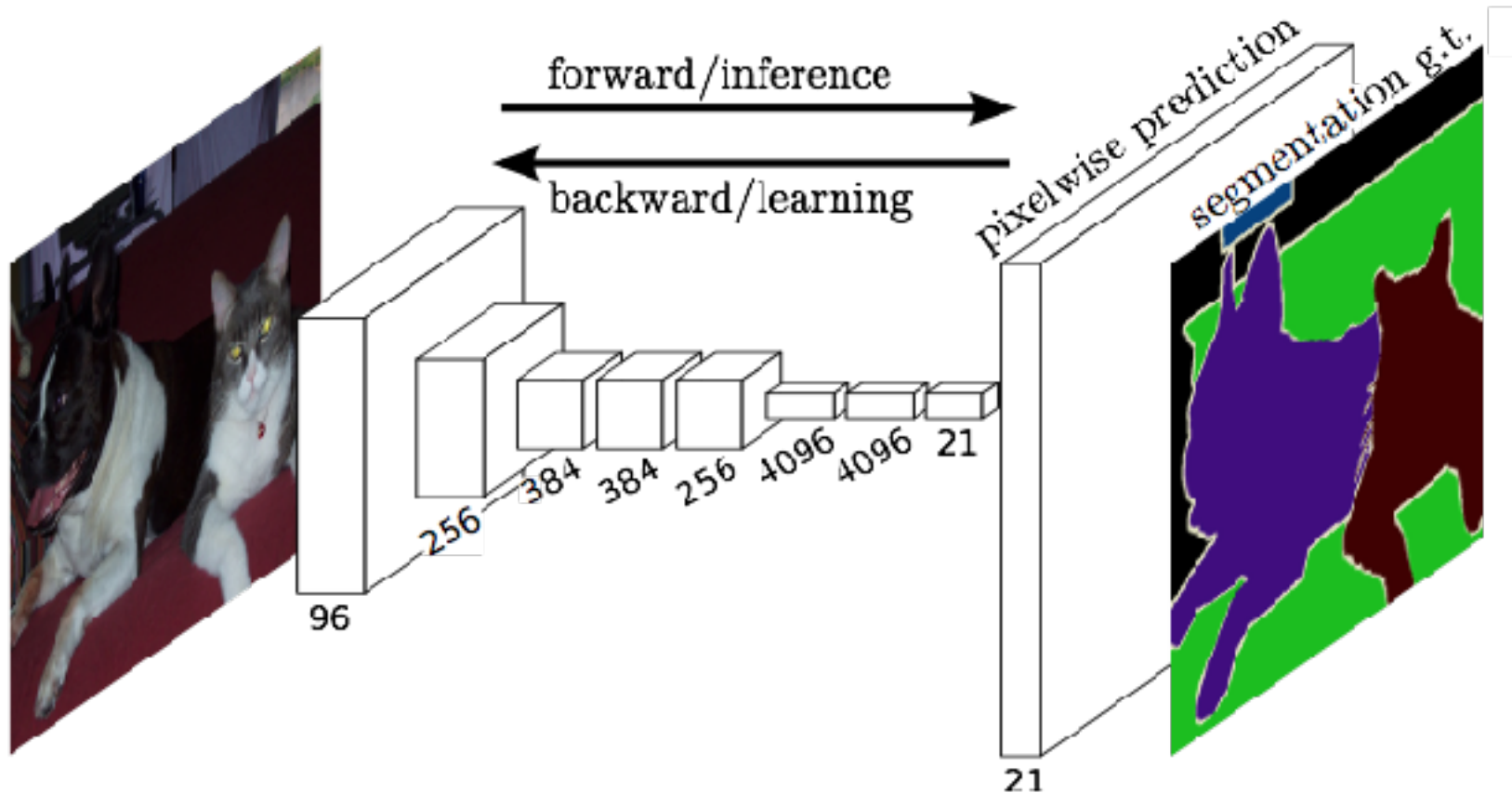# Semantic Segmentation and
# Image Processing

# with Convnets

# Overview

- Methods where output is also an image
  - Fully Convolutional Nets [Long et al., CVPR 2015]
  - Depth, normals and semantic labels from a single image [Eigen ICCV 2015]

- Image processing with Convnets
  - Image colorization [Zhang et al. ECCV 2016]

# A Fuller Understanding of Fully Convolutional Networks

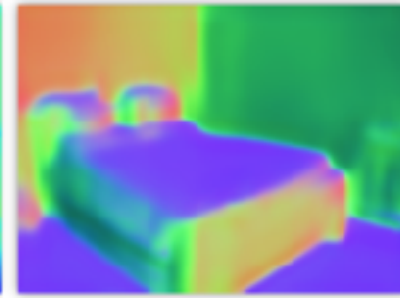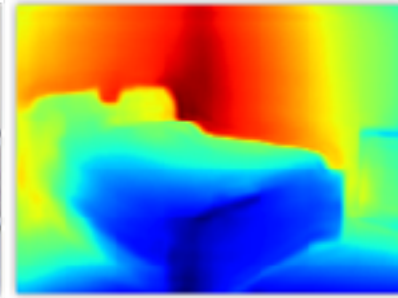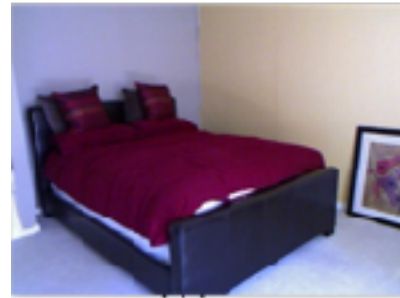Evan Shelhamer*   Jonathan Long*   Trevor Darrell
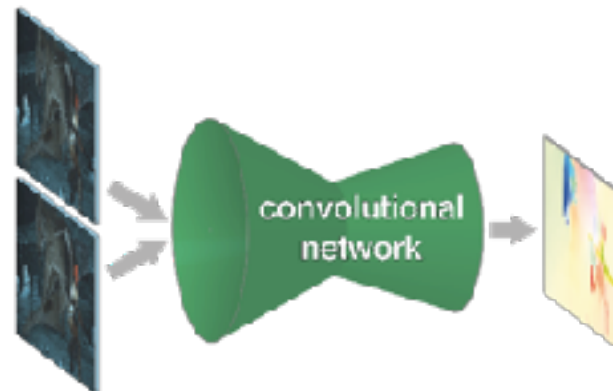
UC Berkeley in CVPR'15, PAMI'16

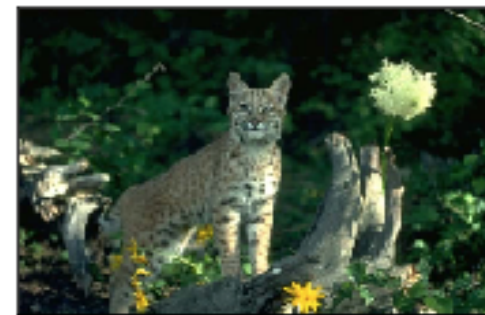# pixels in, pixels out

semantic
segmentation

monocular depth + normals Eigen & Fergus 2015

colorization
Zhang et al.2016

convolutional
network
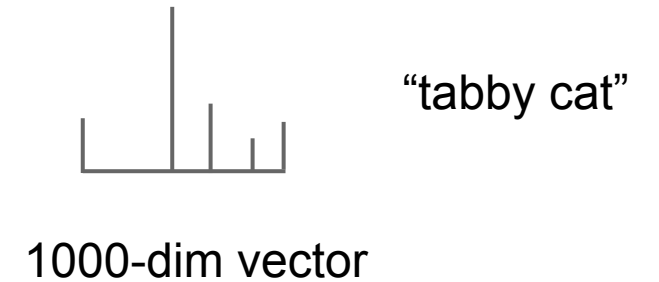
optical flow Fischer et al. 2015

boundary prediction Xie & Tu 2015

# convnets perform classification



< 1 millisecond

1000-dim vector

"tabby cat"

end-to-end learning

# lots of pixels, little time?

~1/10 second

???

end-to-end learning

# a classification network



convolution     fully connected

227 × 227    55 × 55    27 × 27    13 × 13

"tabby cat"

# becoming fully convolutional

convolution

227 × 227    55 × 55    27 × 27    13 × 13    1 × 1

# becoming fully convolutional



convolution

H × W    H/4 × W/4    H/8 × W/8    H/16 × W/16    H/32 × W/32

# upsampling output



convolution

H × W     H/4 × W/4    H/8 × W/8     H/16 × W/16      H/32 × W/32     H × W

# end-to-end, pixels-to-pixels network

# end-to-end, pixels-to-pixels network

# spectrum of deep features

## combine where (local, shallow) with what (global, deep)



image

intermediate layers

fuse features into deep jet

(cf. Hariharan et al. CVPR15 "hypercolumn")

# skip layers



interp + sum

2x conv7
pool4

interp + sum

skip to fuse layers!

4x conv7
2x pool4
pool3

end-to-end, joint learning
of semantics and location

image  conv1  pool1  conv2  pool2  conv3  pool3  conv4  pool4  conv5  pool5  conv6-7

dense output

14

# skip layer refinement



| input image | stride 32 | stride 16 | stride 8 | ground truth |
|---|---|---|---|---|
| | no skips | 1 skip | 2 skips | |

# skip FCN computation



- Stage 1 (60.0ms)
- Stage 2 (18.7ms)
- Stage 3 (23.0ms)

A multi-stream network that fuses features/predictions across layers

| FCN | SDS* | Truth | Input |
|---|---|---|---|



Relative to prior state-of-the-art SDS:

30% relative improvement for mean IoU

286× faster

*Simultaneous Detection and Segmentation Hariharan et al. ECCV14

# SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation



Max pooling indices transferred to decoder to improve output resolution

https://arxiv.org/abs/1511.00561

# UNet: Convolutional Networks for Biomedical Image Segmentation



Segmentation of a 512x512 image takes less than a second on a recent GPU

https://arxiv.org/abs/1505.04597

# Further Resources

http://blog.qure.ai/notes/semantic-segmentation-deep-learning-review

# Overview

- Methods where output is now an image
  - Fully Convolutional Nets [Long et al., CVPR 2015]
  - Depth, normals and semantic labels from a single image [Eigen ICCV 2015]

- Image processing with Convnets
  - Image colorization [Zhang et al. ECCV 2016]

# Beyond Object Classification with Convolutional Networks

David Eigen (NYU -> Clarifai)

Rob Fergus (Facebook / NYU)

# Motivation



Input Image



Semantic Map

- Understand input scene
  - Semantic
  - Geometric

# Motivation



Input Image

Semantic Map

Depth

- Understand input scene
  - Semantic
  - Geometric

# Motivation



Input Image

Depth

Se...

