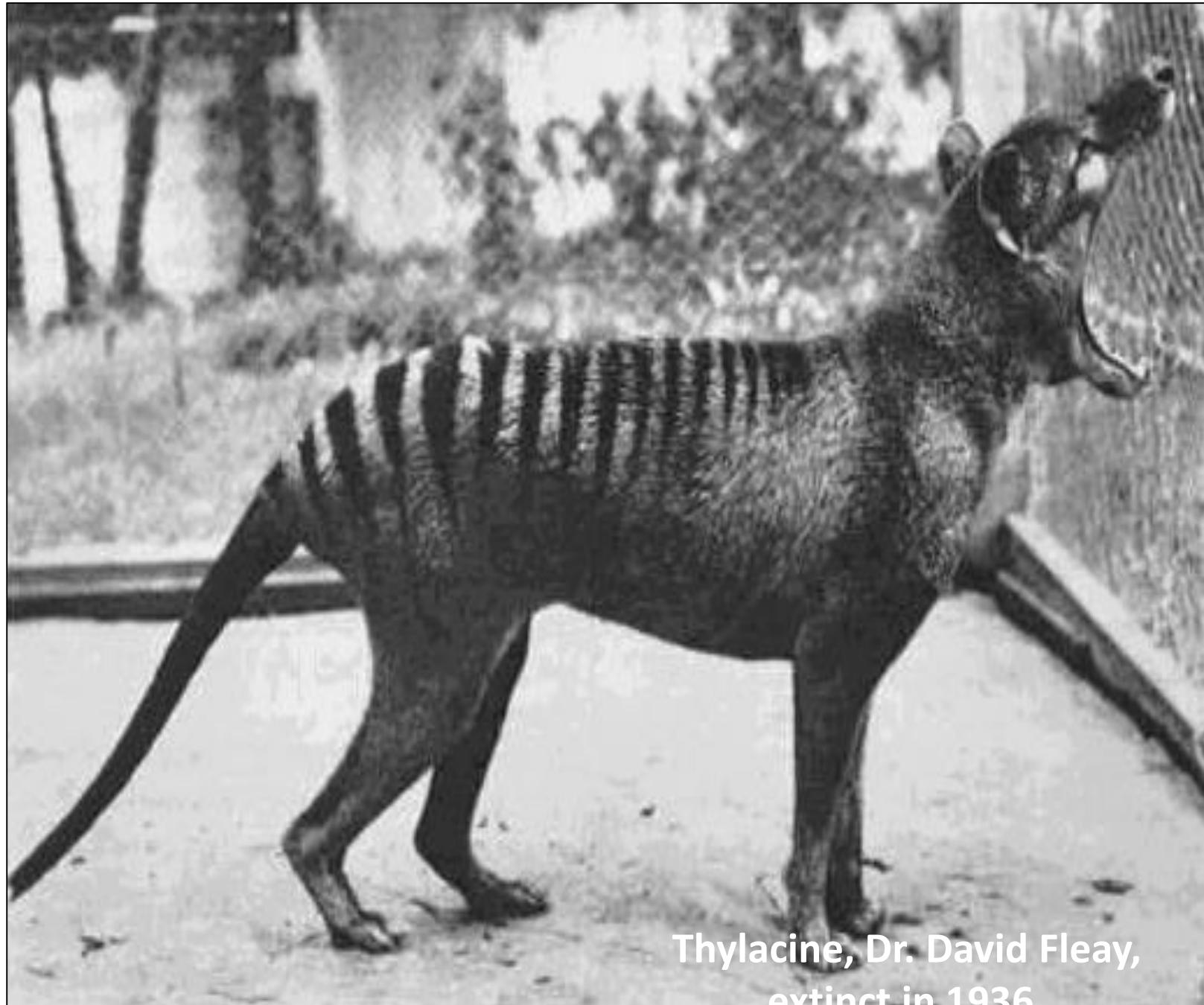
FCNs. Long et al. In CVPR,  
2015.

# Dataset & Task Generalization on PASCAL VOC



Does the  
method work on  
*legacy* black and  
white photos?





