

# Colorful Image Colorization

Richard Zhang, Phillip Isola, Alexei (Alyosha) Efros

[richzhang.github.io/colorization](https://richzhang.github.io/colorization)



Ansel Adams, Yosemite Valley Bridge



Ansel Adams, Yosemite Valley Bridge – Our Result



Grayscale image:  $L$  channel

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





Grayscale image:  $L$  channel

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

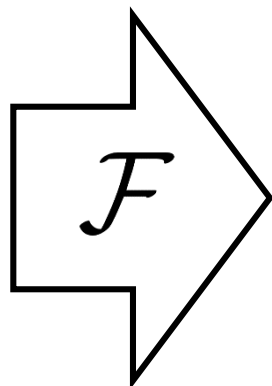
$L$



Color information:  $ab$  channels

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

$ab$

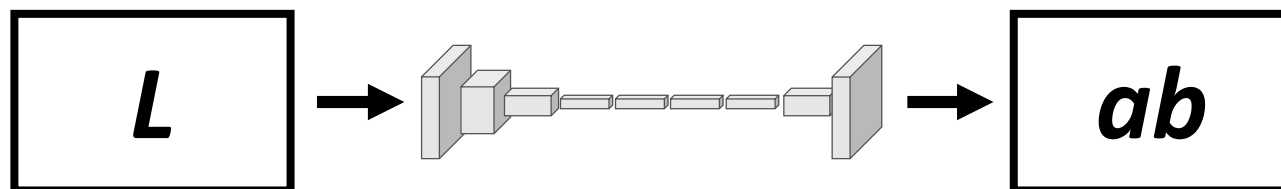


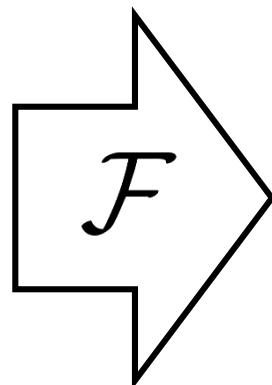
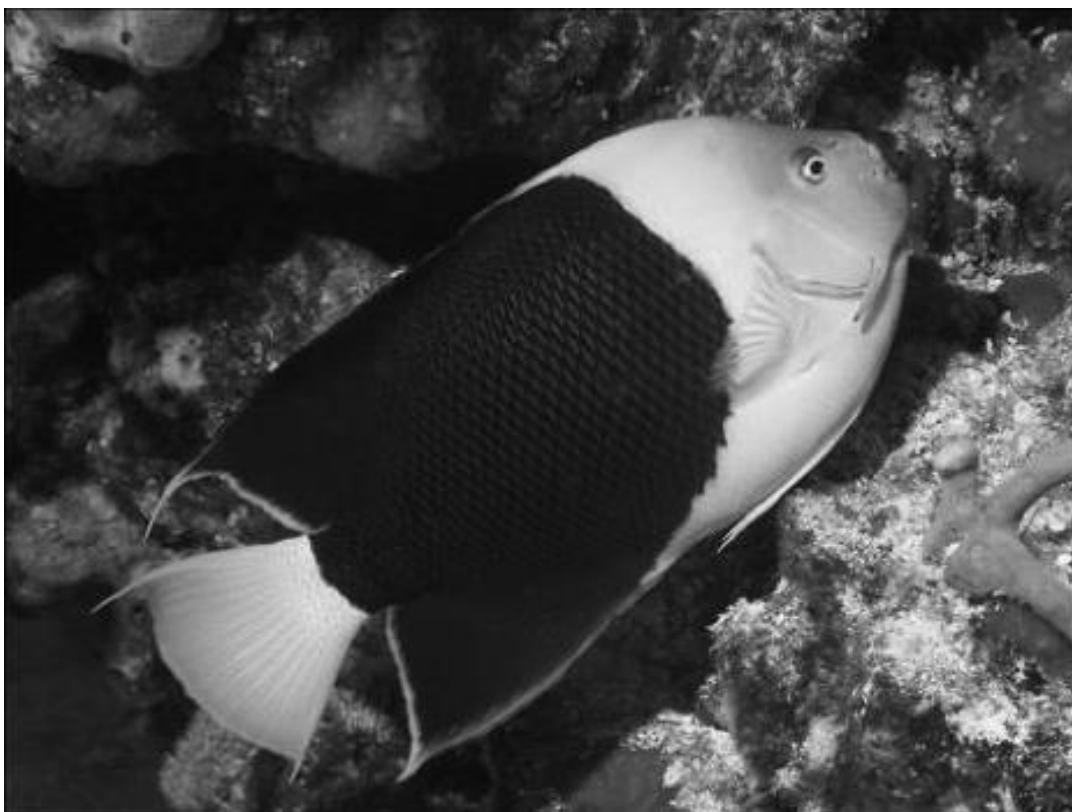
Grayscale image:  $L$  channel

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

Color information:  $ab$  channels

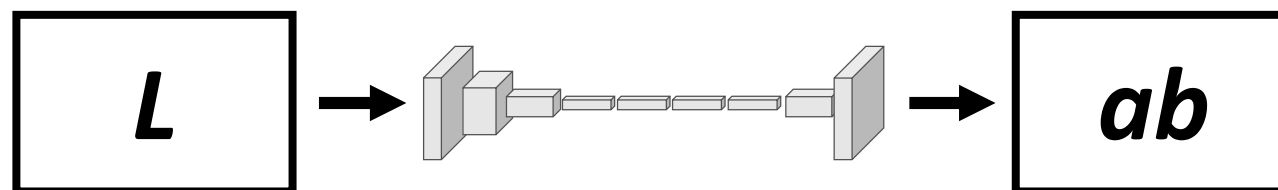
$$\hat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$$

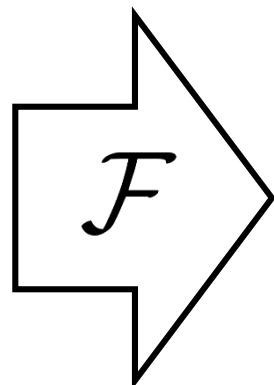




Grayscale image:  $L$  channel

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



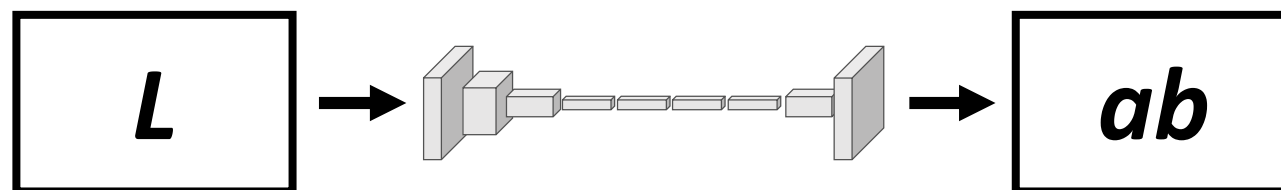


Grayscale image:  $L$  channel

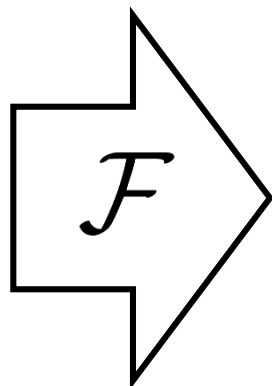
$$\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$$

Concatenate (L,ab)