Normals

- Understand input scene
  - Semantic
  - Geometric

# Motivation



Input Image

Depth

Se

Normals

- **<u>Predict Pixel Maps from a Single Image</u>**

# Architecture

Input: 320x240

Output 1: 19x14

# Architecture

Input: 320x240



96 256 384 384 256 4096 64

upsample

64 128 64 64

Output 2: 75x55

# Architecture

Input: 320x240



96  256  384  384  256  4096  64

upsample

64  128  64  64

upsample

63  64  64  64

conv+pool  concat  convolutions

Output: 147x109

# Architecture

Input: 320x240



**upsample**

**upsample**

conv+pool　　**concat**　　convolutions

# Architecture

Input: 320x240



96 256 384 384 256 4096 64

upsample

64 128 64 64

upsample

64 64+C 64 64

conv+pool  concat  convolutions

# Losses

Depth: $$d = D - D^*$$ D = log predicted depth, D* = log true depth

$$L_{depth}(D, D^*) = \frac{1}{n}\sum_i d_i^2 - \frac{1}{2n^2}\left(\sum_i d_i\right)^2 + \frac{1}{n}\sum_i [(\nabla_x d_i)^2 + (\nabla_y d_i)^2]$$

Norm

Labels

# Training

- Pre-train Alexnet/VGGnet scale 1 with Imagenet
- Scale 2 & 3 random initialization
- Joint train layers 1 & 2 for each task
  - Loss on output of layer 2

- Fix layers 1 & 2, train layer 3

- For depth & normals task, share scale 1
  - But separate scale 2 & 3's
  - 1.6x speedup

# Evaluation



- NYU Depth dataset
  - RGB, Depth
    and per-pixel labels
  - Indoor scenes

- Supervised training
  of models

- Compare to range of other methods
  - Also on SIFTFlow and PASCAL VOC'11

# Depths Comparison



| | Eigen NIPS'14 (2 scales) | Ours | Ground Truth |

# Depth Comparison

- m3d = Make3D [Saxena & Ng 2006]



| | input | m3d | Ours (2-scale) | ground truth |

| **Depth Prediction** | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | Ladicky[20] | Karsch[14] | Baig [1] | Liu [18] | Eigen[4] | Ours(A) | Ours(VGG) |
| $\delta < 1.25$ | 0.542 | – | 0.597 | 0.614 | 0.614 | 0.697 | **0.769** |
| $\delta < 1.25^2$ | 0.829 | – | – | 0.883 | 0.888 | 0.912 | **0.950** |
| $\delta < 1.25^3$ | 0.940 | – | – | 0.971 | 0.972 | 0.977 | **0.988** |
| abs rel | – | 0.350 | 0.259 | 0.230 | 0.214 | 0.198 | **0.158** |
| sqr rel | – | – | – | – | 0.204 | 0.180 | **0.121** |
| RMS(lin) | – | 1.2 | 0.839 | 0.824 | 0.877 | 0.753 | **0.641** |
| RMS(log) | – | – | – | – | 0.283 | 0.255 | **0.214** |
| sc-inv. | – | – | 0.242 | – | 0.219 | 0.202 | **0.171** |

# Surface Normals

# Surface Normals

| Surface Normal Estimation (GT [6]) | | | | | |
|---|---|---|---|---|---|
| | **Angle Distance** | | **Within $t^\circ$ Deg.** | | |
| | **Mean** | **Median** | $11.25^\circ$ | $22.5^\circ$ | $30^\circ$ |
| 3DP [6] | 34.2 | 30.0 | 18.5 | 38.6 | 50.0 |
| Ladicky &al [16] | 32.5 | 22.3 | 27.4 | 50.2 | 60.1 |
| Fouhey &al [7] | 35.1 | 19.2 | 37.6 | 53.3 | 58.9 |
| Wang &al [33] | 26.6 | 15.3 | 40.1 | 61.4 | 69.0 |
| Ours (AlexNet) | 23.1 | 15.1 | 39.4 | 63.6 | 72.7 |
| Ours (VGG) | **20.5** | **13.2** | **44.0** | **68.5** | **77.2** |

| Surface Normal Estimation (GT [27]) | | | | | |
|---|---|---|---|---|---|
| | **Angle Distance** | | **Within $t^\circ$ Deg.** | | |
| | **Mean** | **Median** | $11.25^\circ$ | $22.5^\circ$ | $30^\circ$ |
| 3DP [6] | 37.7 | 34.1 | 14.0 | 32.7 | 44.1 |
| Ladicky &al [16] | 35.5 | 25.5 | 24.0 | 45.6 | 55.9 |
| Wang &al [33] | 28.8 | 17.9 | 35.2 | 57.1 | 65.5 |
| Ours (AlexNet) | 25.9 | 18.2 | 33.2 | 57.5 | 67.7 |
| Ours (VGG) | **22.2** | **15.3** | **38.6** | **64.0** | **73.9** |

# Results: Normals

**Angle from Ground Truth**

# Output from each scale



input

depth

far      near

coarse ----→ fine

normals

# Semantic Labels: NYUD



| RGB input | 4-Class Prediction | 13-Class Prediction | 13-Class Ground Truth |

# Results: NYUD 40 Classes

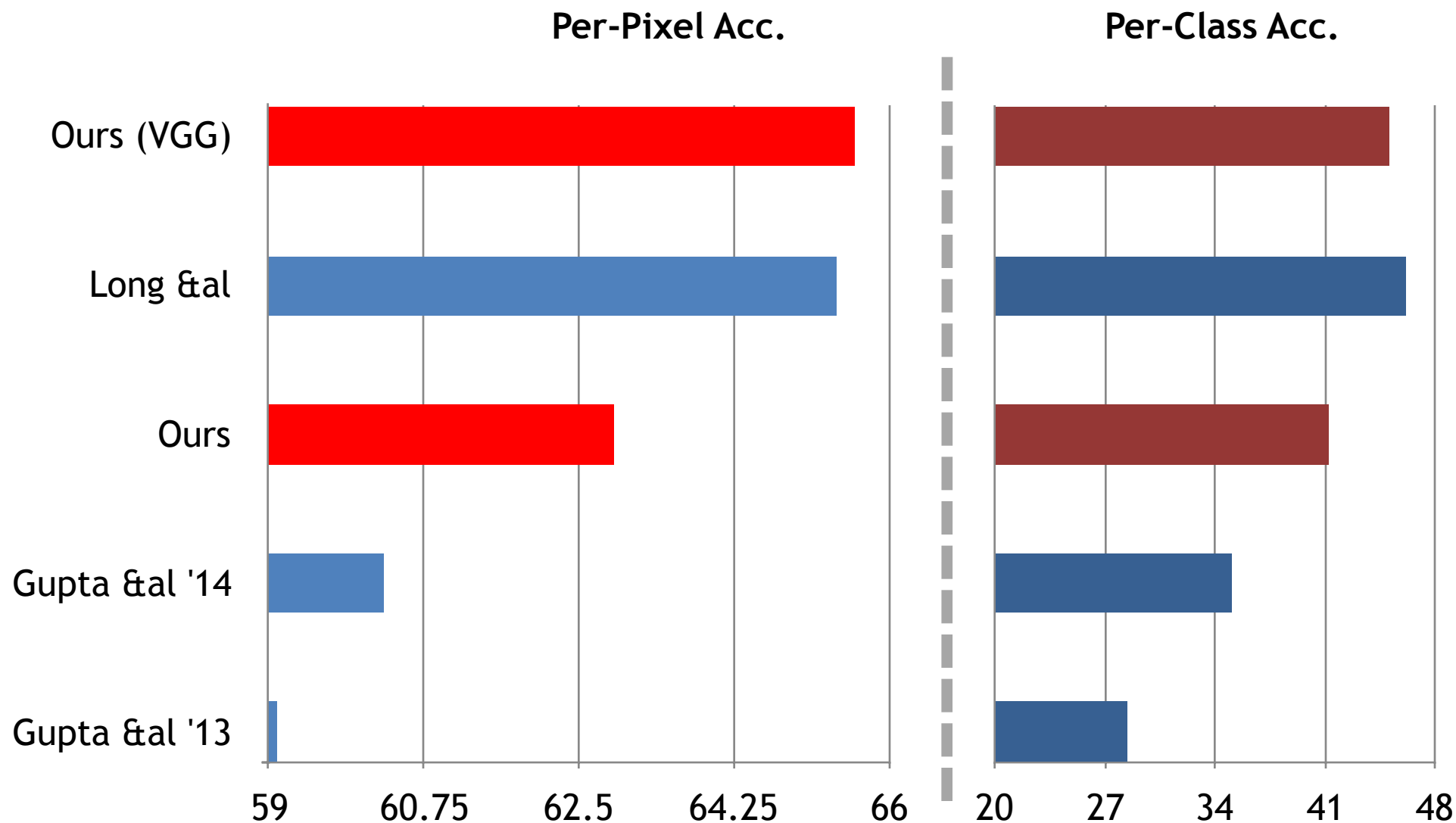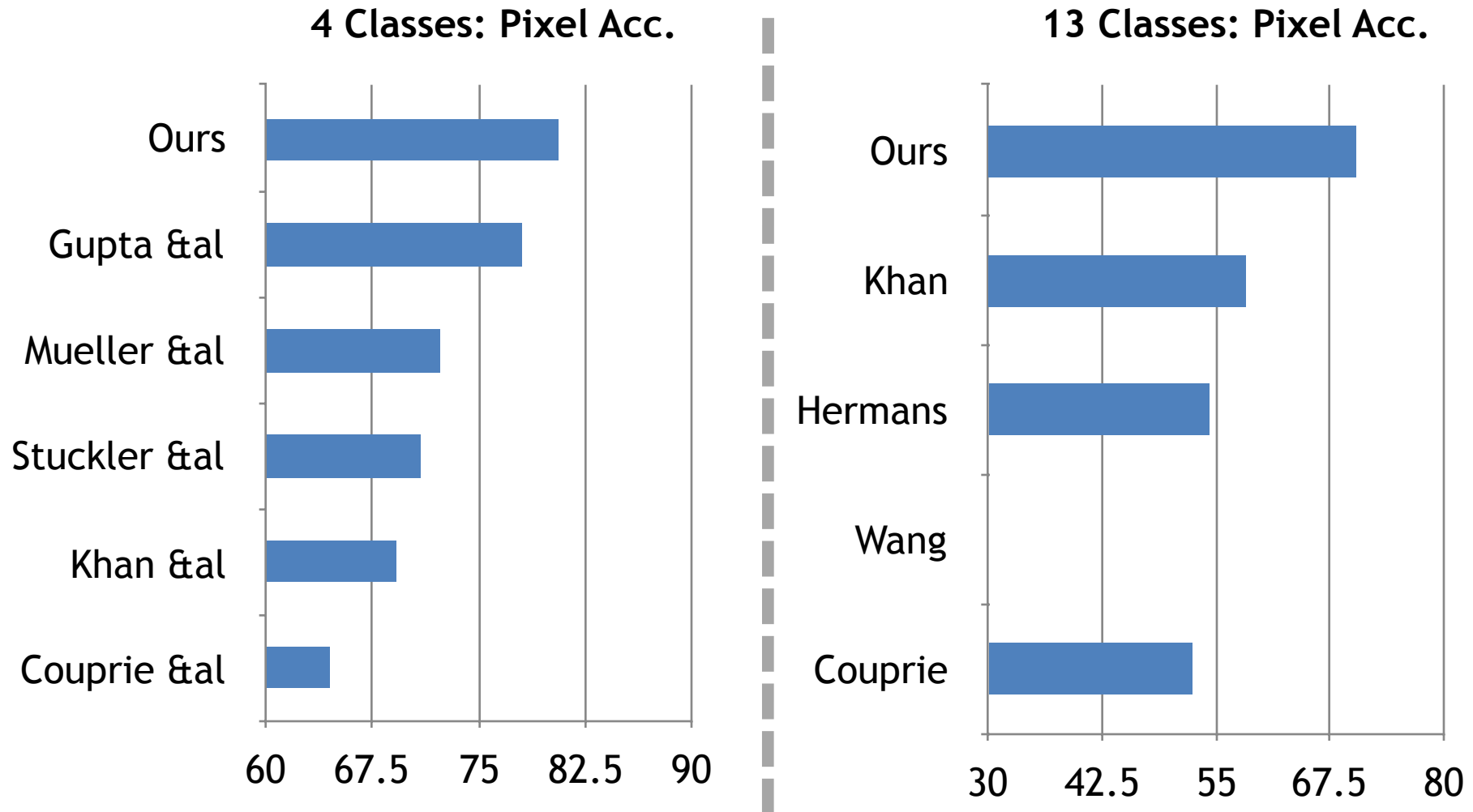- Use RGB + ground truth depth & normals as inputs