Thylacine, Dr. David Fleay,  
extinct in 1936



Thylacine, Dr. David Fleay,  
extinct in 1936



Amateur Family  
Photo 1956



Amateur Family  
Photo 1956



**Henri Cartier-Bresson, Sunday on the Banks of  
the River Seine, 1938.**

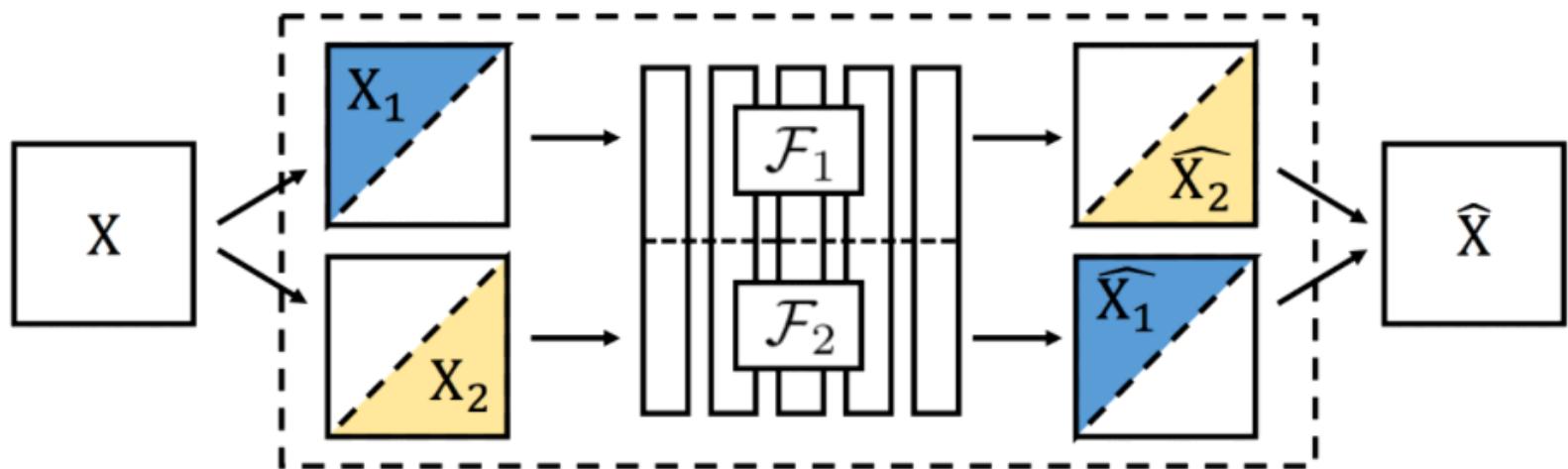


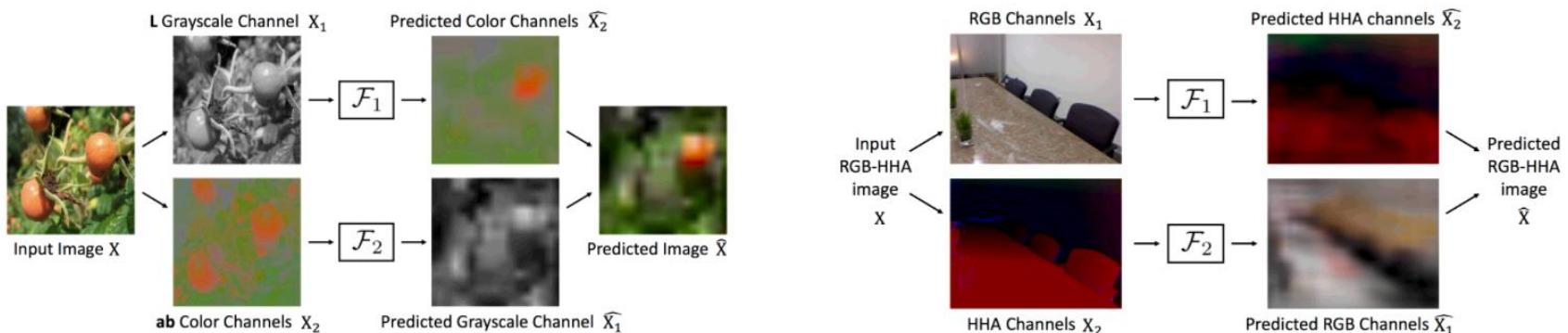
**Henri Cartier-Bresson, Sunday on the Banks of  
the River Seine, 1938.**

# Split-Brain Autoencoders: Unsupervised Learning by Cross-Channel Prediction

Richard Zhang Phillip Isola Alexei A. Efros

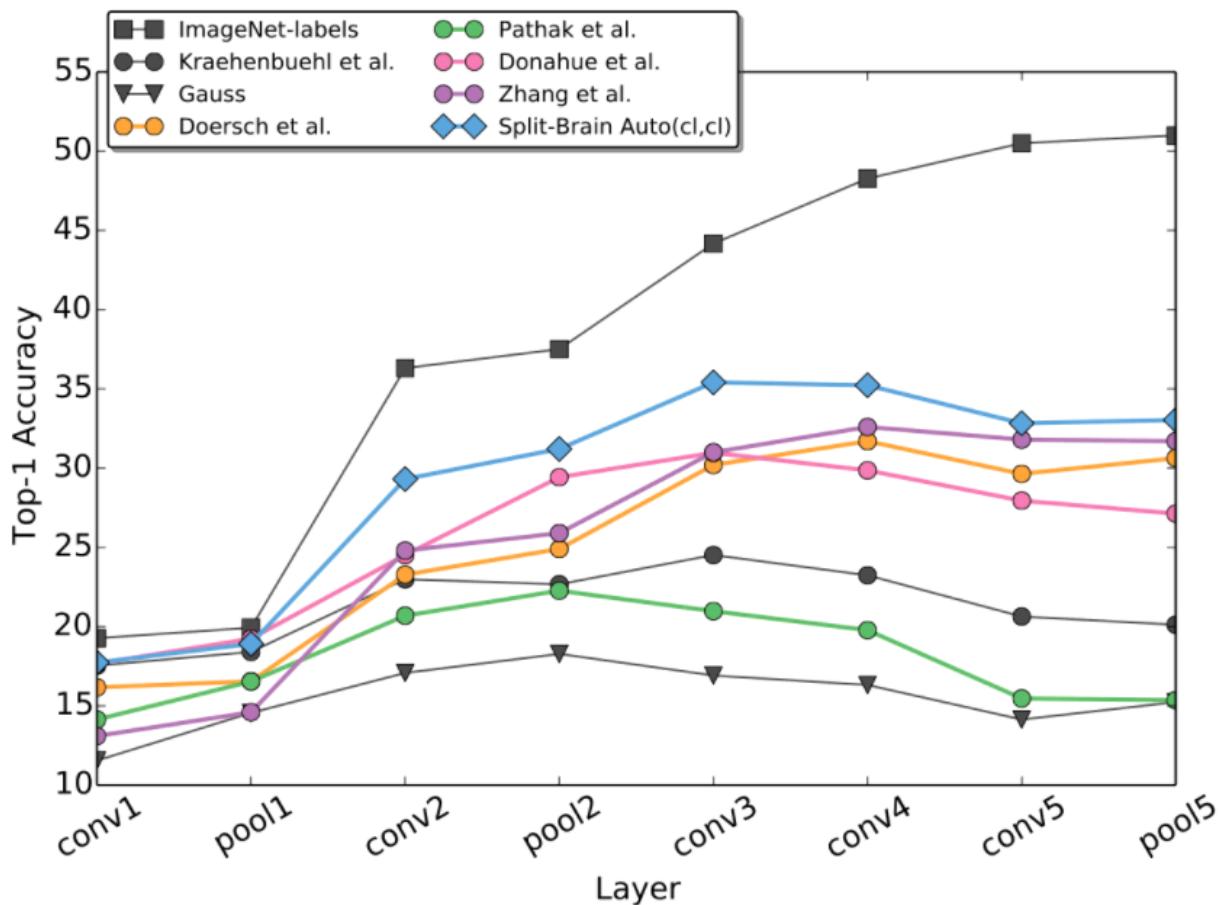
Department of EECS, University of California, Berkeley





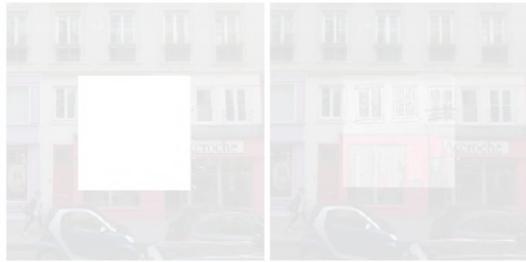
We show that we can learn features in an unsupervised framework, simply by predicting *channels of raw data* from other channels of raw data.