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





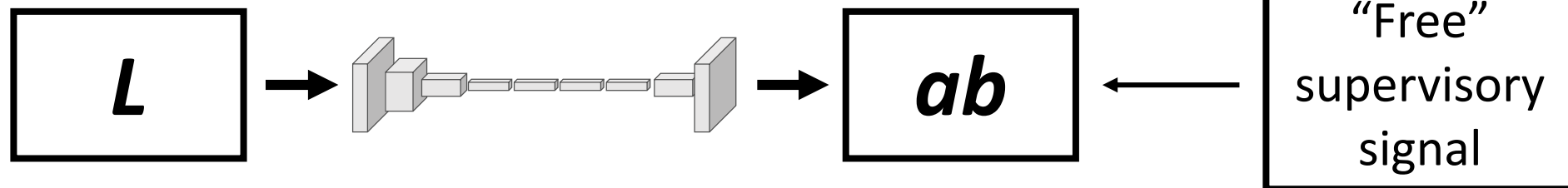


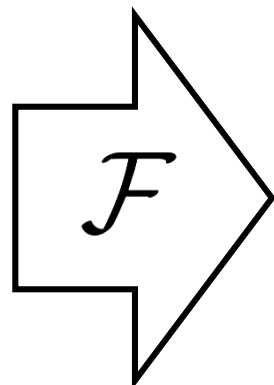
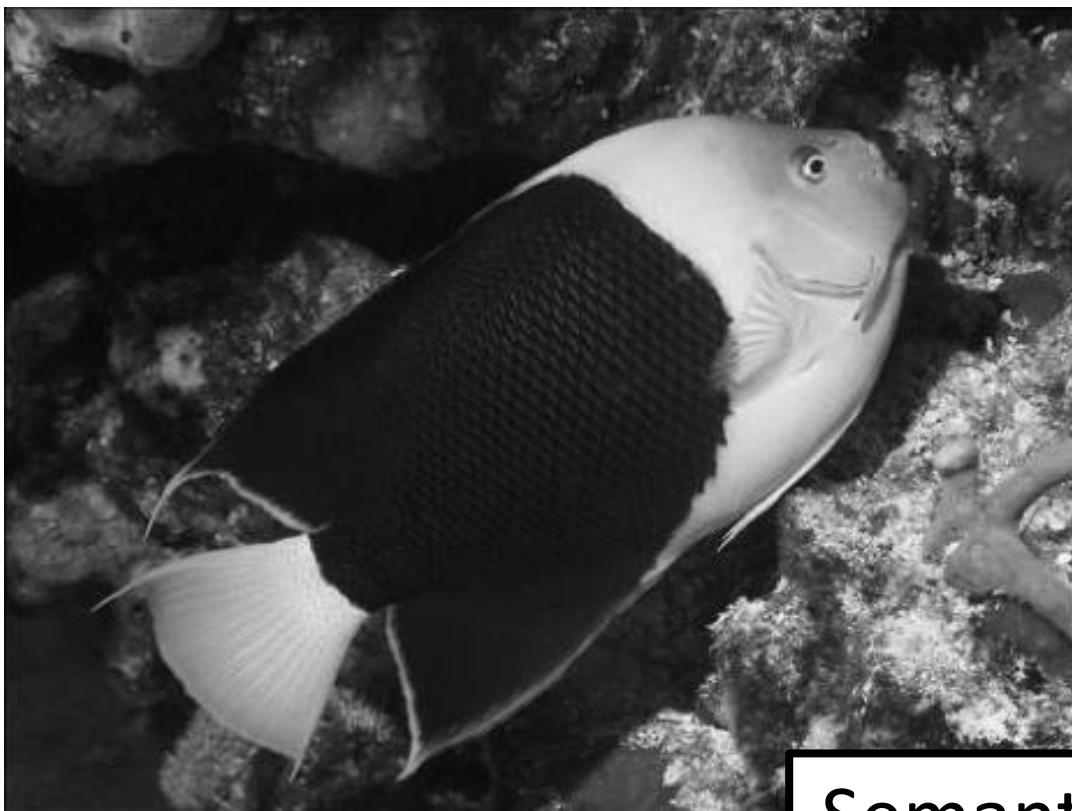
Grayscale image:  $L$  channel

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

Concatenate ( $L, ab$ )

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

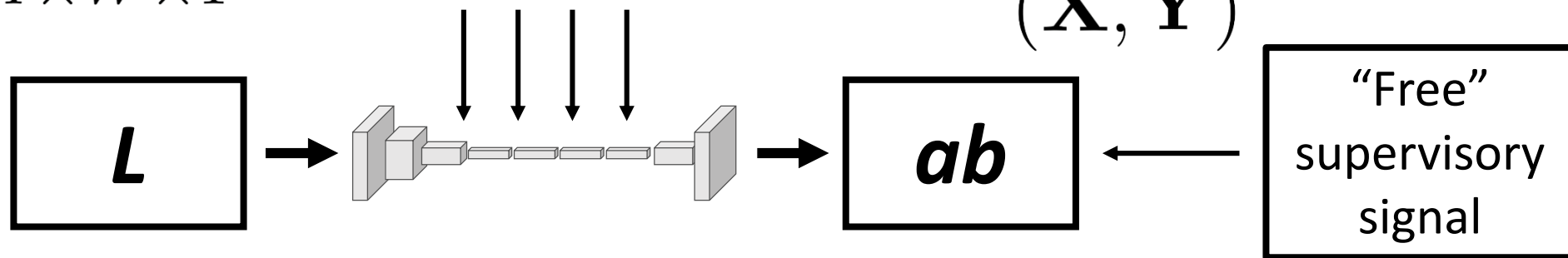




Semantics? Higher-level abstraction?

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

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



# Inherent Ambiguity



Grayscale

# Inherent Ambiguity



Our Output

# Inherent Ambiguity



Our Output



Ground Truth

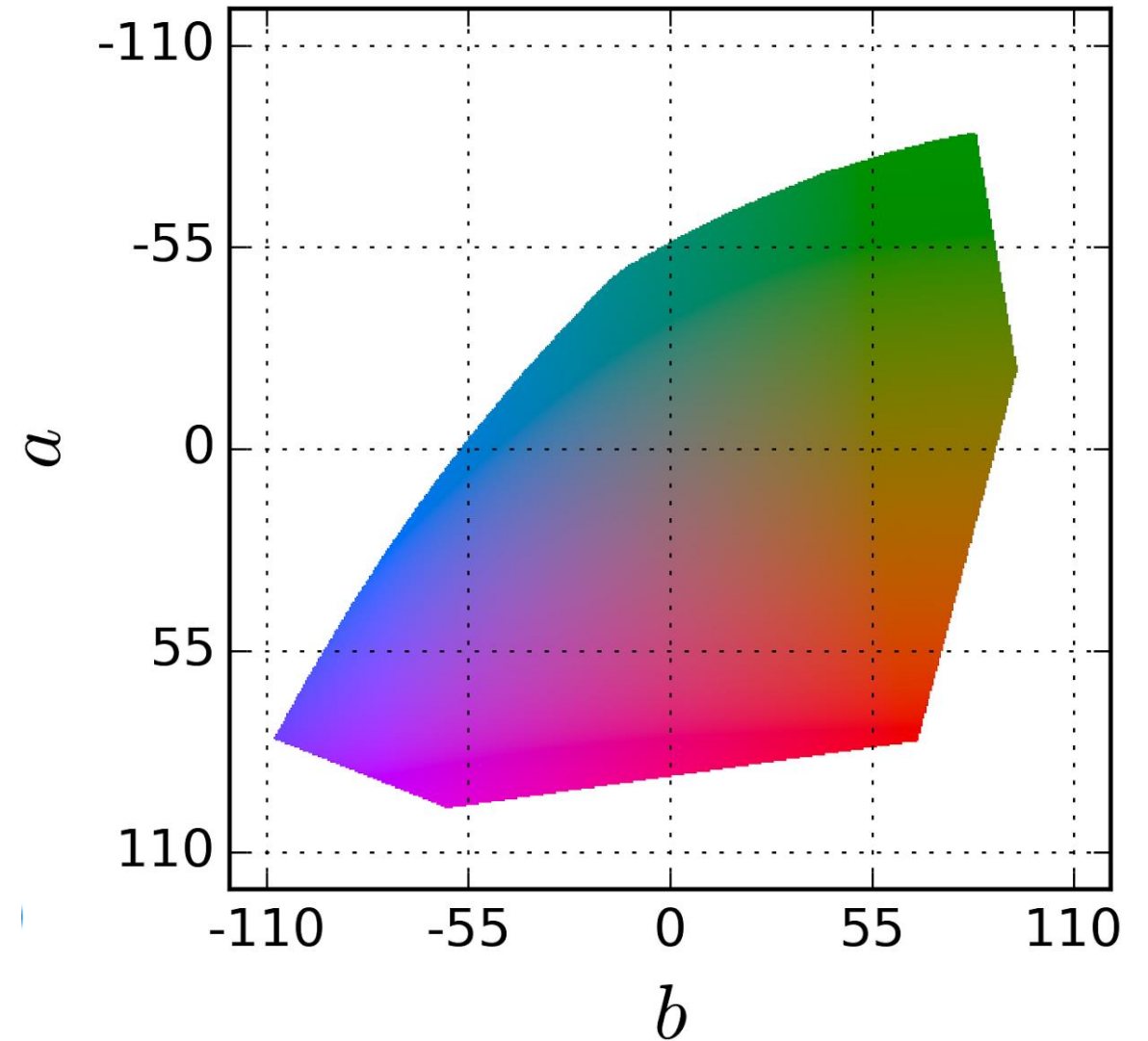
# Better Loss Function

- Regression with L2 loss inadequate

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

## Colors in *ab* space

(continuous)