# Results: NYUD Labels

- Use RGB + ground truth depth & normals as inputs



**4 Classes: Pixel Acc.**

**13 Classes: Pixel Acc.**

# Semantic Labels: Pascal VOC'11

| Pascal VOC Semantic Segmentation | | | | |
|---|---|---|---|---|
| | Pix. Acc. | Per-Cls Acc. | Freq. Jaccard | Av. Jaccard |
| Long&al [19] | **90.3** | **75.9** | **83.2** | **62.7** |
| Ours (VGG) | **90.3** | 72.4 | 82.9 | 62.2 |

# Contribution from different scales

- On NYU Depth

| Contributions of Scales | | | | | | |
|---|---|---|---|---|---|---|
| | **Depth** | **Normals** | **4-Class** | | **13-Class** | |
| | | | RGB+D+N | RGB | RGB+D+N | RGB |
| | Pixelwise Error lower is better | | Pixelwise Accuracy higher is better | | | |
| Scale 1 only | 0.218 | 29.7 | 71.5 | 71.5 | 58.1 | 58.1 |
| Scale 2 only | 0.290 | 31.8 | 77.4 | 67.2 | 65.1 | 53.1 |
| Scales 1 + 2 | 0.216 | 26.1 | 80.1 | 74.4 | 69.8 | 63.2 |
| Scales 1 + 2 + 3 | 0.198 | 25.9 | 80.6 | 75.3 | 70.5 | 64.0 |

- Depth & normals: scale 1 most important
- Semantic labels: scale 2 most important
  (if D & N are available)

# Using Predicted Depths

- Use <u>predicted</u> depth/normals as input?

**Per-Pixel Acc.**  **Per-Class Acc.**

Scales 1+2:
- 69.8 (RGB + GT D&N)
- 65 (RGB + Pred D&N)
- 63.2 (RGB only)

Scale 2 only:
- 65.1 (RGB + GT D&N)
- 58.7 (RGB + Pred D&N)
- 53.1 (RGB only)

Per-Class Acc. Scales 1+2:
- 58.9 (RGB + GT D&N)
- 49.5 (RGB + Pred D&N)
- 50.6 (RGB only)

Per-Class Acc. Scale 2 only:
- 52.3 (RGB + GT D&N)
- 43.8 (RGB + Pred D&N)
- 38.3 (RGB only)

Per-Pixel axis: 35, 43.75, 52.5, 61.25, 70
Per-Class axis: 30, 37.5, 45, 52.5, 60

- NYU Depth 13-class

RGB only  RGB + Pred D&N  RGB + GT D&N

# Summary

- Relatively simple multi-scale model gives good results for depth, normals & labels

- Coarse interpretation of scene important for understanding depth/normals

- See ICCV 2015 paper: "Predicting Depth, Surface Normals and Semantic Labels with a Common
Multi-Scale Convolutional Architecture", D. Eigen and R. Fergus, arXiv 1411.4734

- Code available

# Overview

- Methods where output is also an image
  - Fully Convolutional Nets [Long et al., CVPR 2015]
  - Depth, normals and semantic labels from a single image [Eigen ICCV 2015]

- Image processing with Convnets
  - Image colorization [Zhang et al. ECCV 2016]

# Denoising with ConvNets

- Burger et al. "Can plain NNs compete with BM3D?" CVPR 2012

Original

Noised

Denoised

# Deblurring with Convnets

- Blind deconvolution
  - Learning to Deblur, Schuler et al., arXiv 1406.7444, 2014



Blurry image with ground truth kernel

Result of [Zho+13] PSNR 23.17

Deblurring result w. noise *agnostic* training PSNR 23.29

Deblurring result w. noise *specific* training **PSNR 23.41**

# Inpainting with Convnets

- Image Denoising and Inpainting with Deep Neural Networks, Xie et al. NIPS 2012.

- Mask-specific inpainting with deep neural networks, Köhler et al., Pattern Recognition 2014



Original        Schmid CVPR'10        Köhler et al. '14

# Removing Local Corruption

- Restoring An Image Taken Through a Window Covered with Dirt or Rain, Eigen et al., ICCV 2013.

# Removing Local Corruption

Restoring An Image Taken
Through a Window Covered with
Dirt or Rain

Rain Sequence
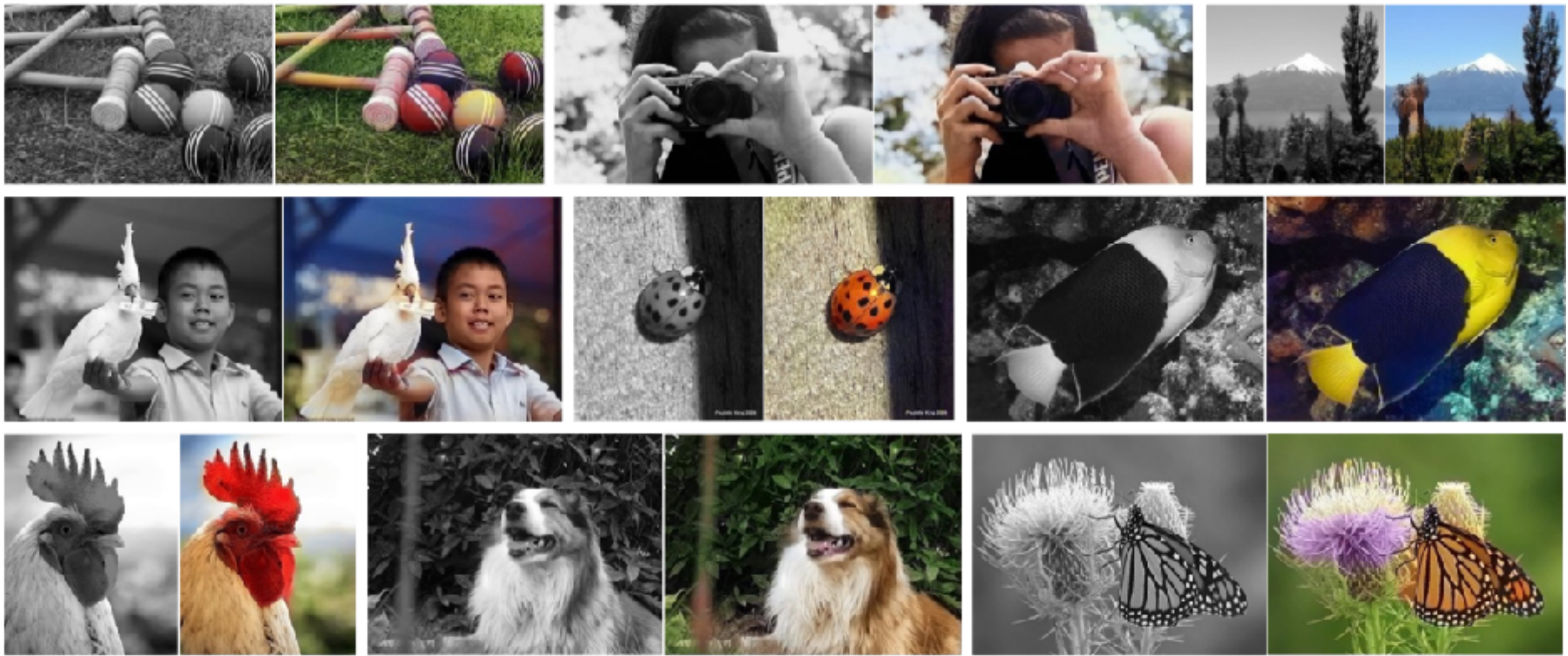
Each frame processed independently

David Eigen, Dilip Krishnan and Rob Fergus

ICCV 2013

# Overview

- Methods where output is now also an image
  - Fully Convolutional Nets [Long et al., CVPR 2015]
  - Depth, normals and semantic labels from a single image [Eigen ICCV 2015]

- Image processing with Convnets
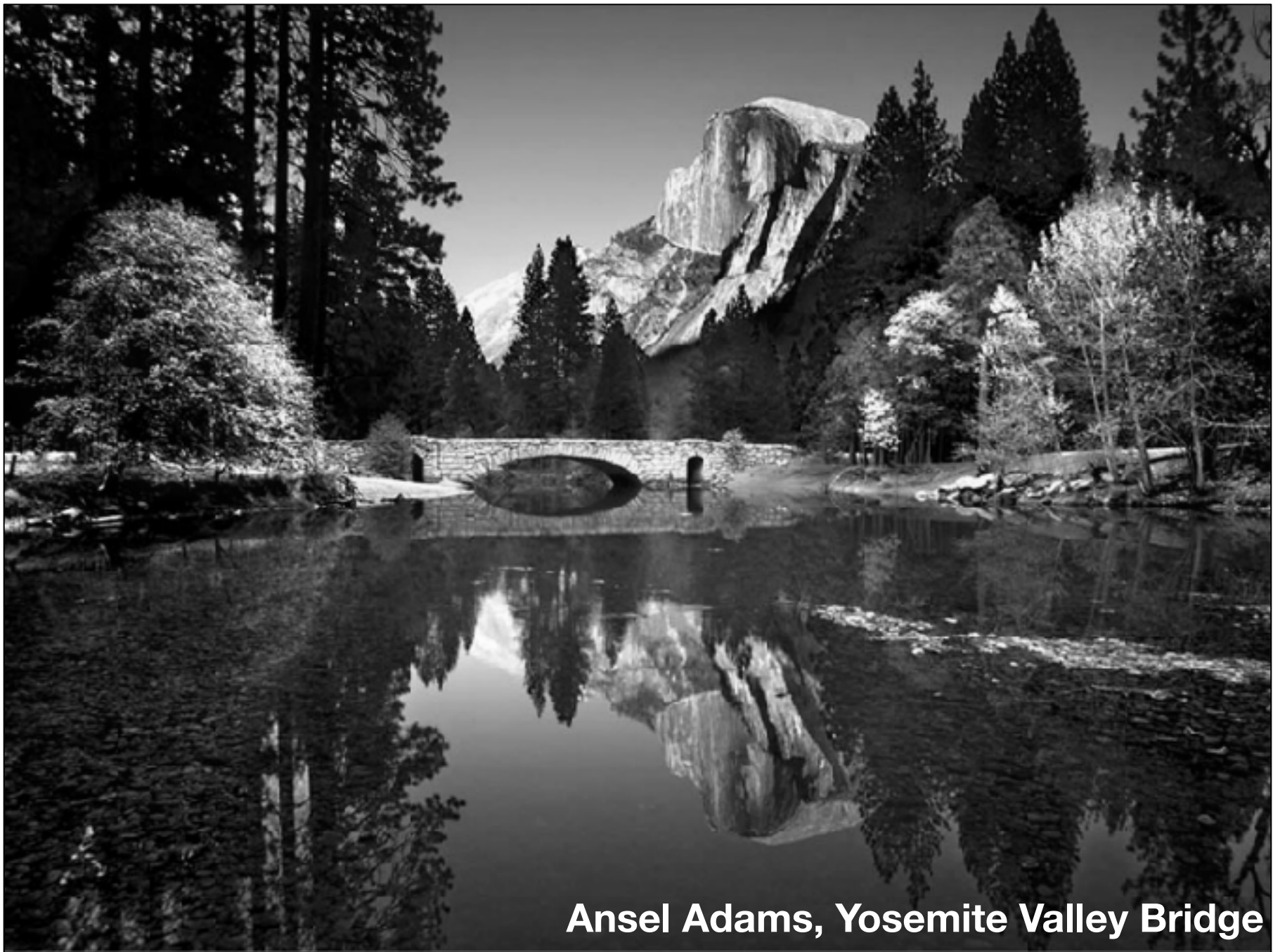  - Image colorization [Zhang et al. ECCV 2016]