# Feature Evaluation Results



# Self Supervision Examples

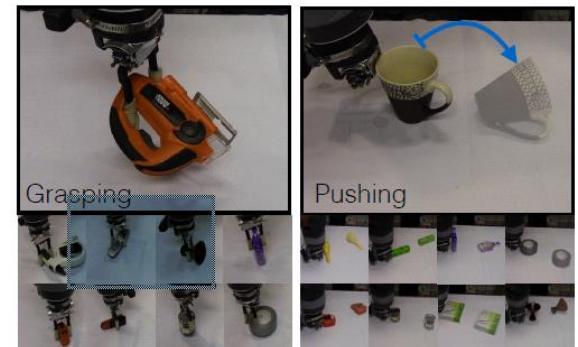
Context



Color

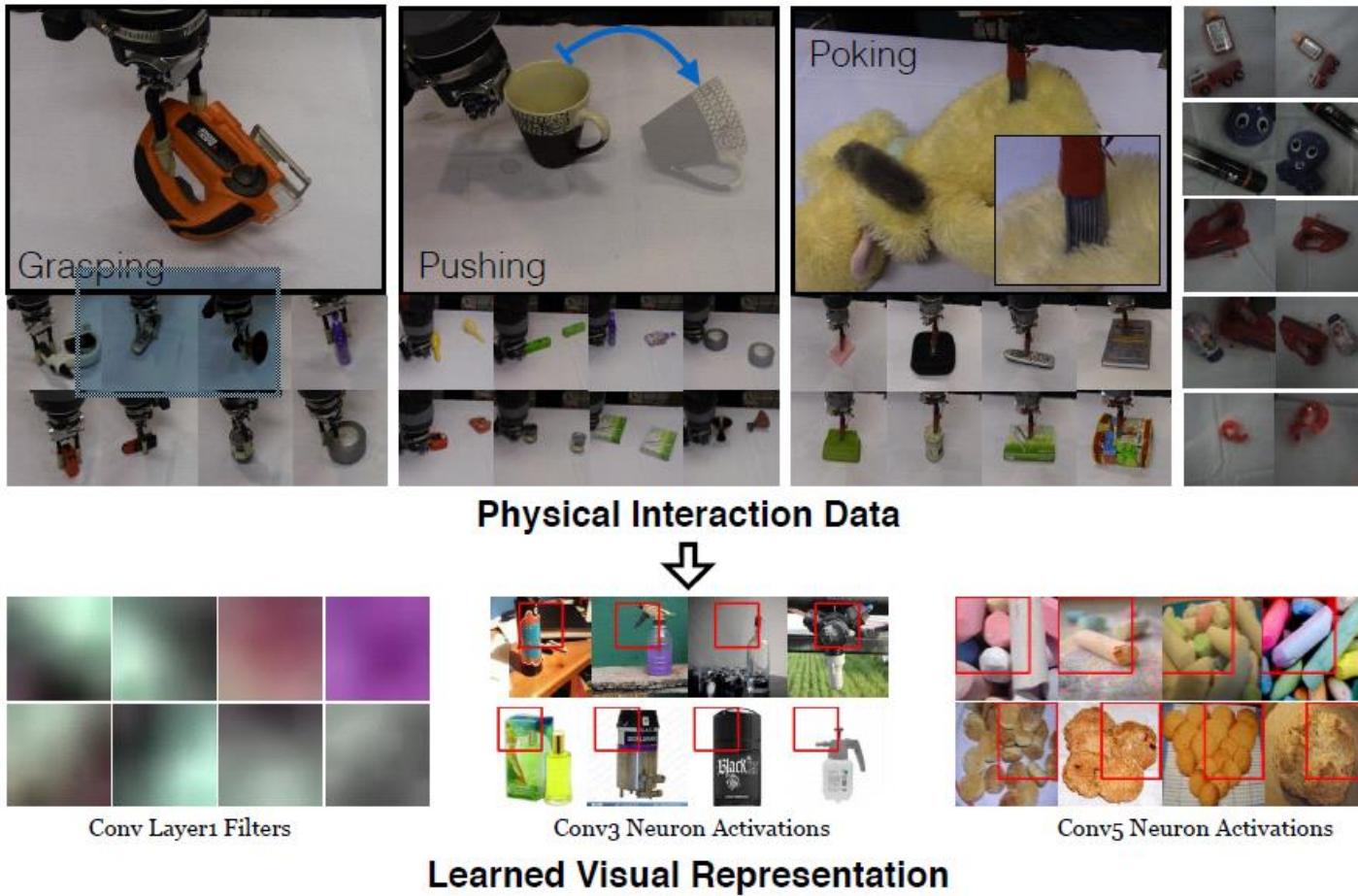


Physical Interactions



# **Using Physical Interactions for Self Supervision**

# The Curious Robot: Learning Visual Representations via Physical Interactions

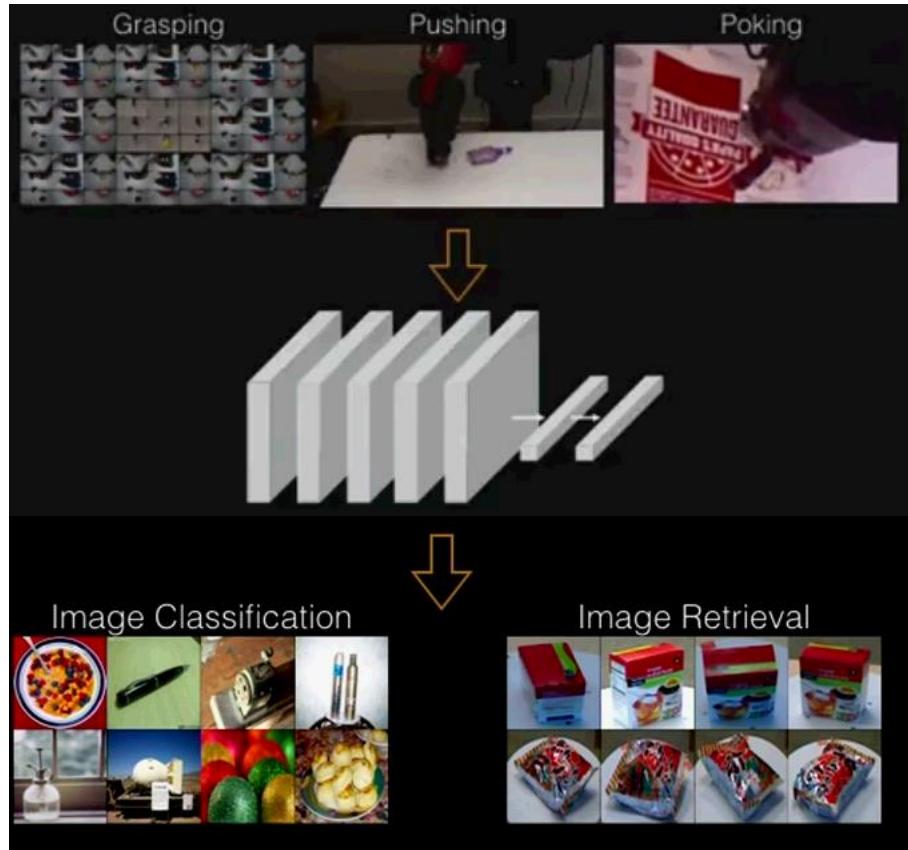


# Learning in Humans

- Use Actions and supervision effects on these actions to train our representations
- Build a similar idea
- Use a robot to learn visual features