# Better Loss Function

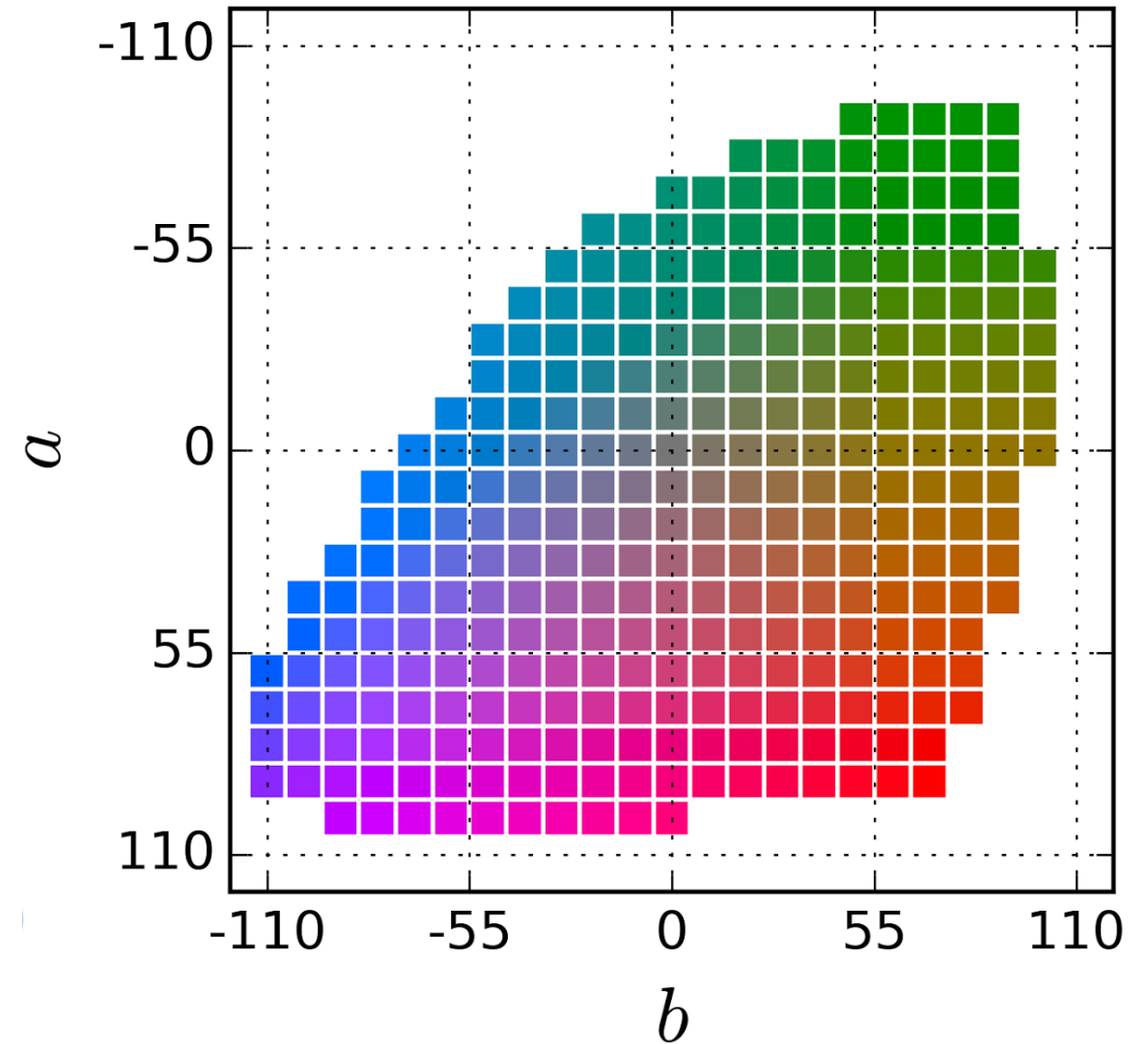
- Regression with L2 loss inadequate

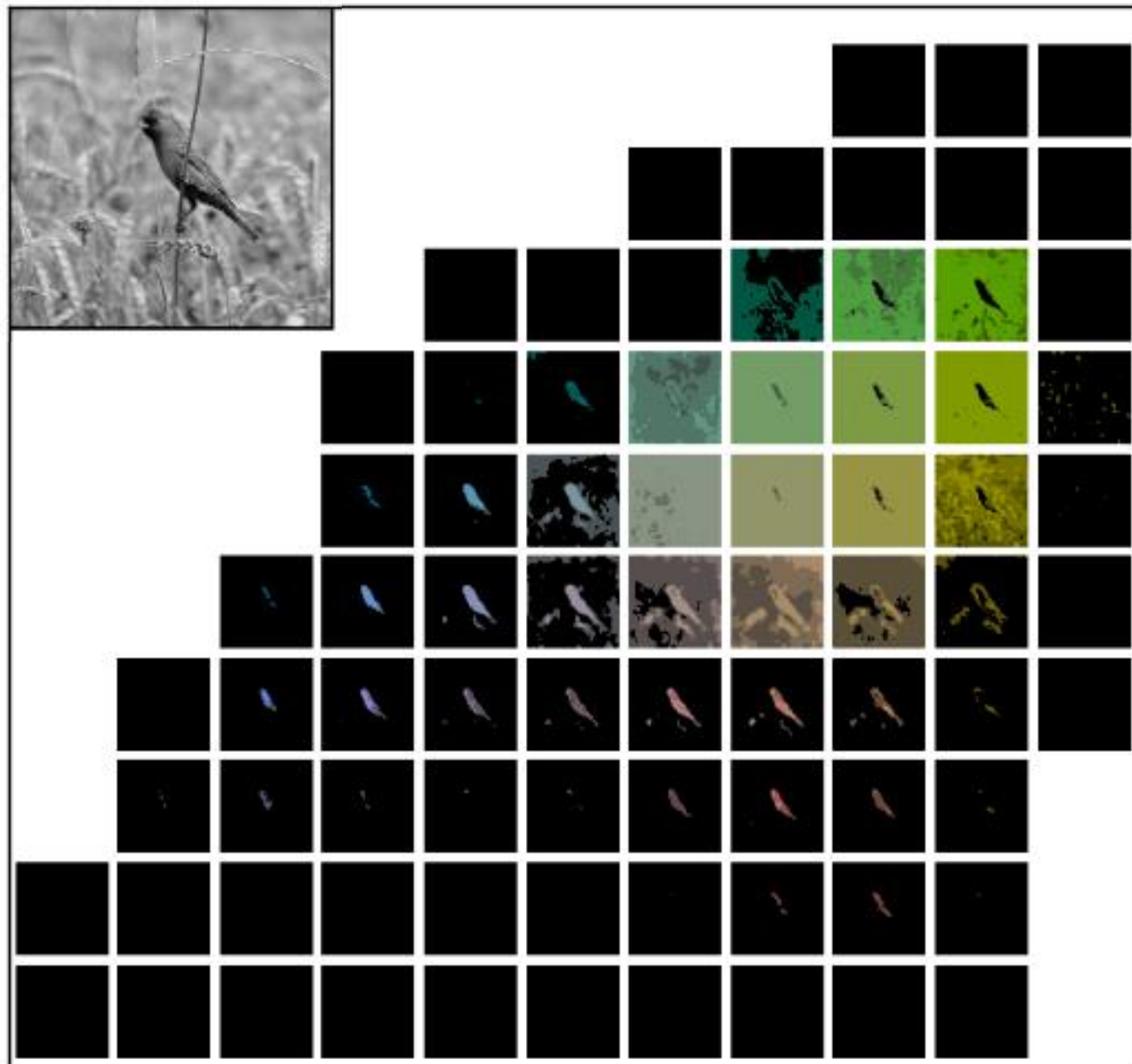
$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$

- Use **multinomial classification**

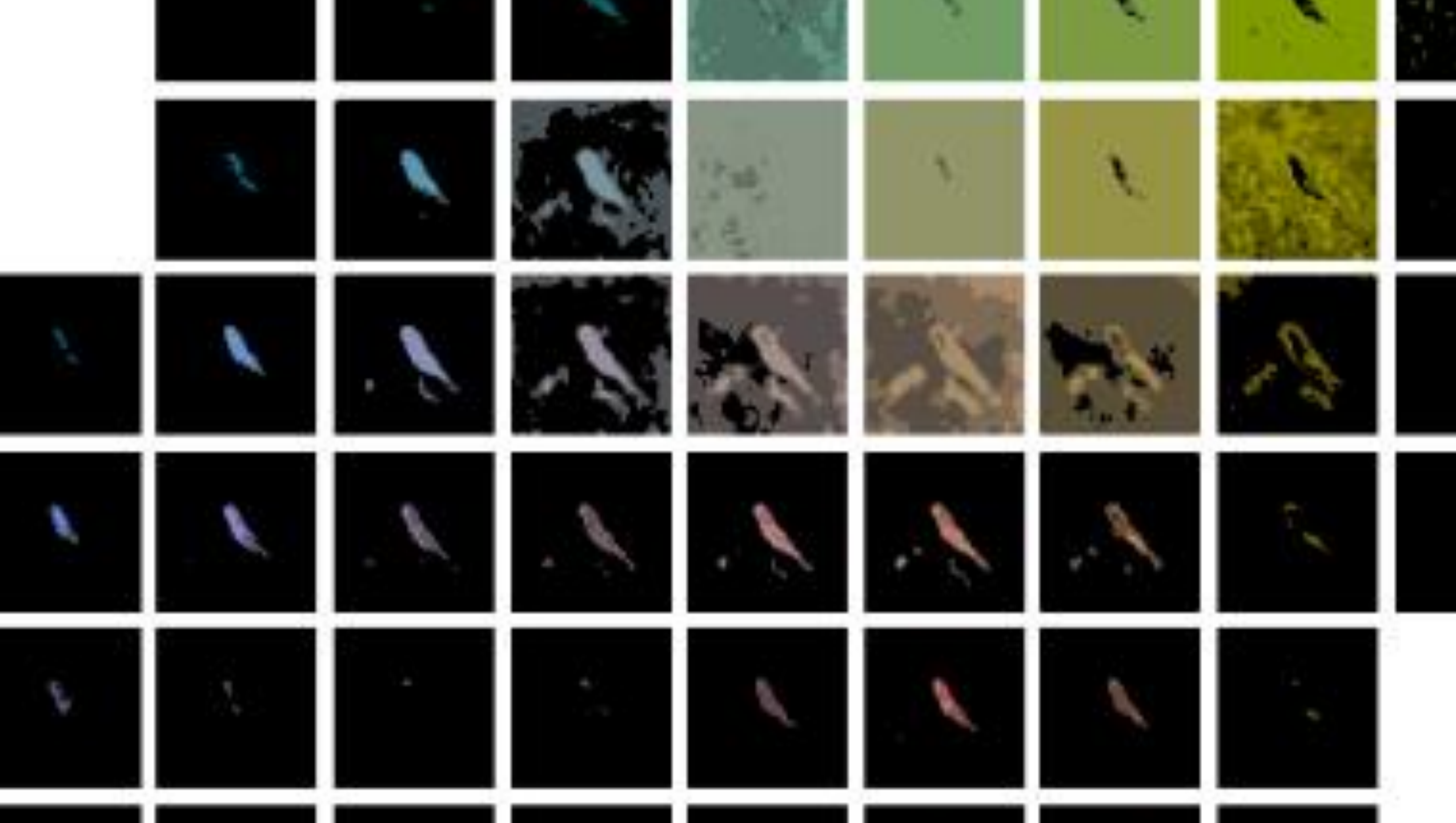
$$L(\hat{\mathbf{Z}}, \mathbf{Z}) = -\frac{1}{HW} \sum_{h,w} \sum_q \mathbf{Z}_{h,w,q} \log(\hat{\mathbf{Z}}_{h,w,q})$$