# Colorful Image Colorization

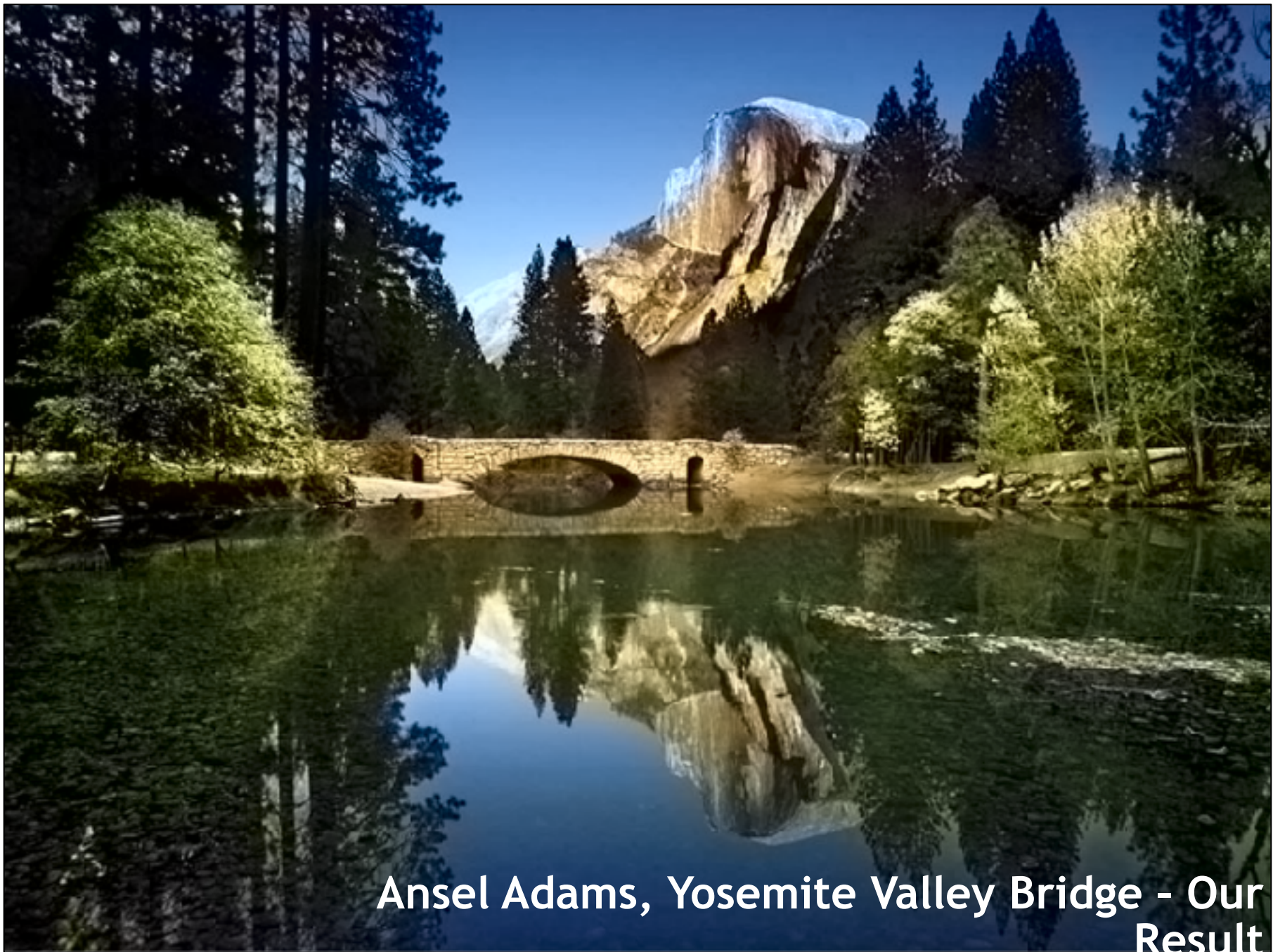Richard Zhang, Phillip Isola, Alexei (Alyosha) Efros

richzhang.github.io/colorization
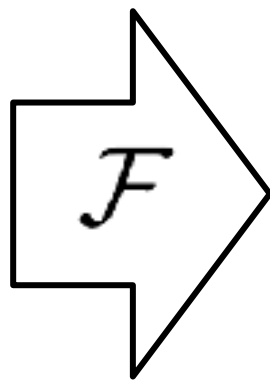
**Ansel Adams, Yosemite Valley Bridge**

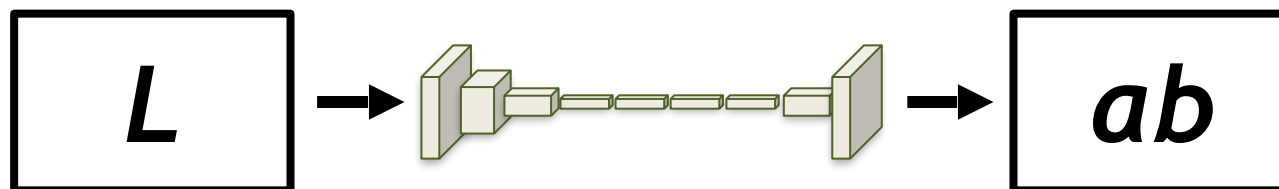Ansel Adams, Yosemite Valley Bridge – Our Result

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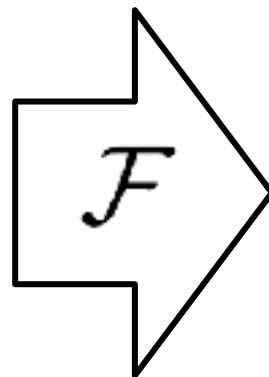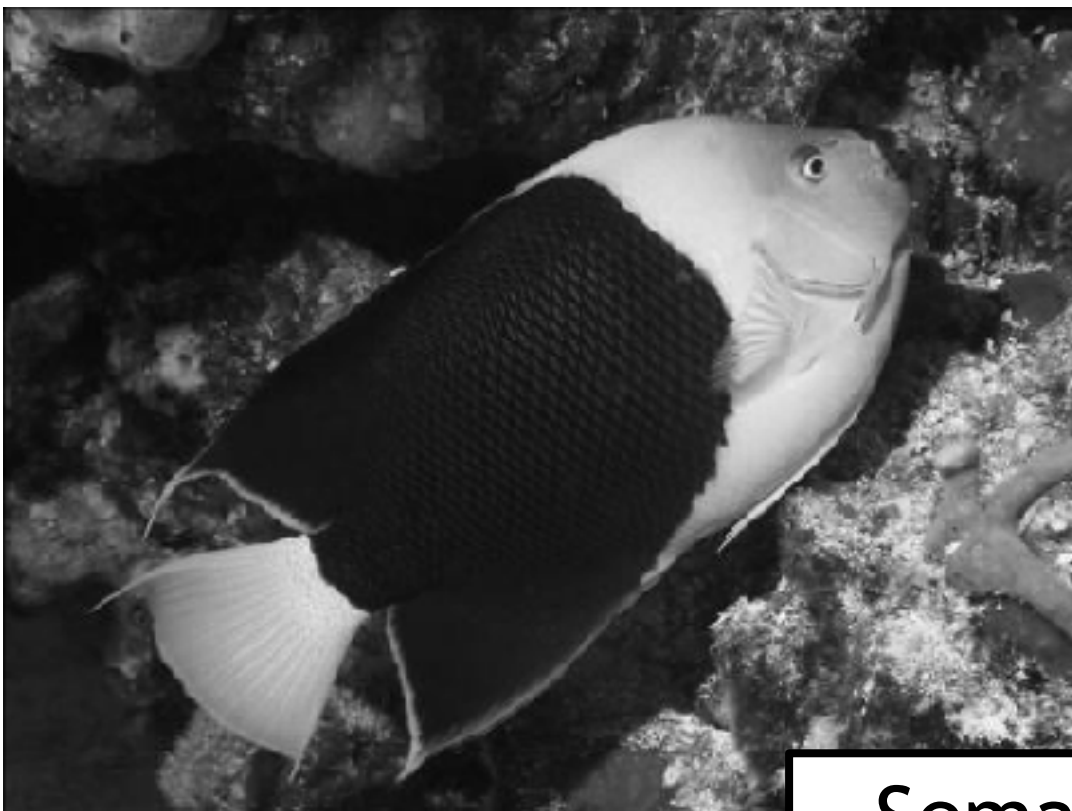$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

Color information: *ab* channel
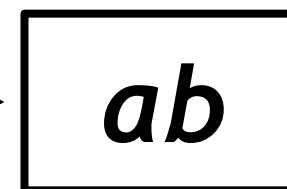$$\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$$

Grayscale image: $L$

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

Semantics? Higher-level abstraction?

ncatenate (L,ab)