# Using robot tasks to learn visual features



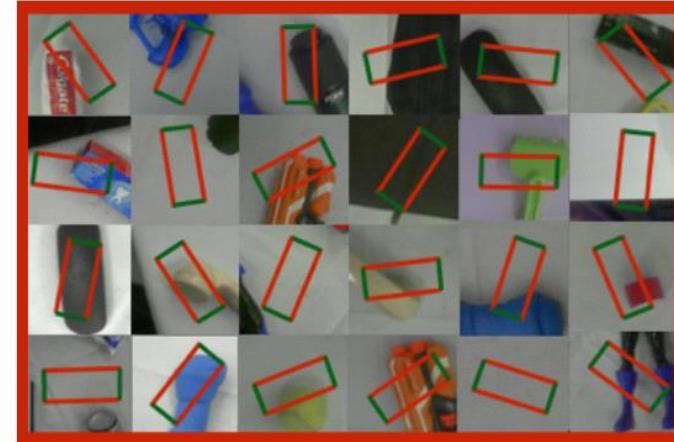
# Grasping Task

- Robot continuously trying to grasp objects
- Grasp configuration lies in 3 dimensions:
  - position of grasp point on the surface of table
  - angle of grasp
- Record successful and unsuccessful grasps

Successful grasps

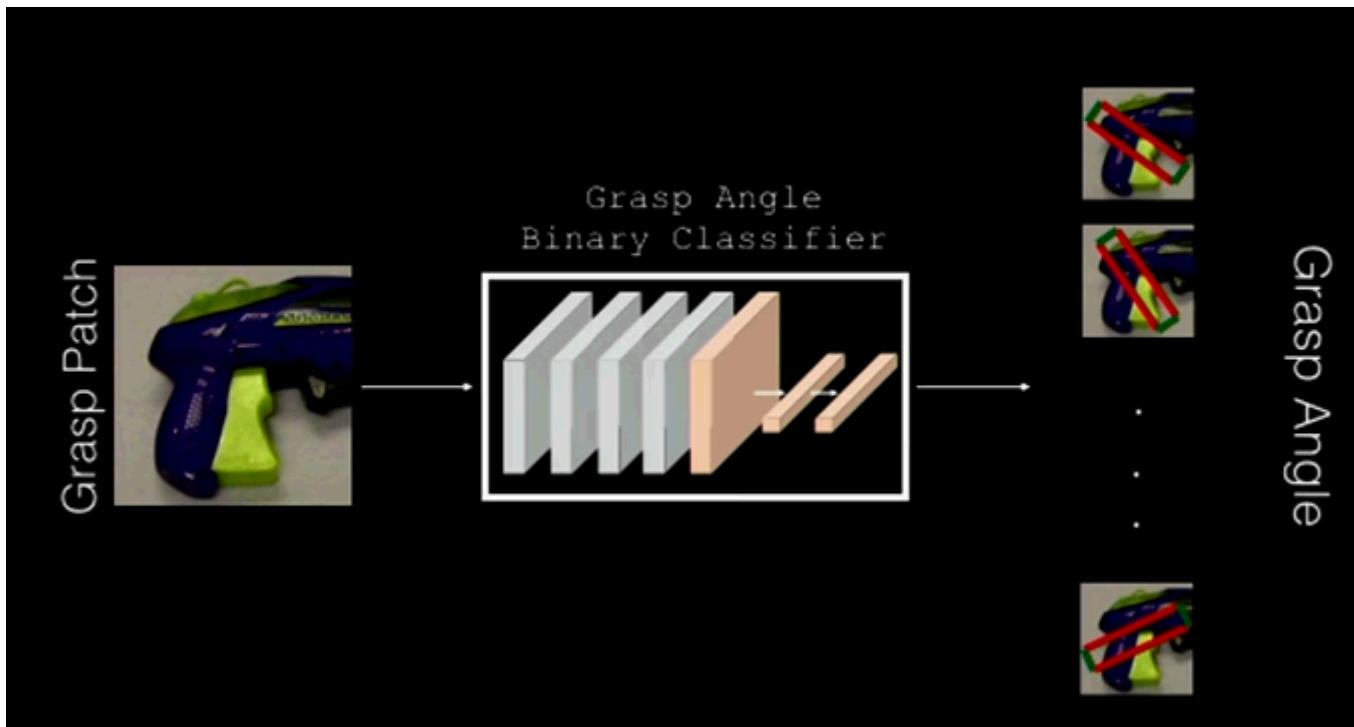


Unsuccessful grasps



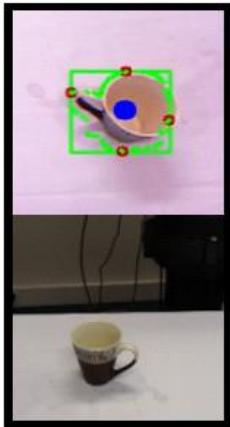
<https://www.youtube.com/watch?v=oSqHc0nLkm8>