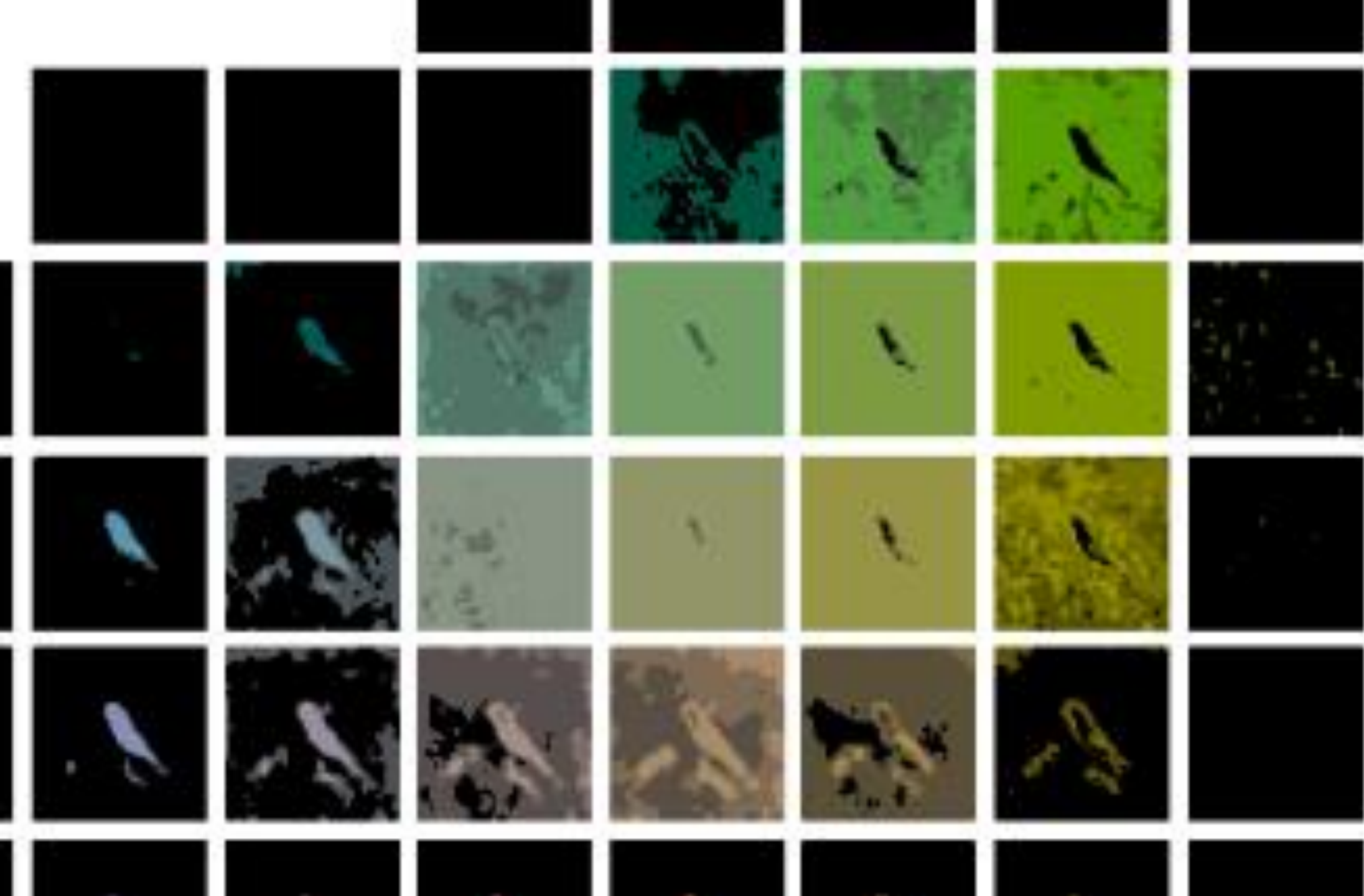
Colors in *ab* space  
(discrete)