$(\mathbf{X}, \widehat{\mathbf{Y}})$

$L$

$ab$

"Free" supervisory signal

# Inherent Ambiguity



Grayscale

# Inherent Ambiguity



Our Output

Ground Truth

# Better Loss Function

- Regression with L2 loss inac

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

# Better Loss Function

**Colors in *ab* space**
(discrete)

- Regression with L2 loss inadequate

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

- Use **multinomial classification**

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

- **Class rebalancing** to encourage

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

# Better Loss Function

log$_{10}$ probability

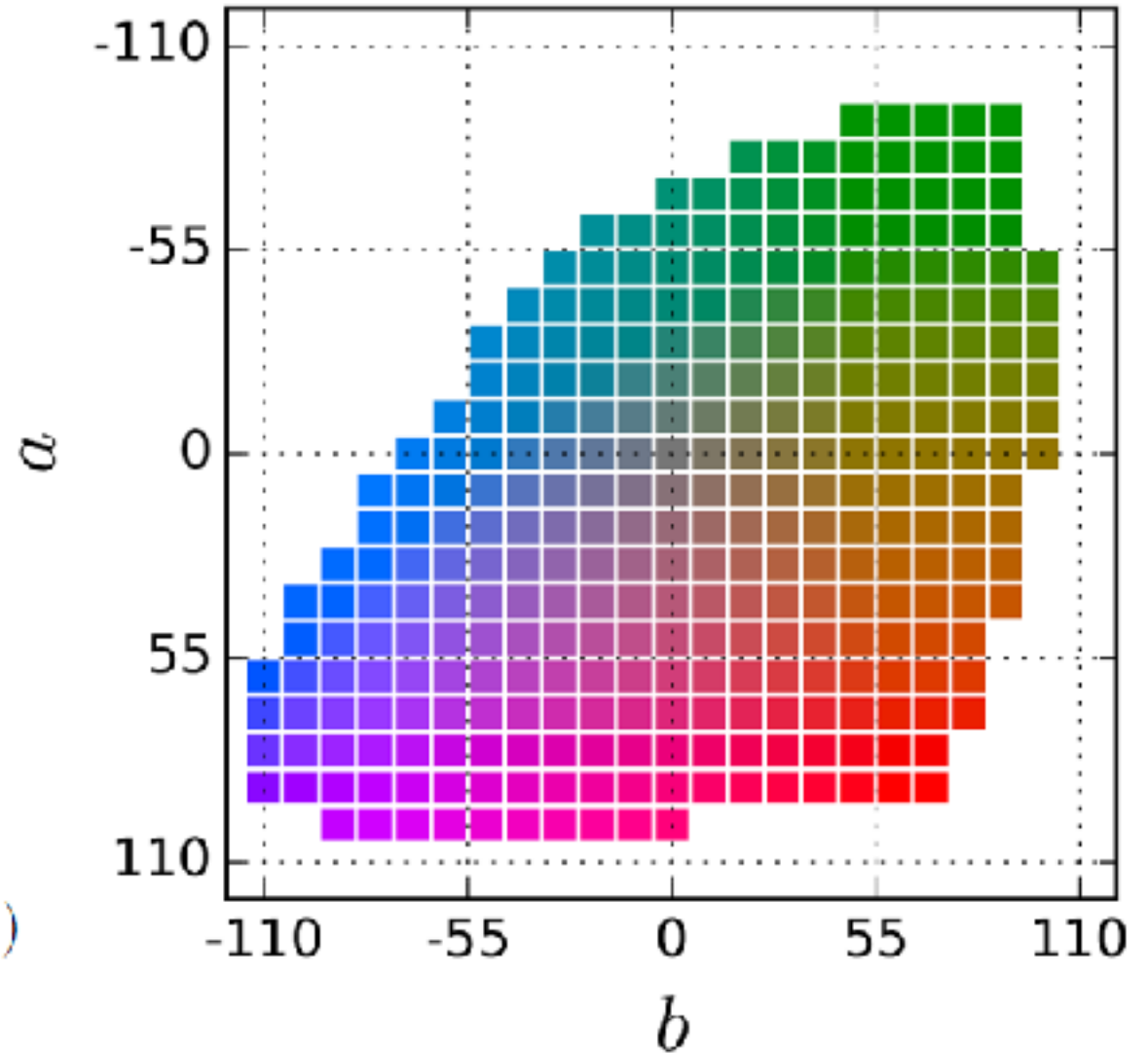**Histogram over *ab* space**

- Regression with L2 loss inadequate

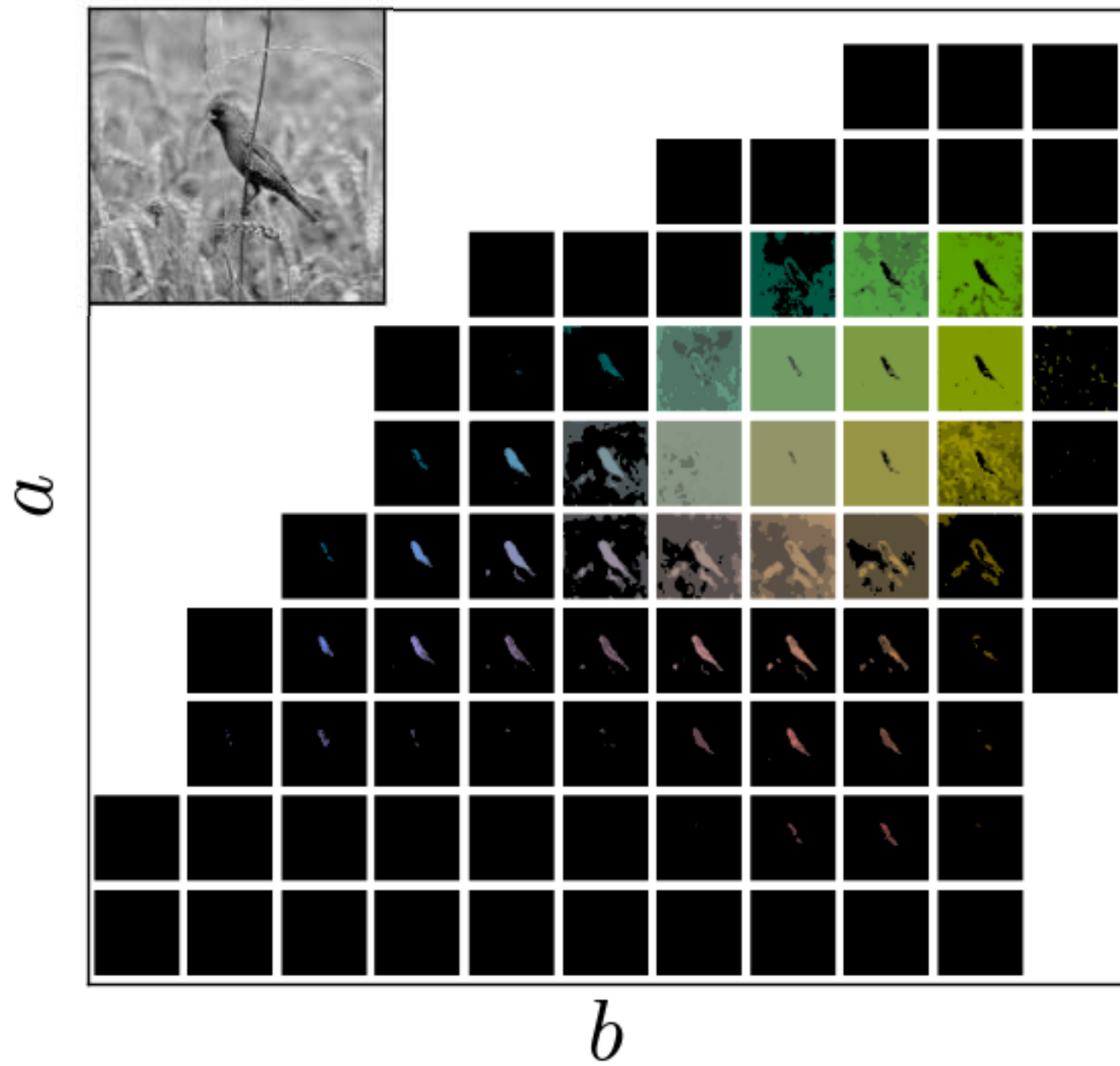$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$
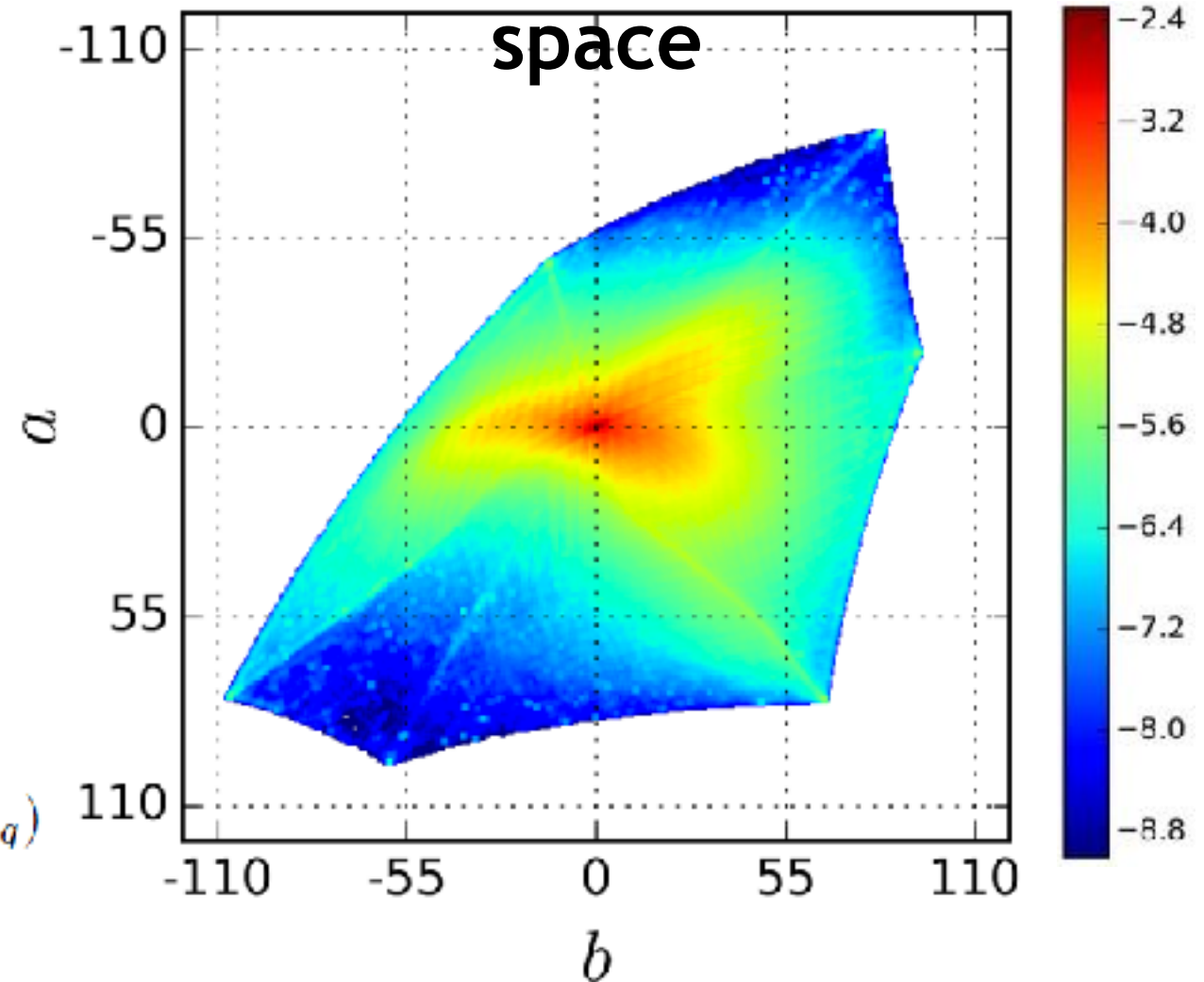
- Use **multinomial classification**