# Grasping Formulation

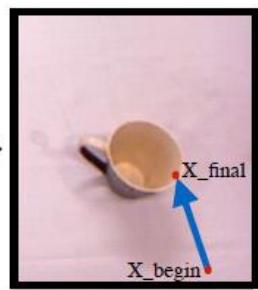


# Pushing Task

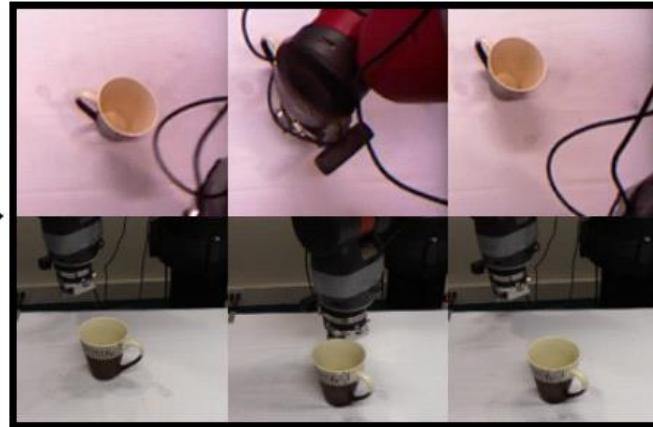
(a) Initial sensing



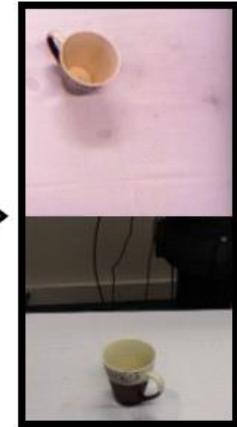
(b) Push select



(c) Plan and execute push action

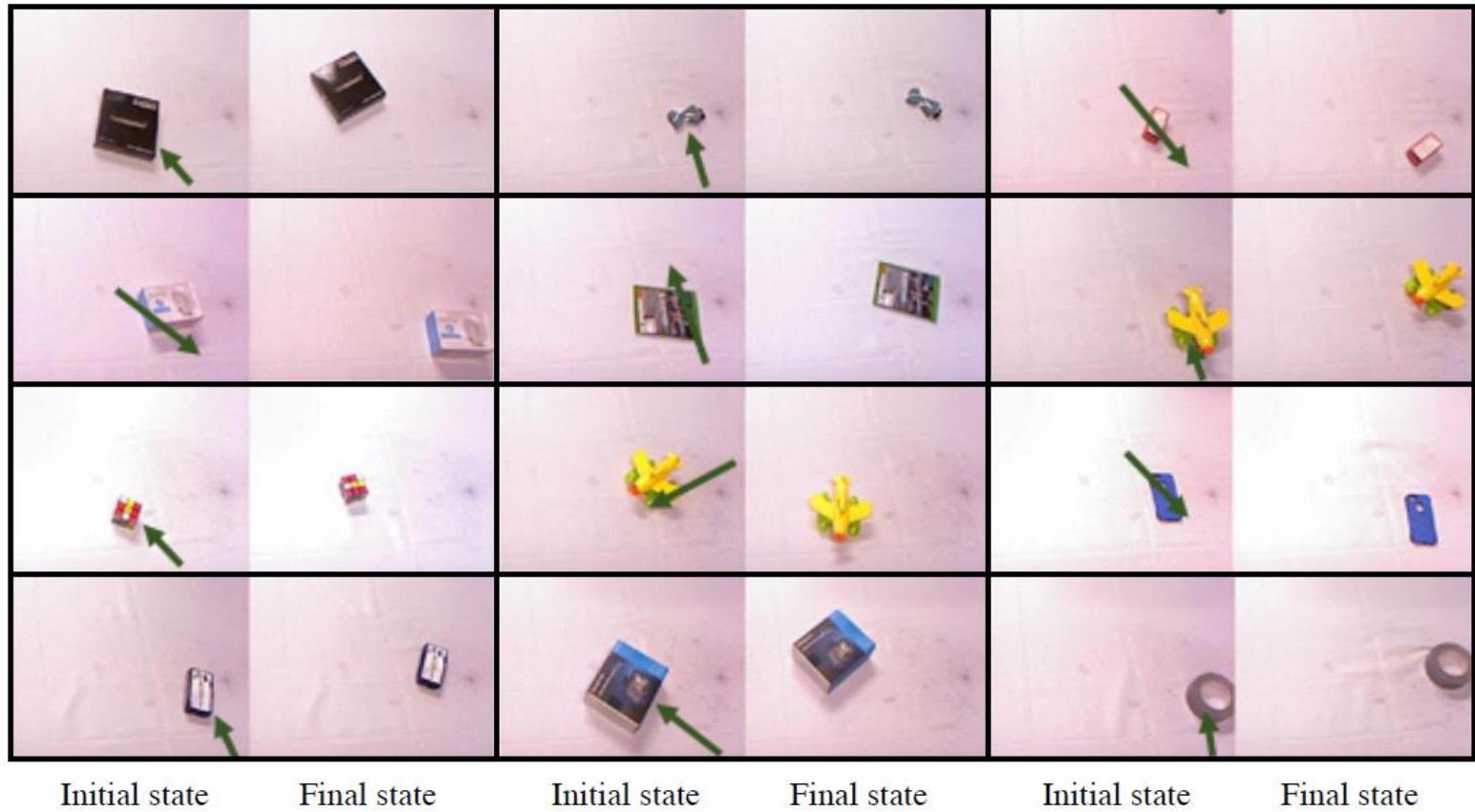


(d) Final sensing

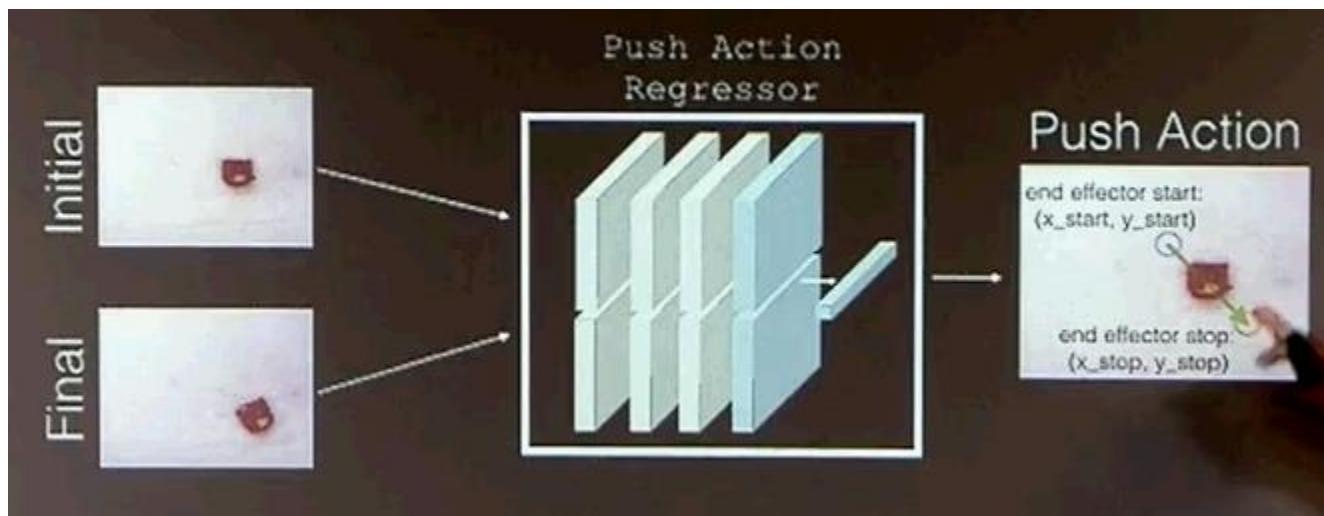


# Pushing Task

Objects and push action pairs



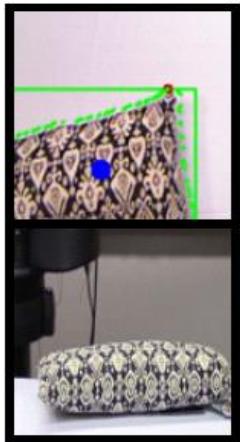
# Planar Pushing Formulation



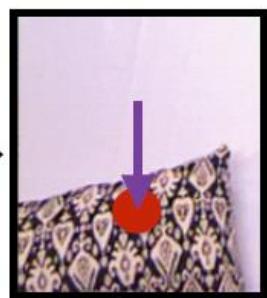