$a$  $b$







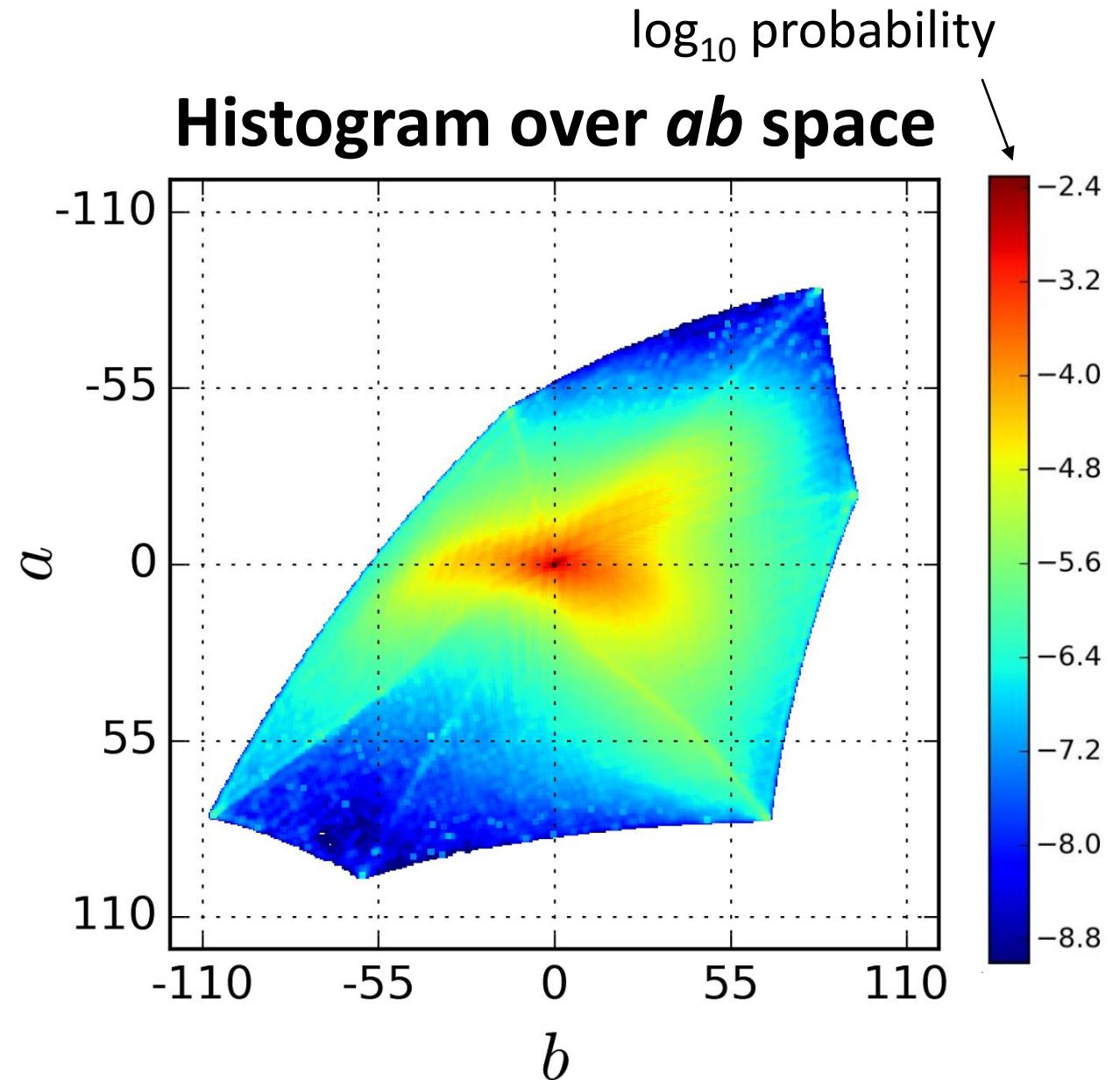
# Better Loss Function

- Regression with L2 loss inadequate

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

- Use **multinomial classification**

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



# Better Loss Function

- Regression with L2 loss inadequate

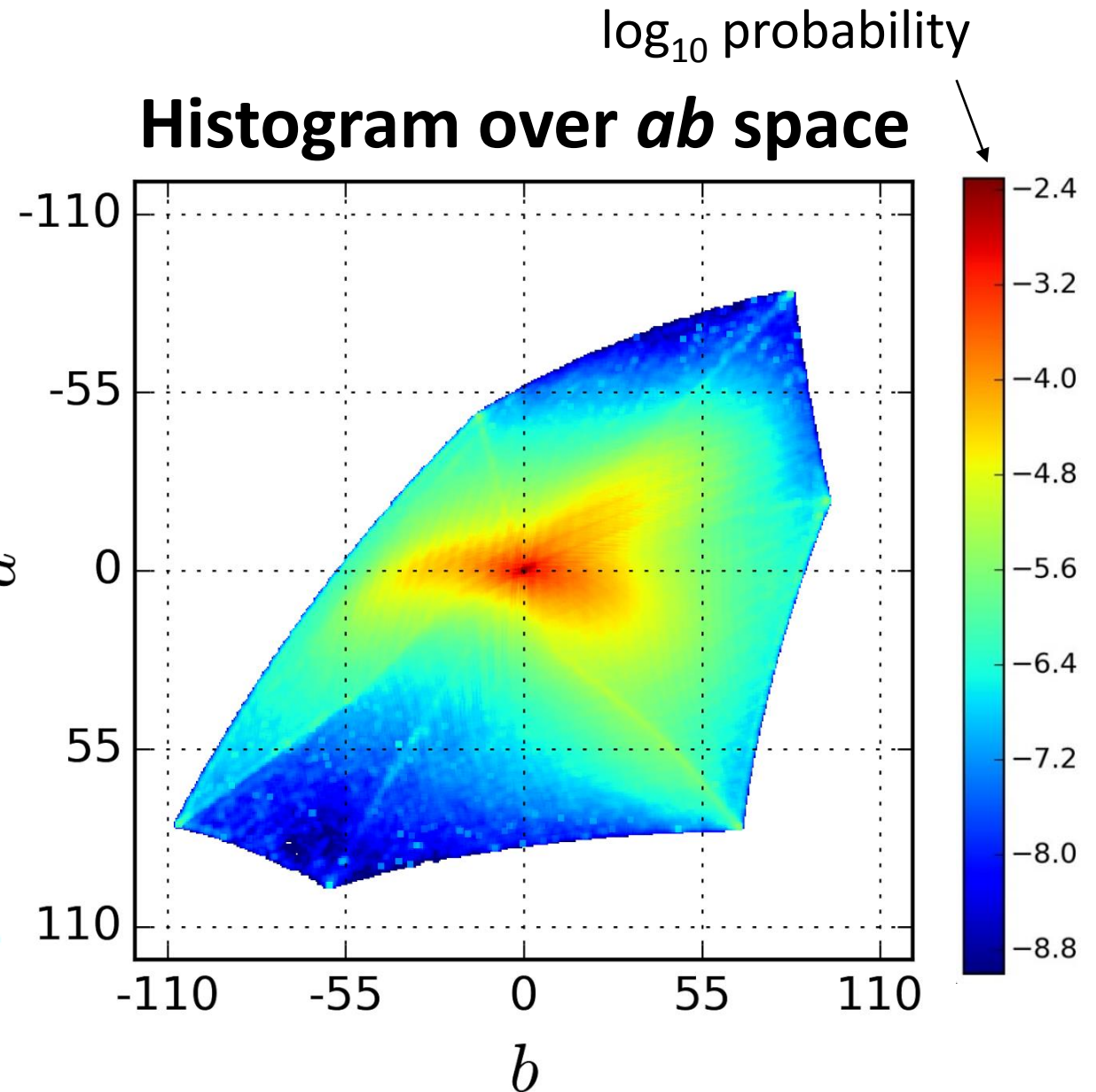
$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$

- 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 learning of *rare* colors

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



Non-

parametric

Hertzmann et al. In SIGGRAPH, 2001.

Welsh et al. In TOG, 2002.

Irony et al. In Eurographics, 2005.

Liu et al. In TOG, 2008.

Chia et al. In ACM 2011.

Gupta et al. In ACM, 2012.



Input gray image

Input objects



Input background



User input

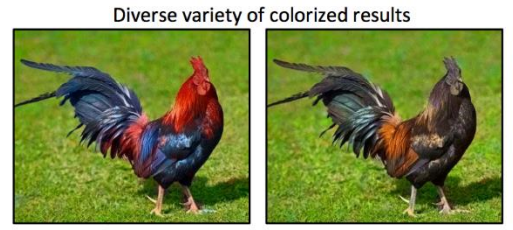
Internet objects



Internet backgrounds



Candidate image selection



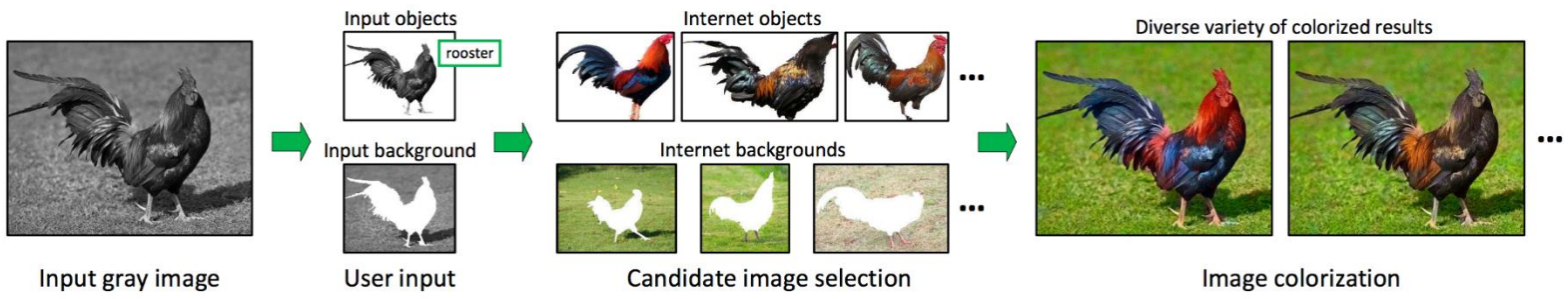
Diverse variety of colorized results

Image colorization

Non-

parametric

Hertzmann et al. In SIGGRAPH, 2001.  
Welsh et al. In TOG, 2002.  
Irony et al. In Eurographics, 2005.  
Liu et al. In TOG, 2008.  
Chia et al. In ACM 2011.  
Gupta et al. In ACM, 2012.

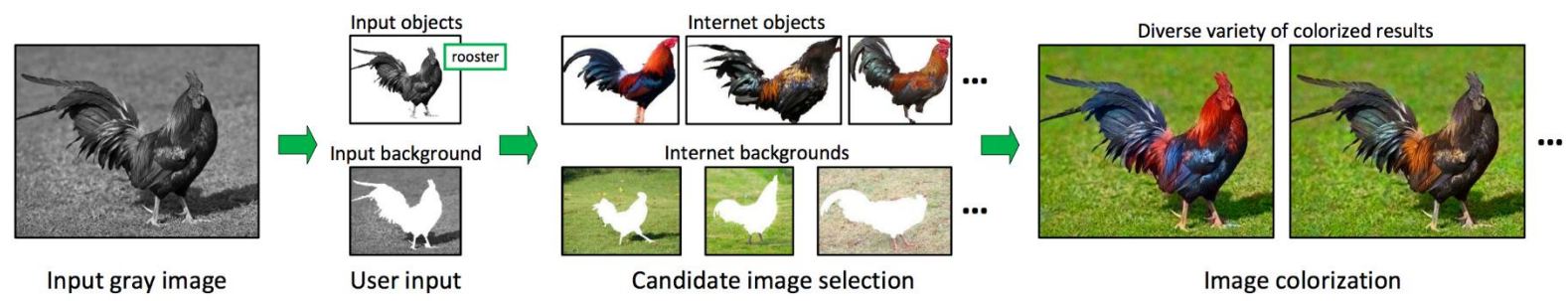


Parametric

L2 Regression

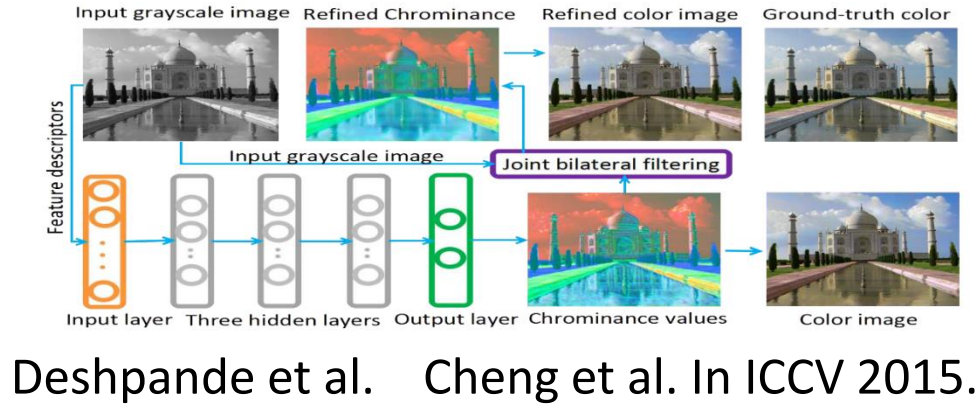
**Non-parametric**

Hertzmann et al. In SIGGRAPH, 2001.  
Welsh et al. In TOG, 2002.  
Irony et al. In Eurographics, 2005.  
Liu et al. In TOG, 2008.  
Chia et al. In ACM 2011.  
Gupta et al. In ACM, 2012.