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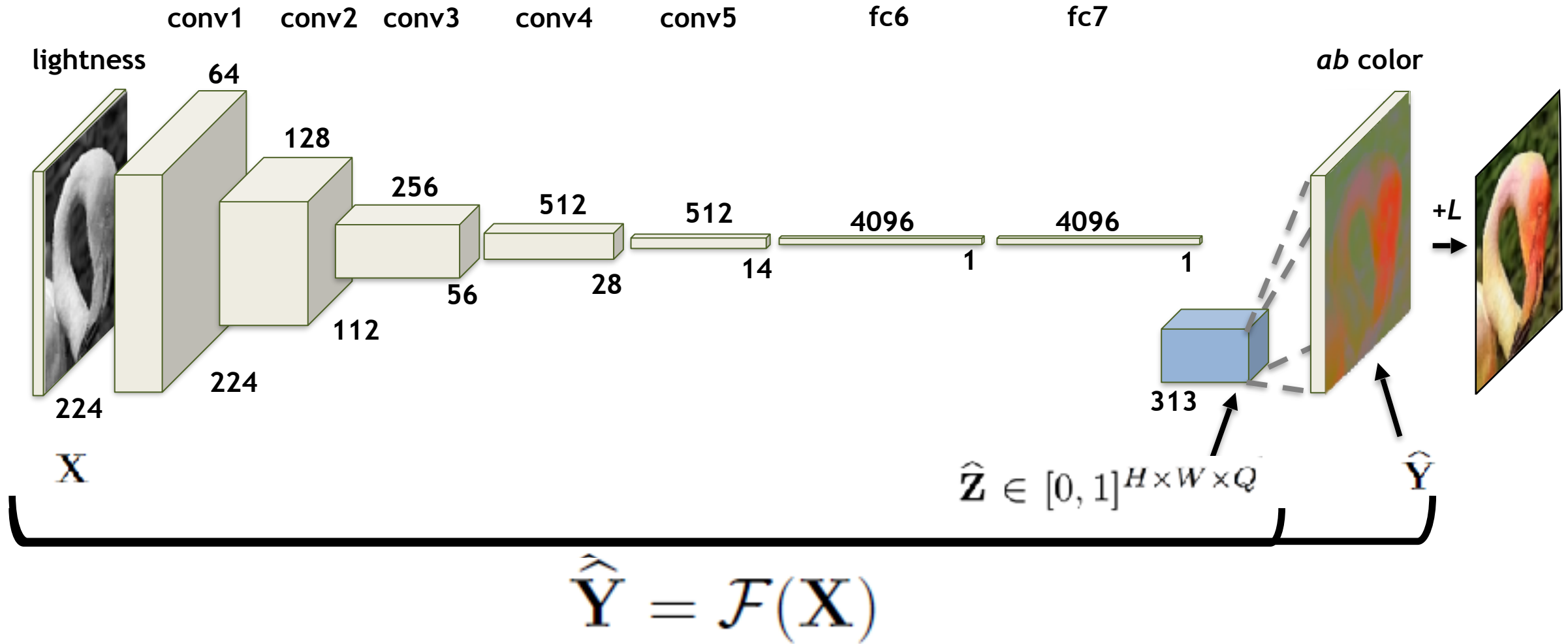
- **Class rebalancing** to encourage

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

# Network Architecture



lightness

conv1    conv2   conv3     conv4      conv5       fc6        fc7

64

128

256

512

512

4096

4096

313

*ab* color

+*L*

224

224

112

56

28

14

1

1

$\mathbf{X}$

$\widehat{\mathbf{Z}} \in [0,1]^{H \times W \times Q}$

$\widehat{\mathbf{Y}}$

$$\widehat{\widehat{\mathbf{Y}}} = \mathcal{F}(\mathbf{X})$$

# Network Architecture



conv1  conv2  conv3  conv4  conv5  conv6  conv7  conv8

à trous [1]/dilated [2] à trous/dilated

lightness  64  128  256  512  512  512  512  256  ab color

224  224  112  56  28  28  28  28  56

313

$\hat{\mathbf{Z}} \in [0, 1]^{H \times W \times Q}$

$\mathbf{X}$  $\hat{\mathbf{Y}}$

$+L$

$\hat{\mathbf{Z}} = \mathcal{G}(\mathbf{X})$

$\hat{\mathbf{Y}} = \mathcal{H}(\hat{\mathbf{Z}})$

[1] Chen *et al*. In arXiv, 2016.
[2] Yu and Koltun. In ICLR,
2016.

Ground Truth | Input | L2 Regression | Class w/ Rebalancing

# Failure Cases

# Biases

# Evaluation

## Visual Quality

## Quantitative

# Evaluation

| | **Visual Quality** | **Representation Learning** |
|---|---|---|
| **Quantitative** | Per-pixel accuracy<br><br>Perceptual realism<br><br>Semantic interpretability | Task generalization<br>ImageNet classification<br><br>Task & dataset generalization<br>PASCAL classification, detection, segmentation |
| **Qualitative** | Low-level stimuli<br><br>Legacy grayscale photos | Hidden unit activations |

# Evaluation

| | **Visual Quality** | **Representation Learning** |
|---|---|---|
| **Quantitative** | Per-pixel accuracy | Task generalization |
| | | ImageNet classification |
| | **Perceptual realism** | |
| | | Task & dataset generalization |
| | Semantic interpretability | PASCAL classification, detection, segmentation |
| **Qualitative** | Low-level stimuli | |
| | | Hidden unit activations |
| | Legacy grayscale photos | |

# Perceptual Realism / Amazon Mechanical Turk Test

clap if "fake" clap if "fake"

Fake, 0% fooled

clap if "fake" clap if "fake"

Fake, 55% fooled

clap if "fake"  clap if "fake"

Fake, 58% fooled

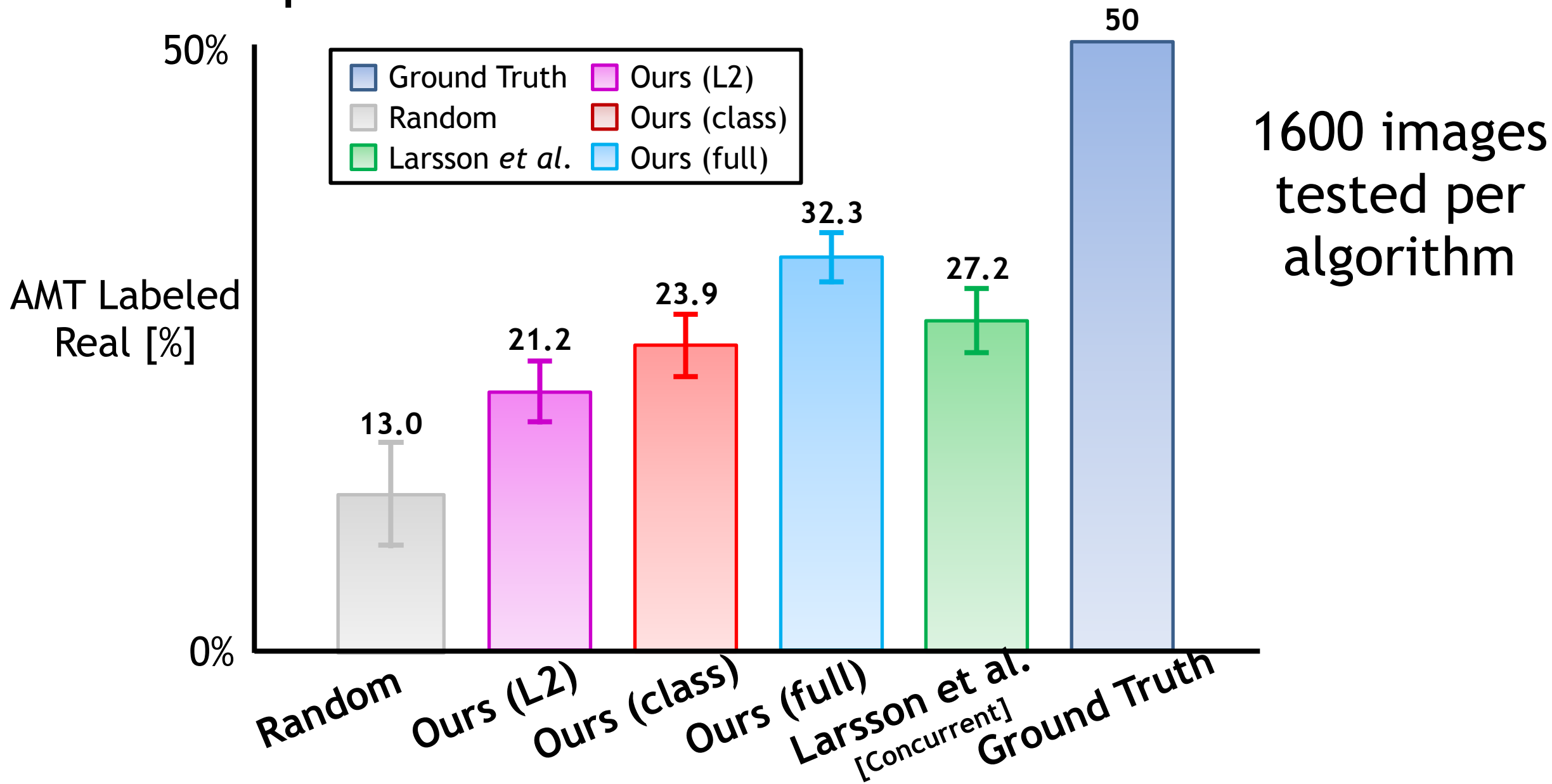from Reddit /u/SherySantucci

Recolorized by Reddit ColorizeBot

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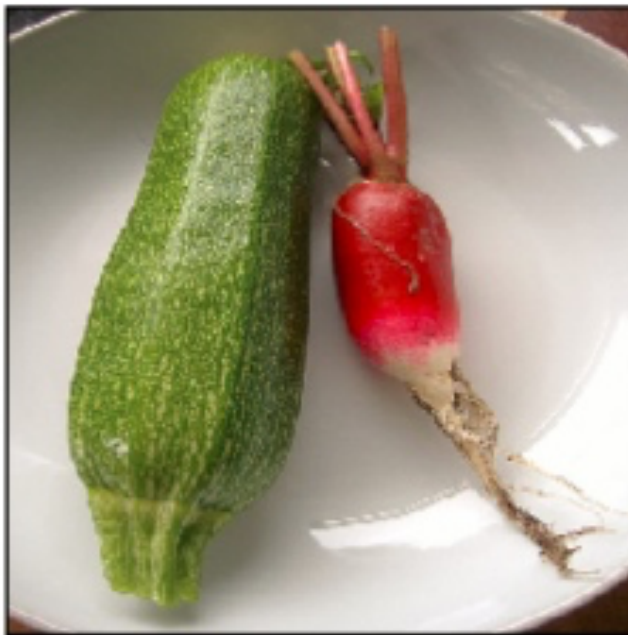
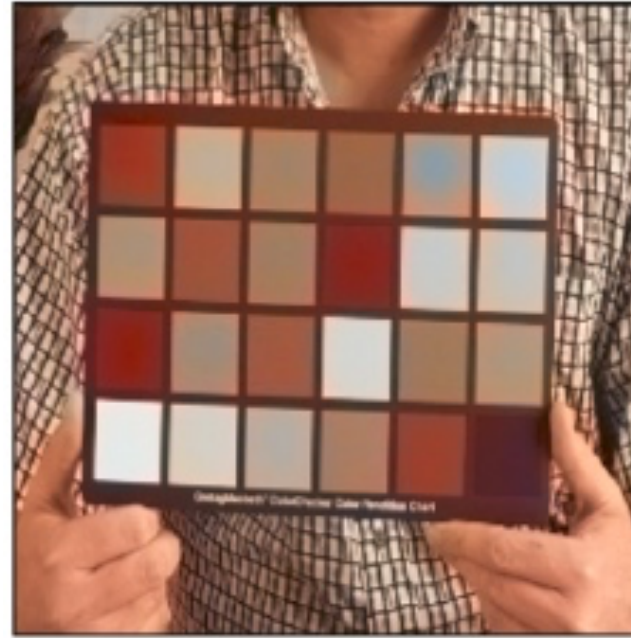Recolorized by Reddit ColorizeBot