# Poking Task

Learning Visual Representations via Physical Interactions

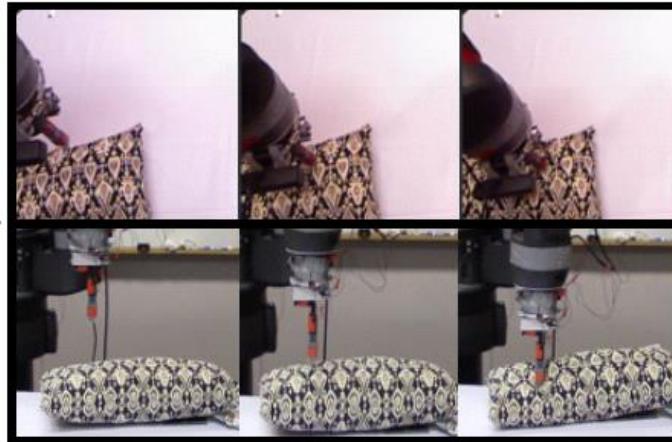
(a) Initial sensing



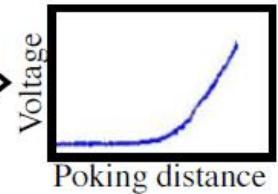
(b) Poke point select



(c) Plan and execute down push action



(d) Tactile sensing

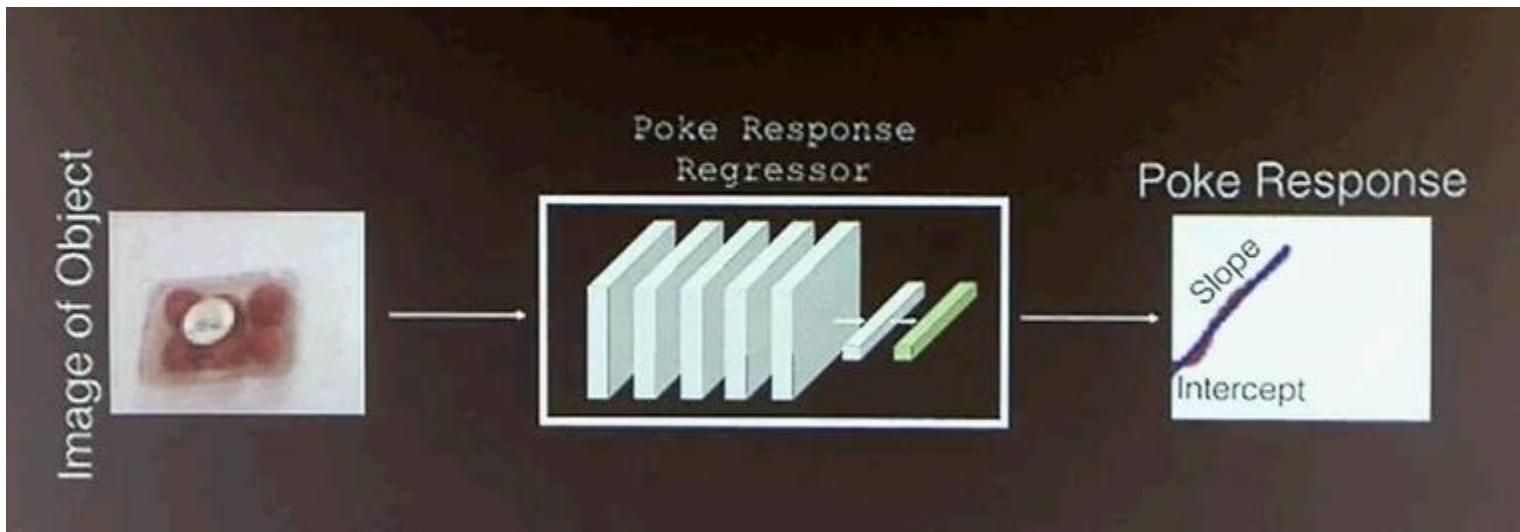


# Poking Task

Objects and poke tactile response pairs



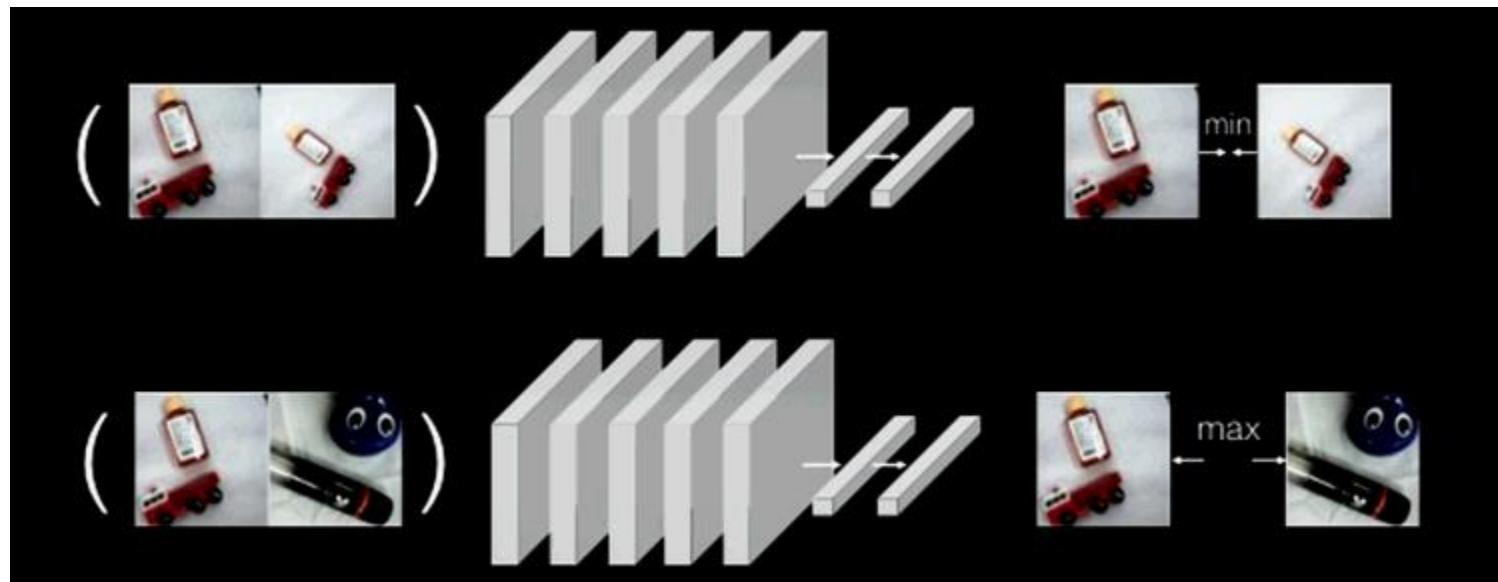
# Poking Formulation



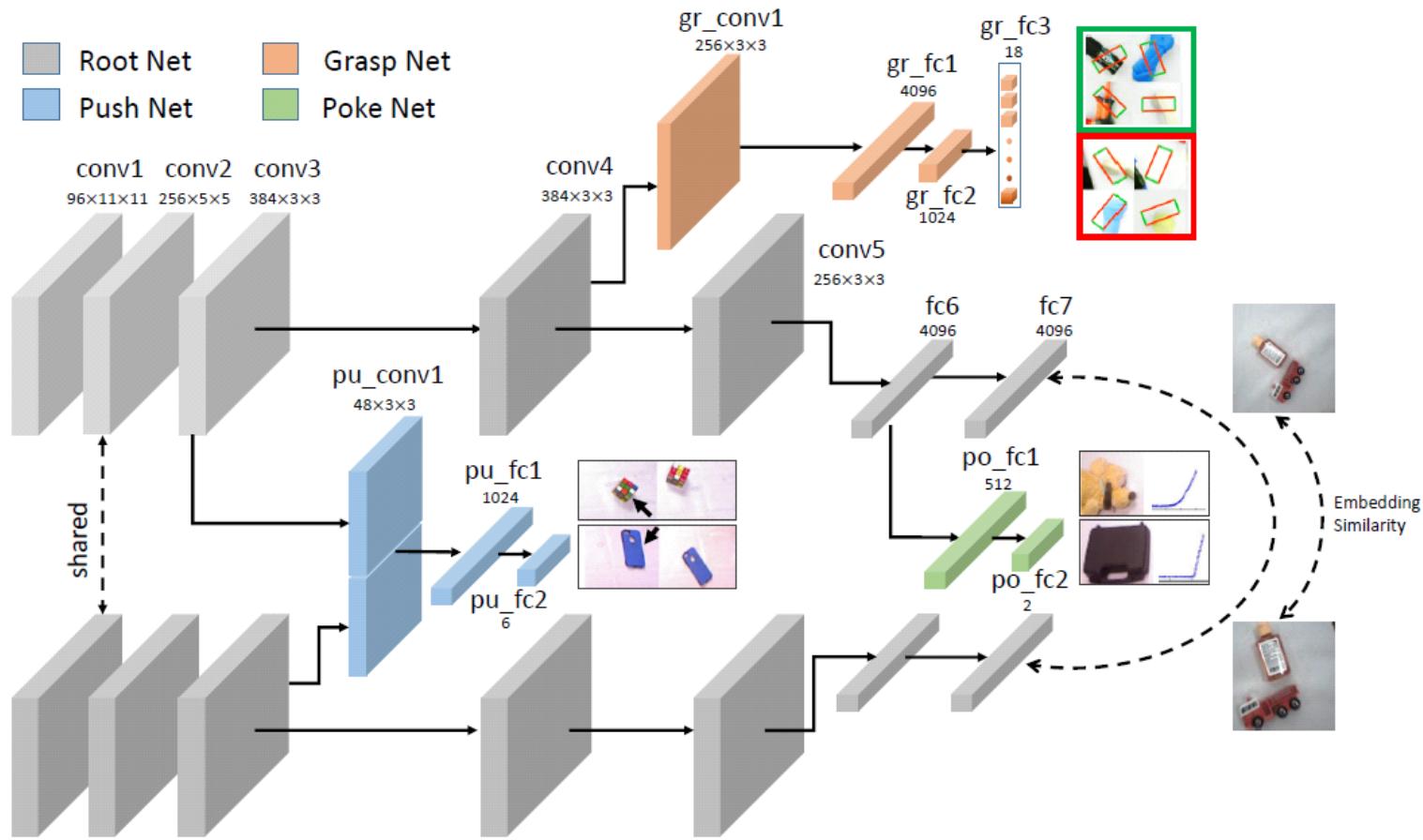