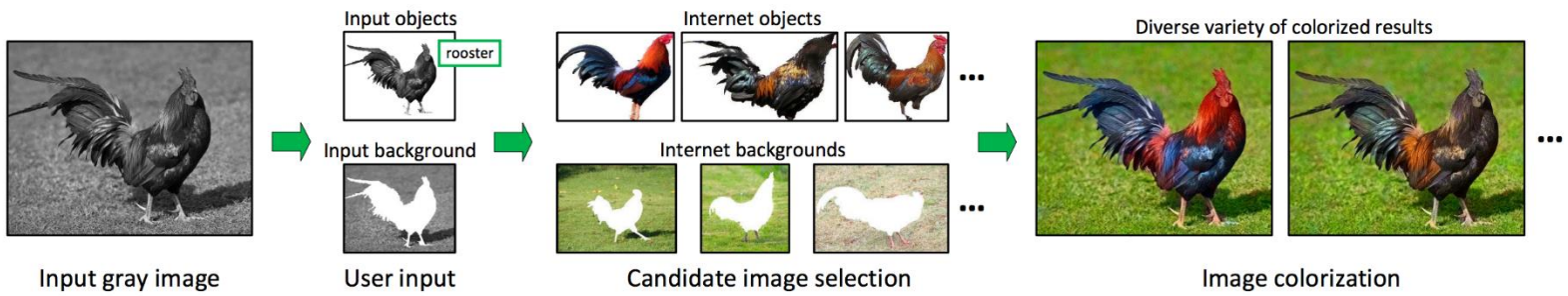
**Hand-engineered Features**

**Parametric**  
**L2 Regression**



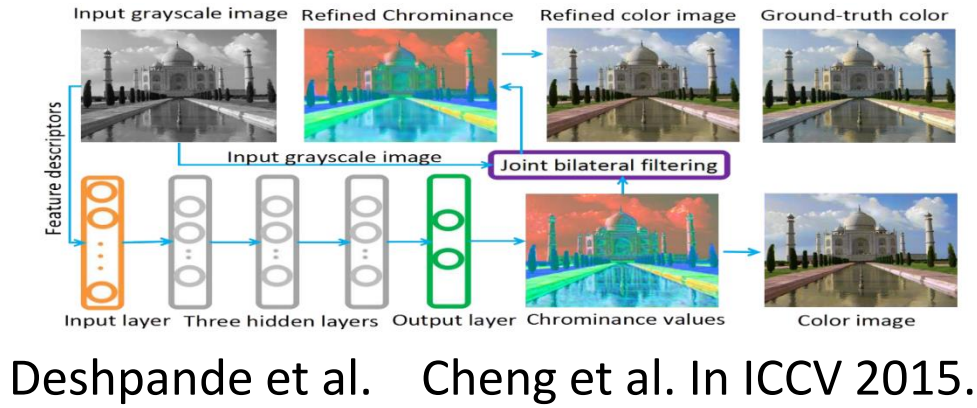
# Non-parametric

Hertzmann et al. In SIGGRAPH, 2001.  
 Welsh et al. In TOG, 2002.  
 Irony et al. In Eurographics, 2005.  
 Liu et al. In TOG, 2008.  
 Chia et al. In ACM 2011.  
 Gupta et al. In ACM, 2012.



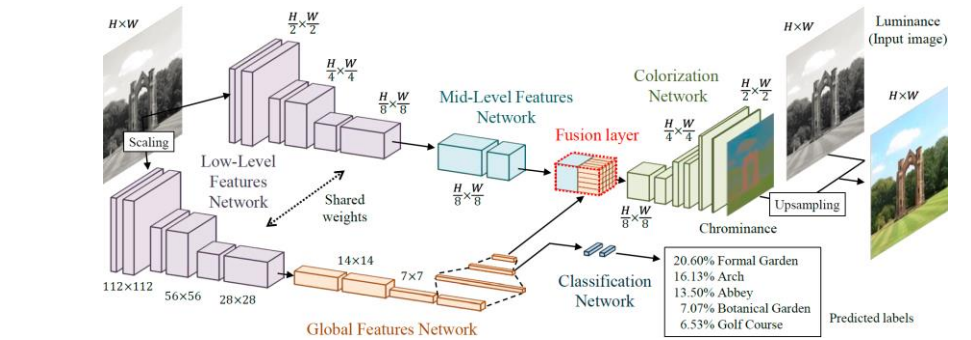
## Hand-engineered Features

### L2 Regression



# Parametric

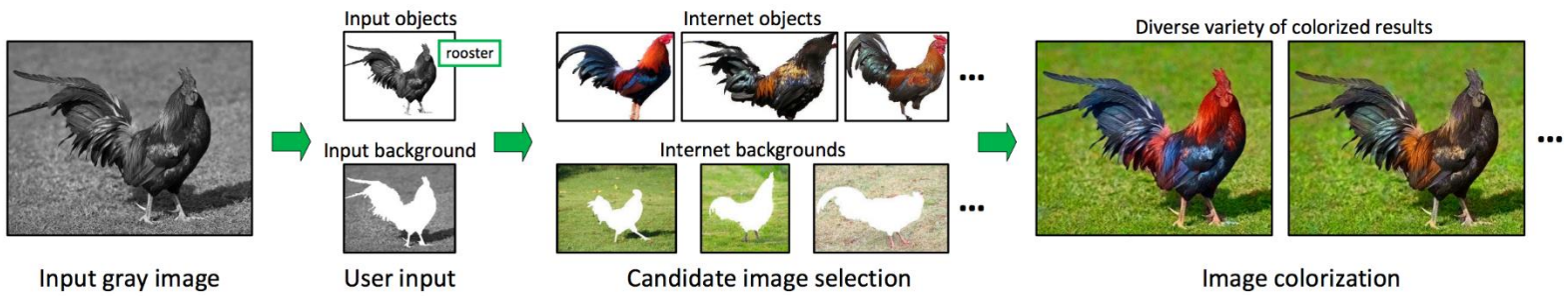
## Deep Networks



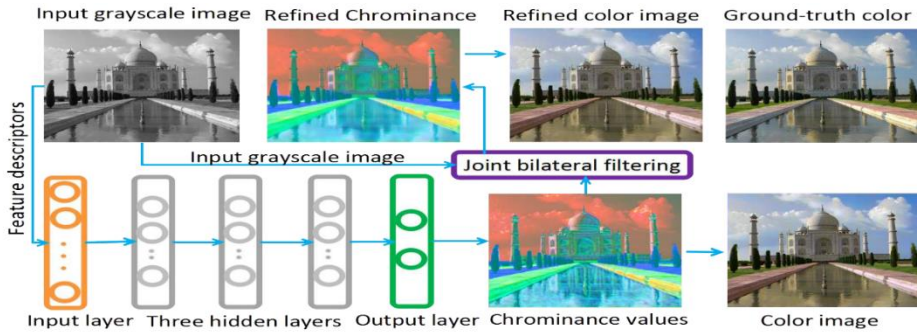


# Non-parametric

Hertzmann et al. In SIGGRAPH, 2001.  
 Welsh et al. In TOG, 2002.  
 Irony et al. In Eurographics, 2005.  
 Liu et al. In TOG, 2008.  
 Chia et al. In ACM 2011.  
 Gupta et al. In ACM, 2012.

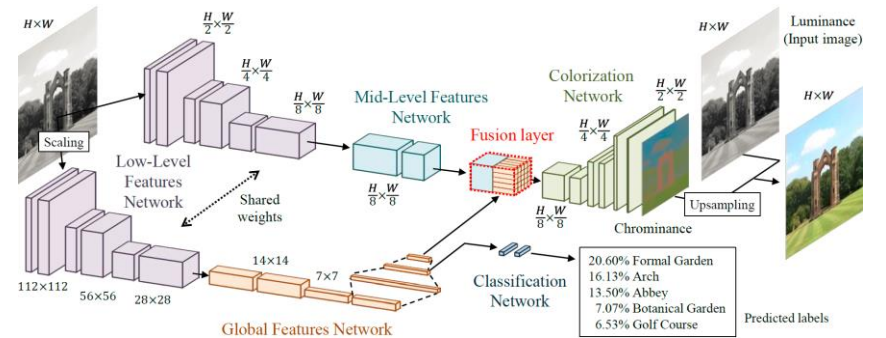


## Hand-engineered Features



Deshpande et al. Cheng et al. In ICCV 2015.

## Deep Networks

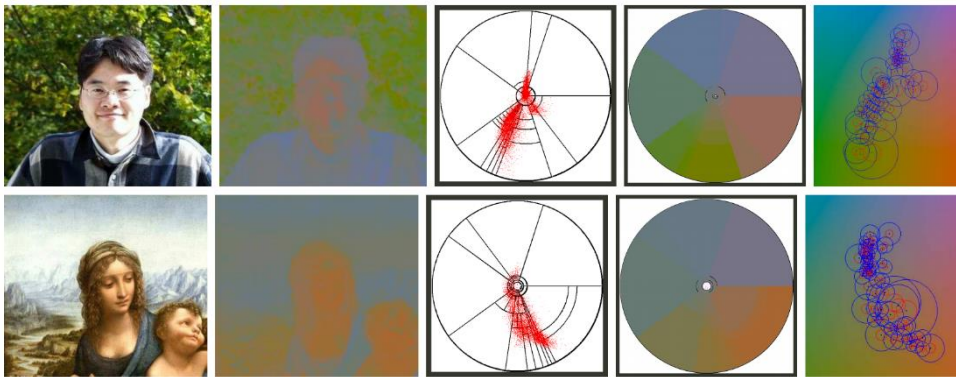


Dahl. Jan 2016. Iizuka et al. In SIGGRAPH, 2016.