# Perceptual Realism Test



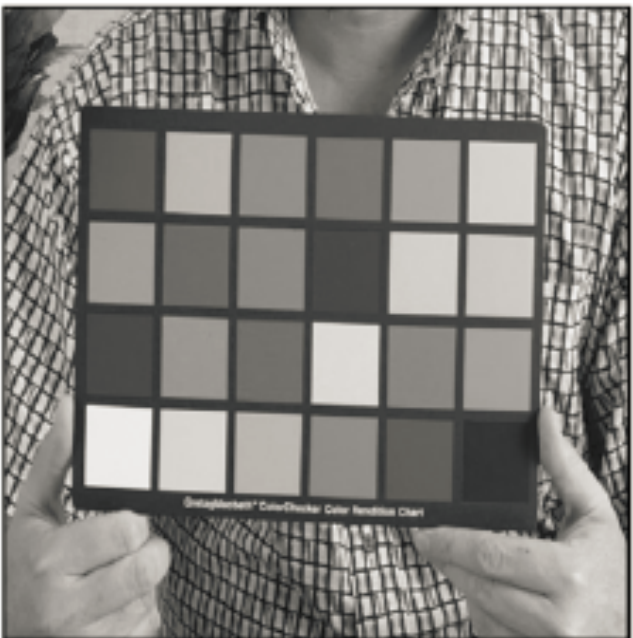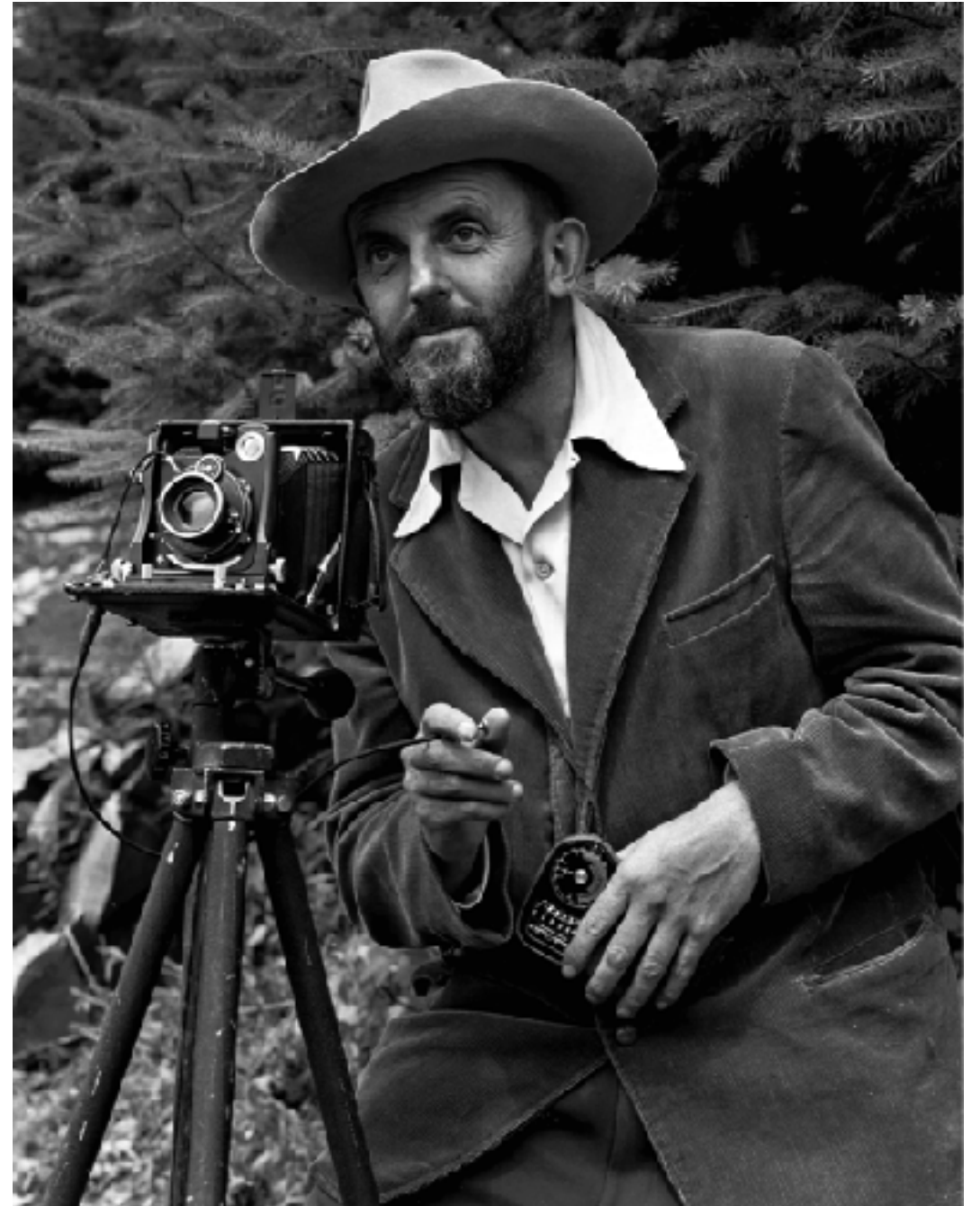1600 images tested per algorithm

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

Thylacine, Dr. David Fleay, extinct in 1936.

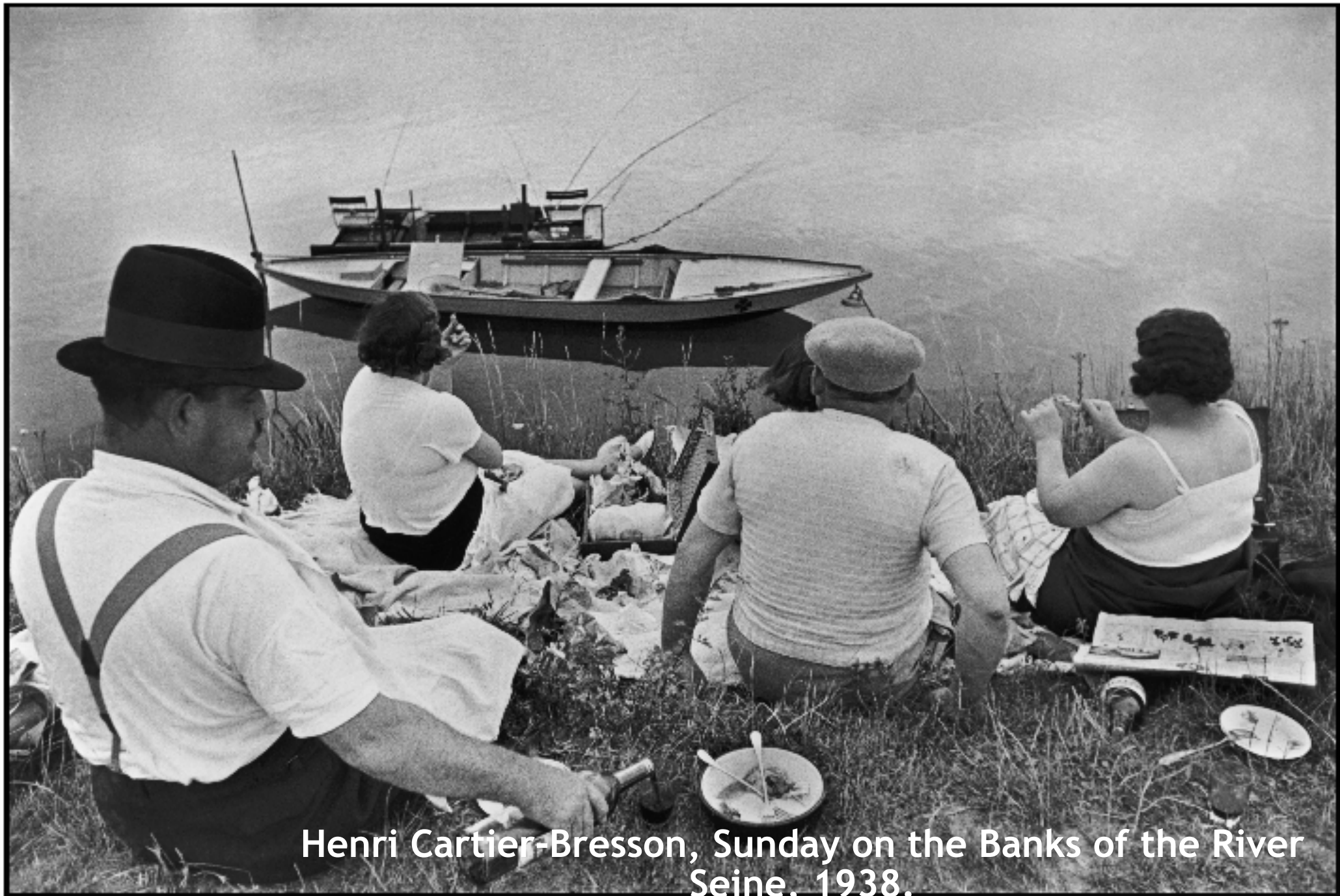Thylacine, Dr. David Fleay, extinct in 1936.