# Pose and Scale Invariance



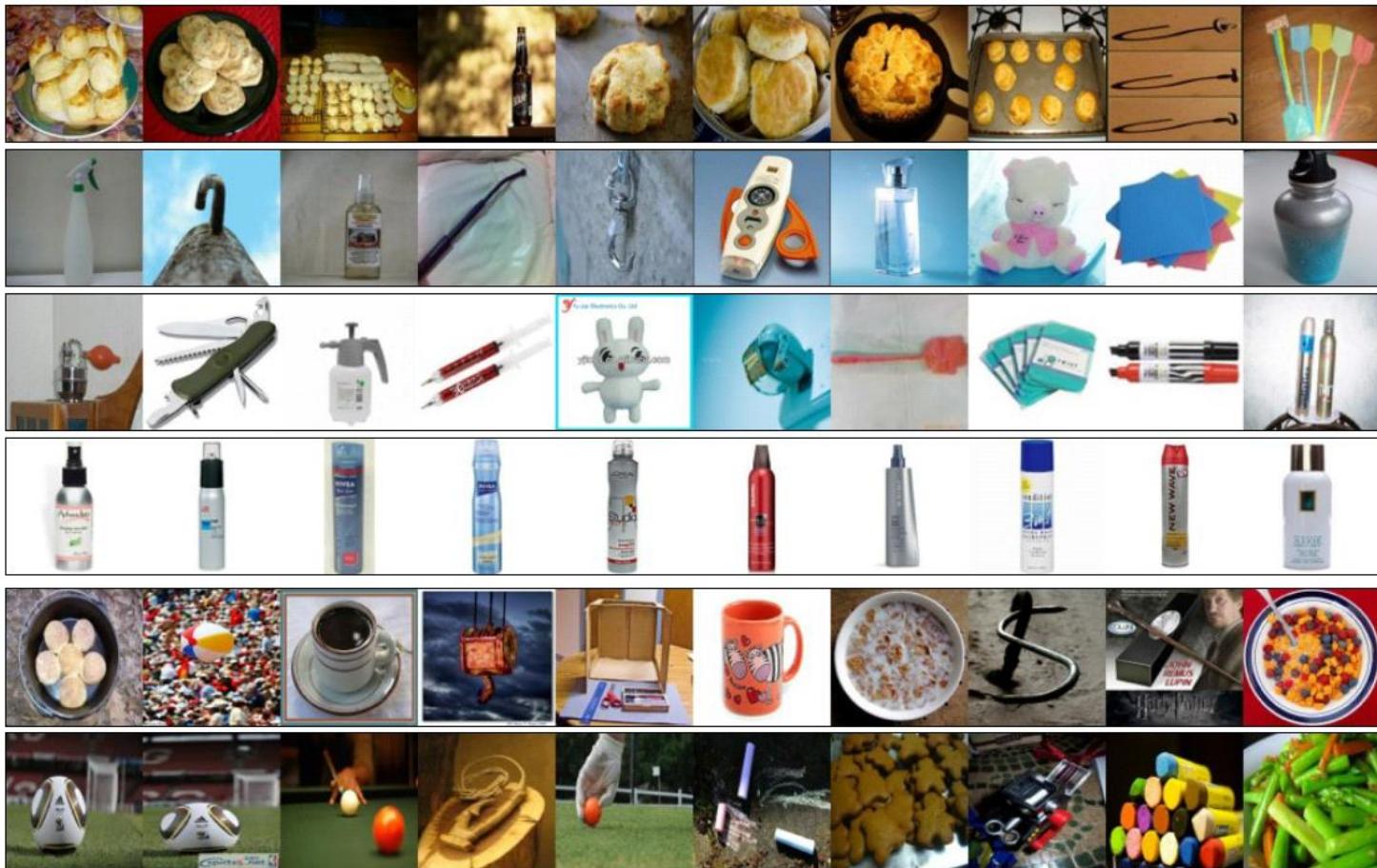
# Pose and Scale Invariance Formulation



# Convolutional Architecture



# Nearest Neighbors



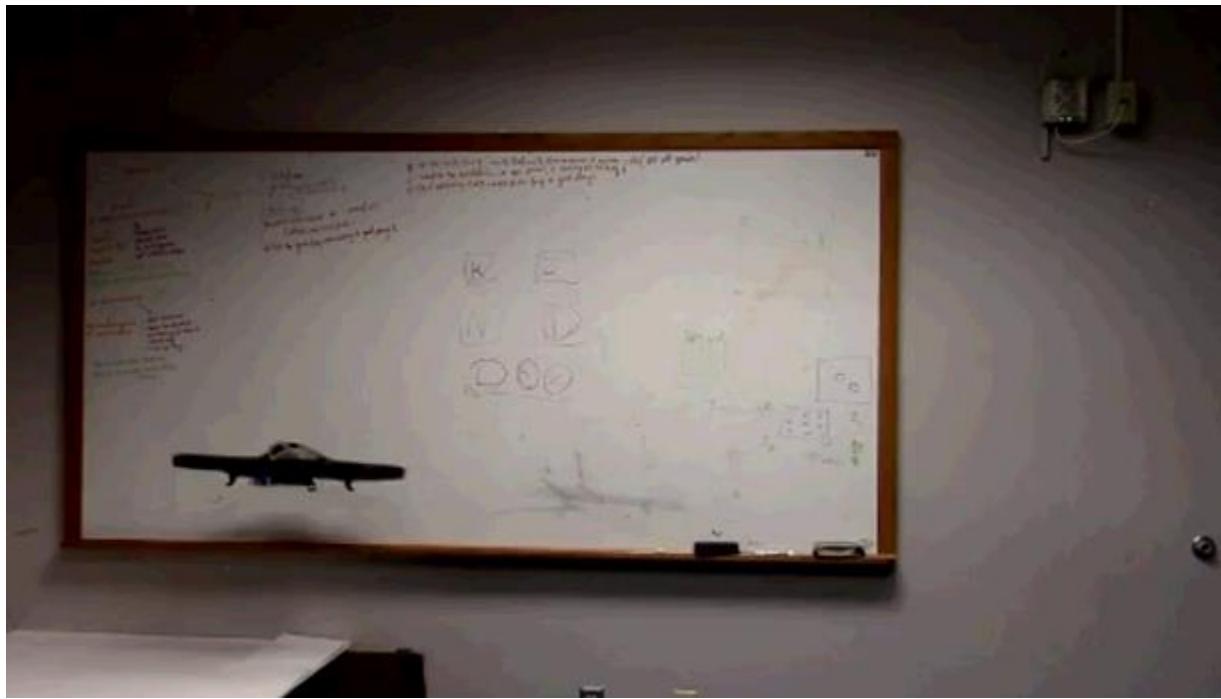
# Classification Accuracy

	Household	UW RGBD	Caltech-256
Root network with random init.	0.250	0.468	0.242
Root network trained on robot tasks ( <b>ours</b> )	0.354	0.693	0.317
AlexNet trained on ImageNet	0.625	0.820	0.656
Root network trained on identity data	0.315	0.660	0.252
Auto-encoder trained on all robot data	0.296	0.657	0.280

# Task Ablation Analysis

	Household	UW RGB-D	Caltech-256
All robot tasks	0.354	0.693	0.317
Except Grasp	0.309	0.632	0.263
Except Push	0.356	0.710	0.279
Except Poke	0.342	0.684	0.289
Except Identity	0.324	0.711	0.297

# Drone: Real World Poking



# Self Supervision Examples

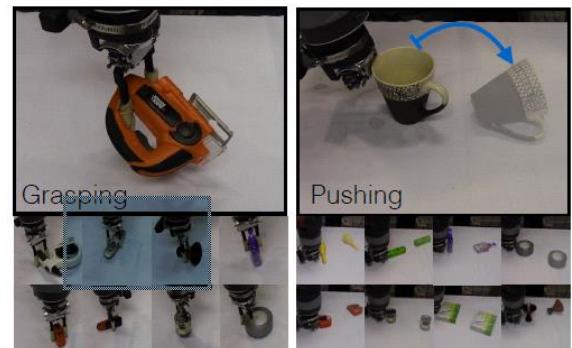
## Context



## Color



## Physical Interactions



Thank You