# Parametric

## L2 Regression

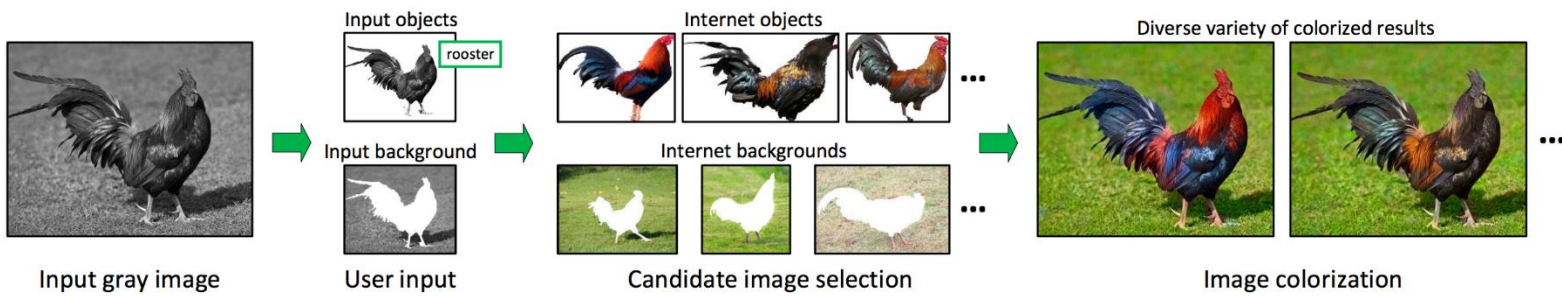
## Classification



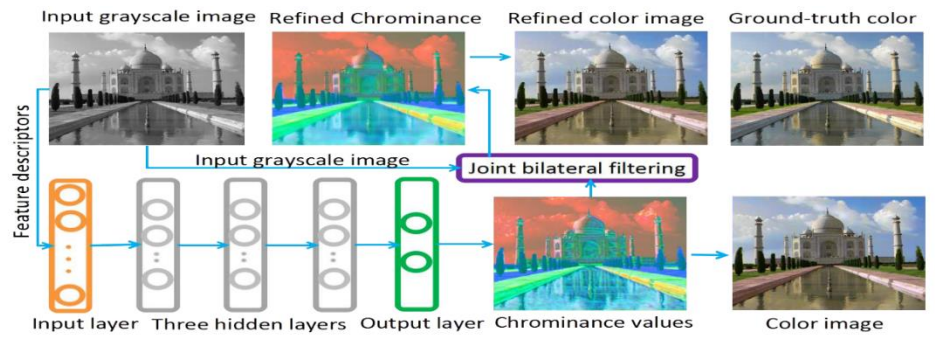
Charpiat et al. In ECCV 2008.

# Non-parametric

Hertzmann et al. In SIGGRAPH, 2001.  
 Welsh et al. In TOG, 2002.  
 Irony et al. In Eurographics, 2005.  
 Liu et al. In TOG, 2008.  
 Chia et al. In ACM 2011.  
 Gupta et al. In ACM, 2012.

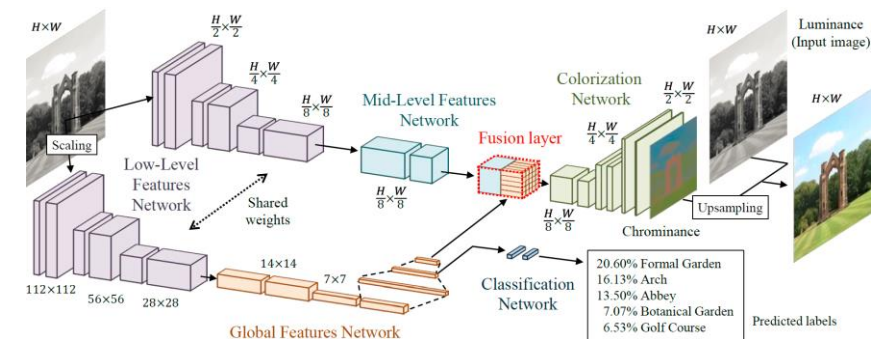


## Hand-engineered Features



Deshpande et al. Cheng et al. In ICCV 2015.

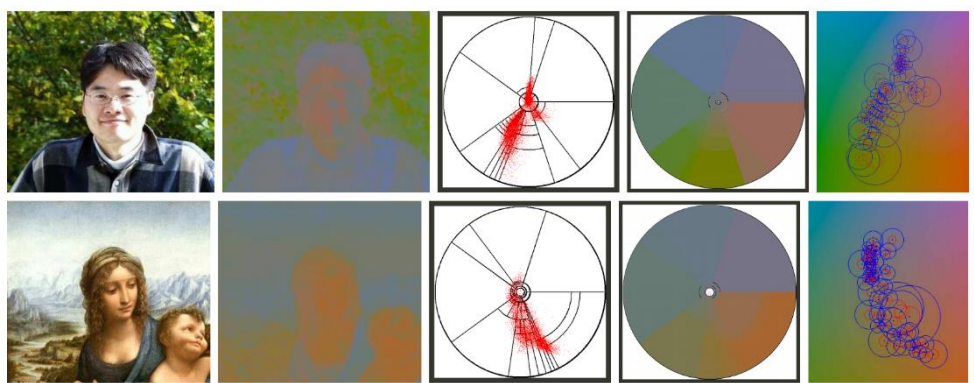
## Deep Networks



Dahl. Jan 2016. Iizuka et al. In SIGGRAPH, 2016.

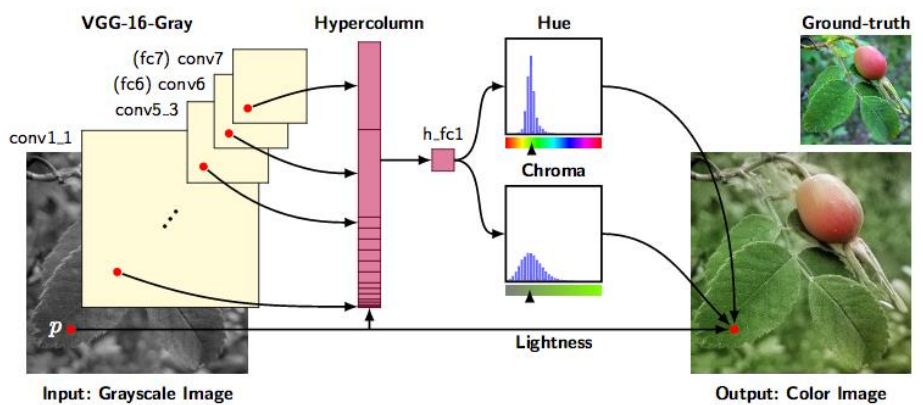
# Parametric

## L2 Regression



Charpiat et al. In ECCV 2008.

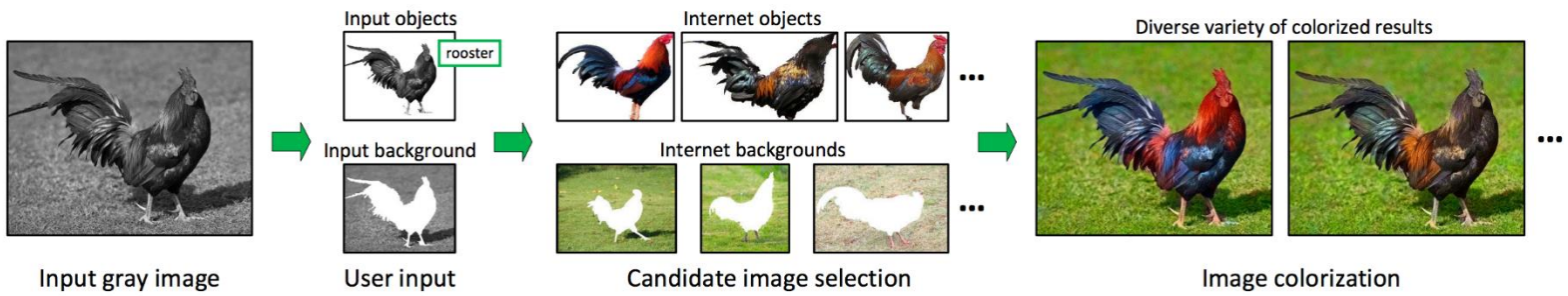
## Classification



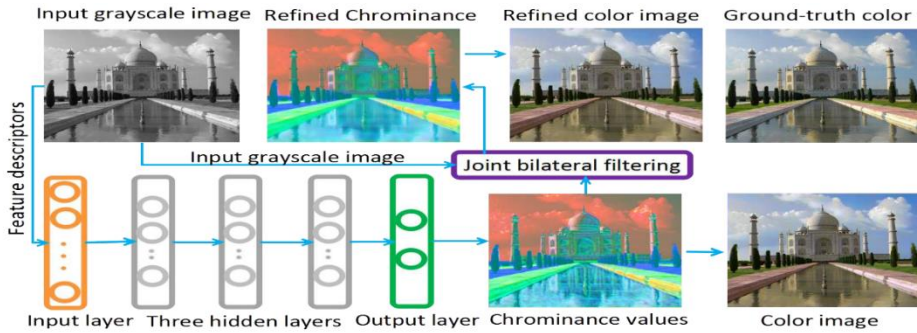
Larsson et al. In ECCV 2016. [Concurrent]

# Non-parametric

Hertzmann et al. In SIGGRAPH, 2001.  
 Welsh et al. In TOG, 2002.  
 Irony et al. In Eurographics, 2005.  
 Liu et al. In TOG, 2008.  
 Chia et al. In ACM 2011.  
 Gupta et al. In ACM, 2012.