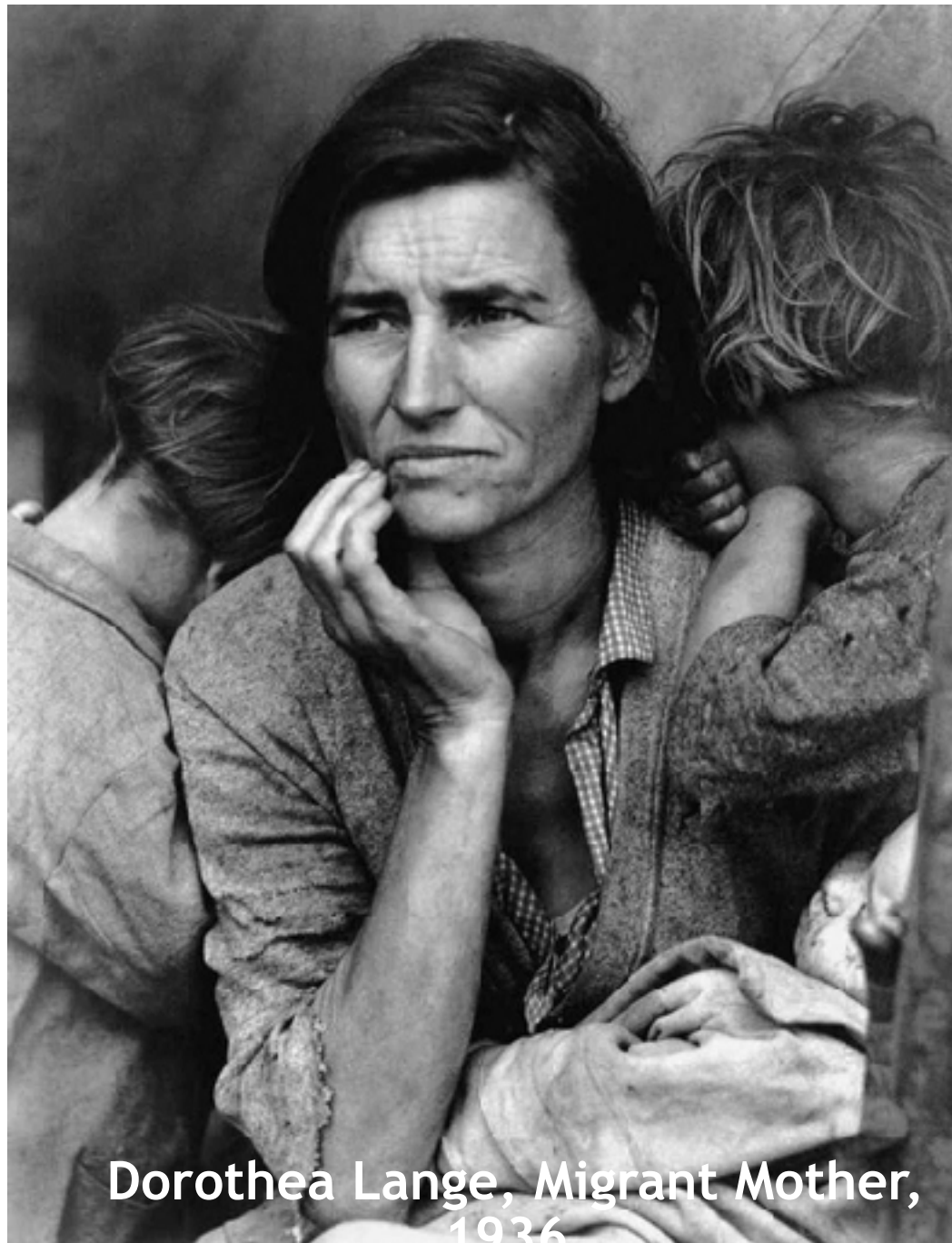Amateur Family Photo, 1956.

Amateur Family Photo, 1956.

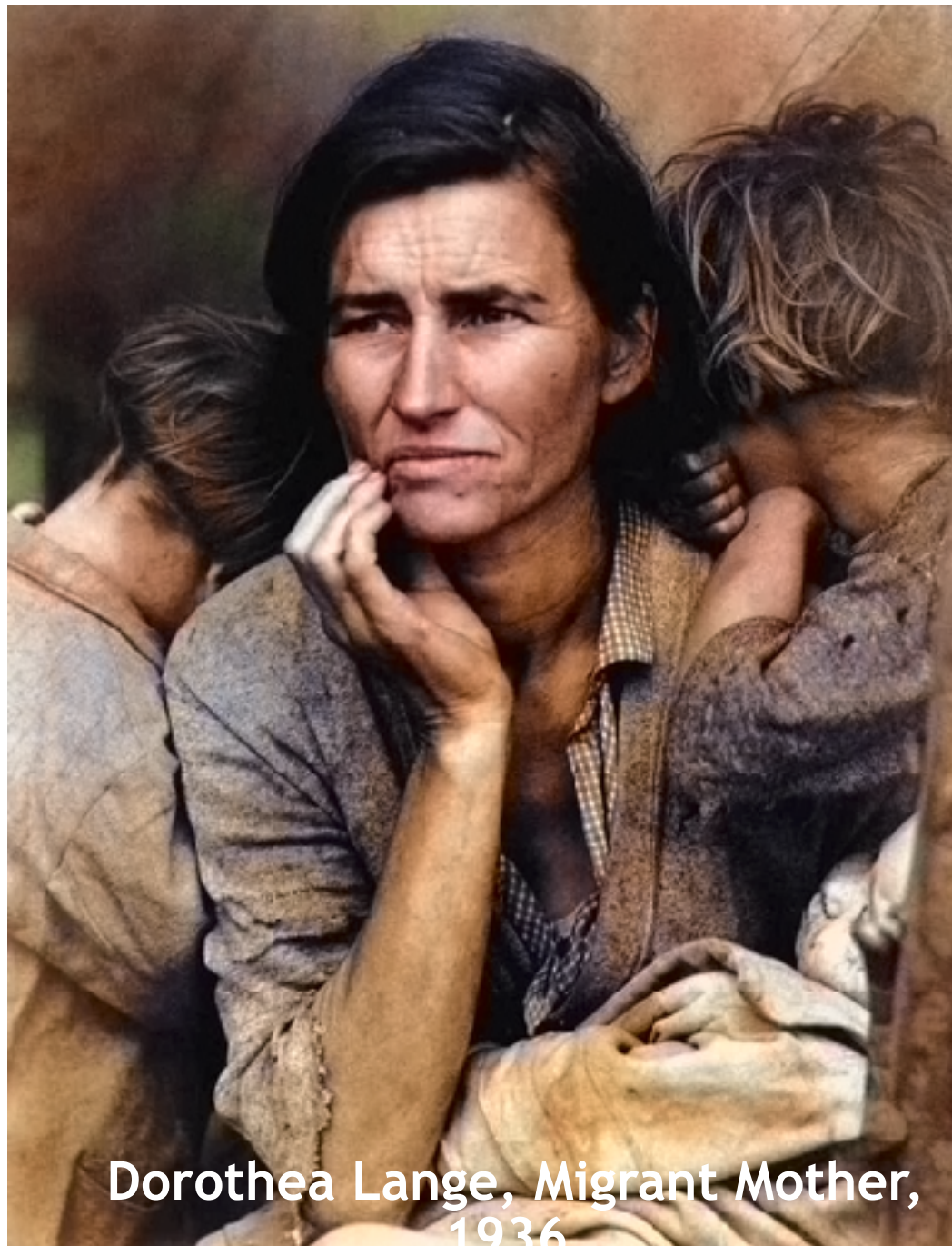Henri Cartier-Bresson, Sunday on the Banks of the River Seine, 1938.

Henri Cartier-Bresson, Sunday on the Banks of the River Seine, 1938.
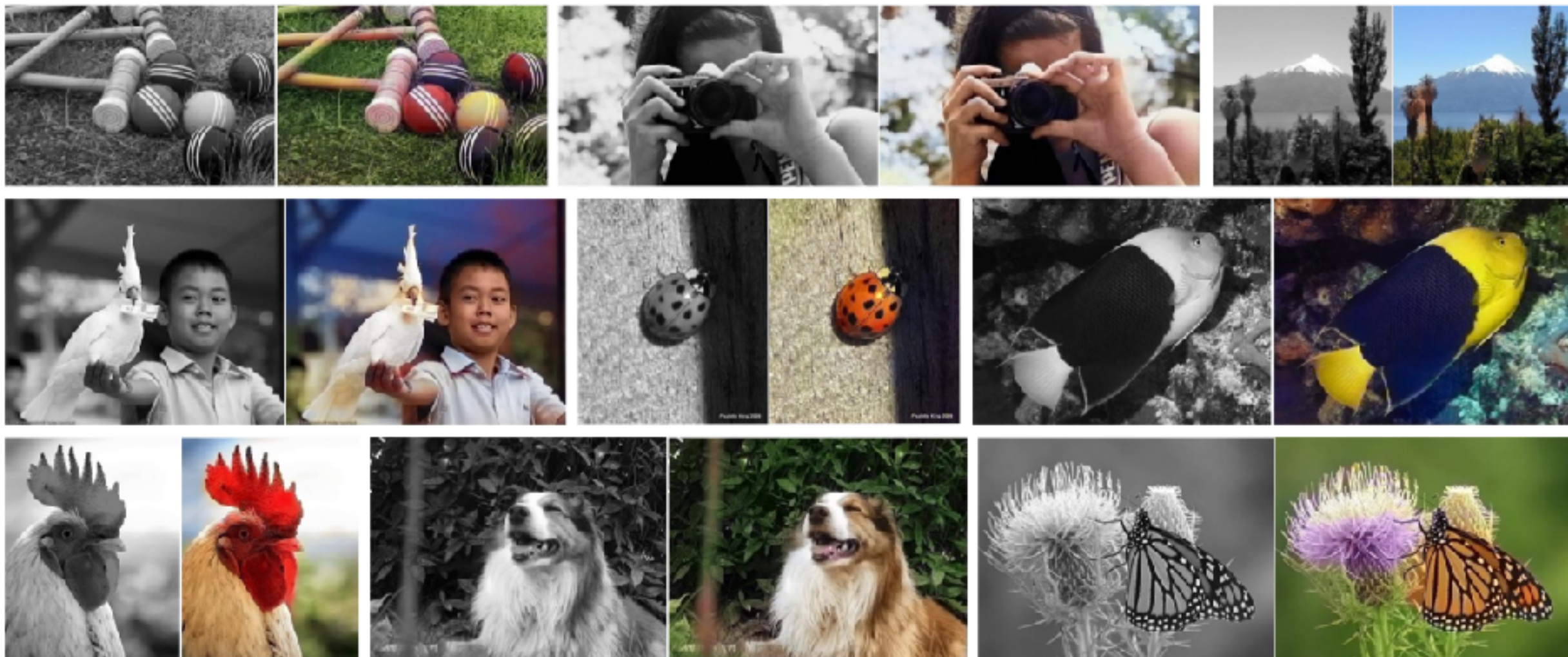
Dorothea Lange, Migrant Mother, 1936

Dorothea Lange, Migrant Mother, 1936

# Additional Information

- Demo
  - http://demos.algorithmia.com/colorize-photos/
- Reddit ColorizeBot
  - Type "colorizebot" under any image post
- Code
  - https://github.com/richzhang/colorization
- Website – full paper, user examples, visualizations
  - http://richzhang.github.io/colorization

For the full paper, additional examples and our model:
richzhang.github.io/colorization