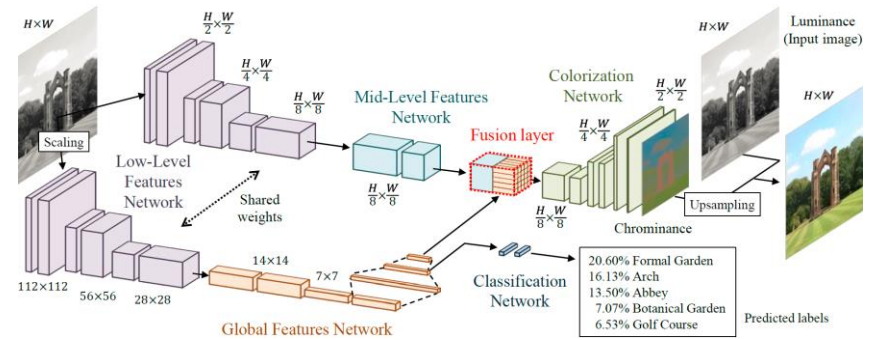


## Hand-engineered Features



Deshpande et al. Cheng et al. In ICCV 2015.

## Deep Networks

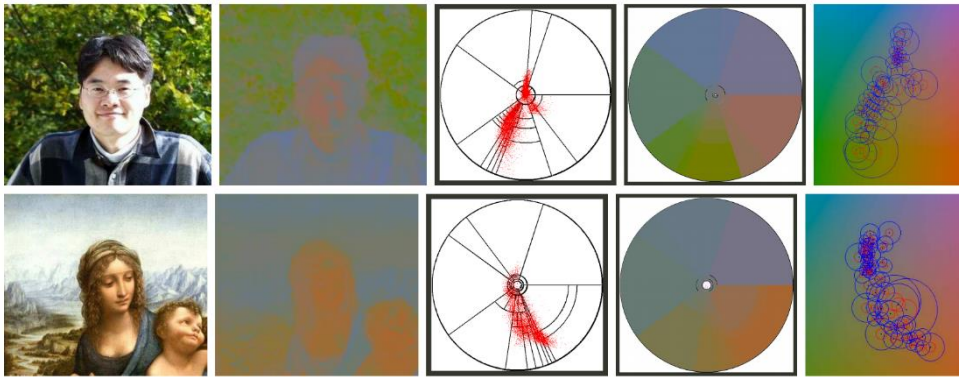


Dahl. Jan 2016. Iizuka et al. In SIGGRAPH, 2016.

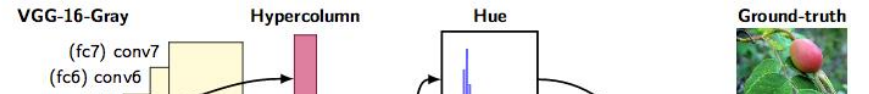
# Parametric

## L2 Regression

## Classification



Charpiat et al. In ECCV 2008.

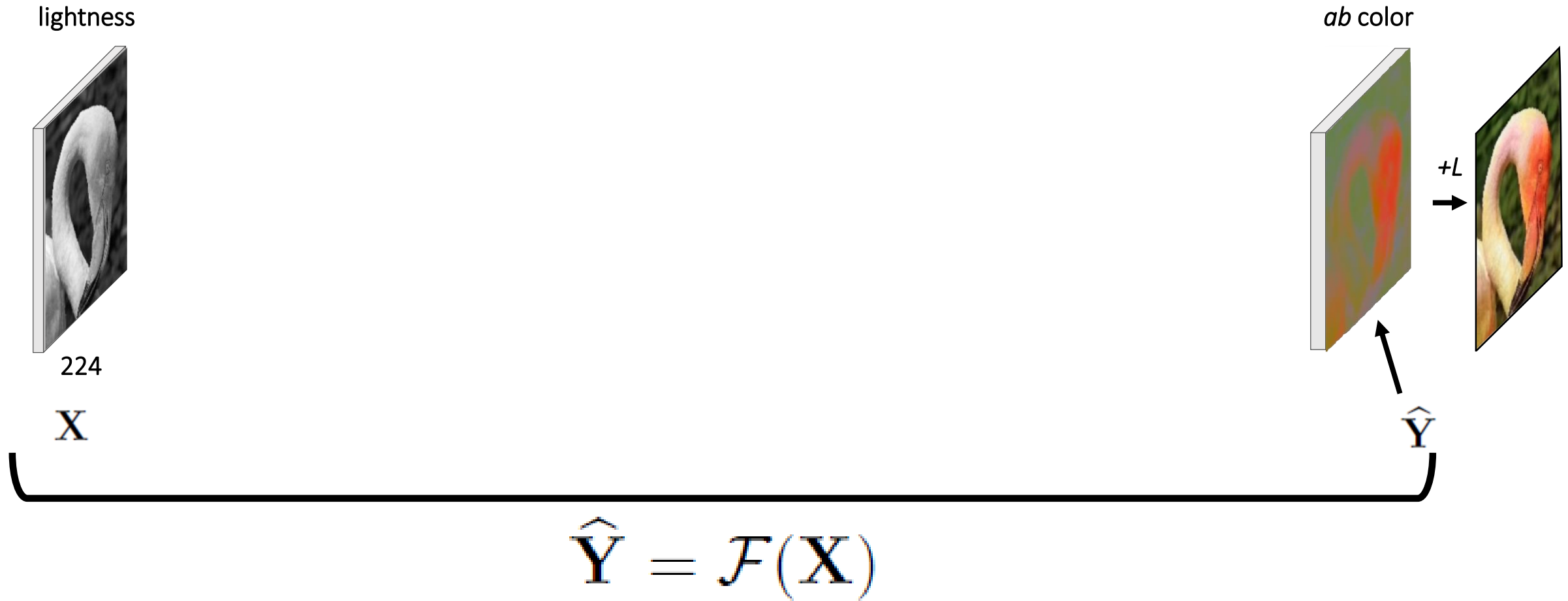


**Upcoming Oral O-3A-04  
 Tomorrow, 9–10 AM**

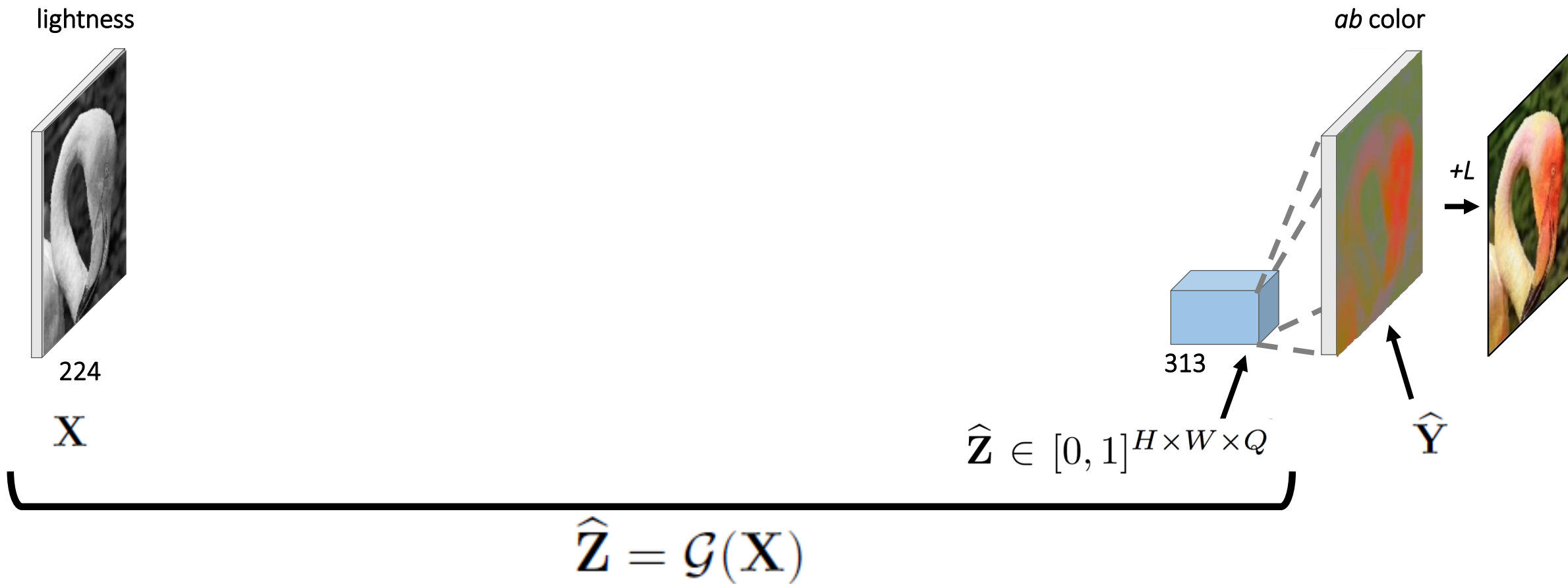
Input: Grayscale Image Output: Color Image

Larsson et al. In ECCV 2016. [Concurrent]

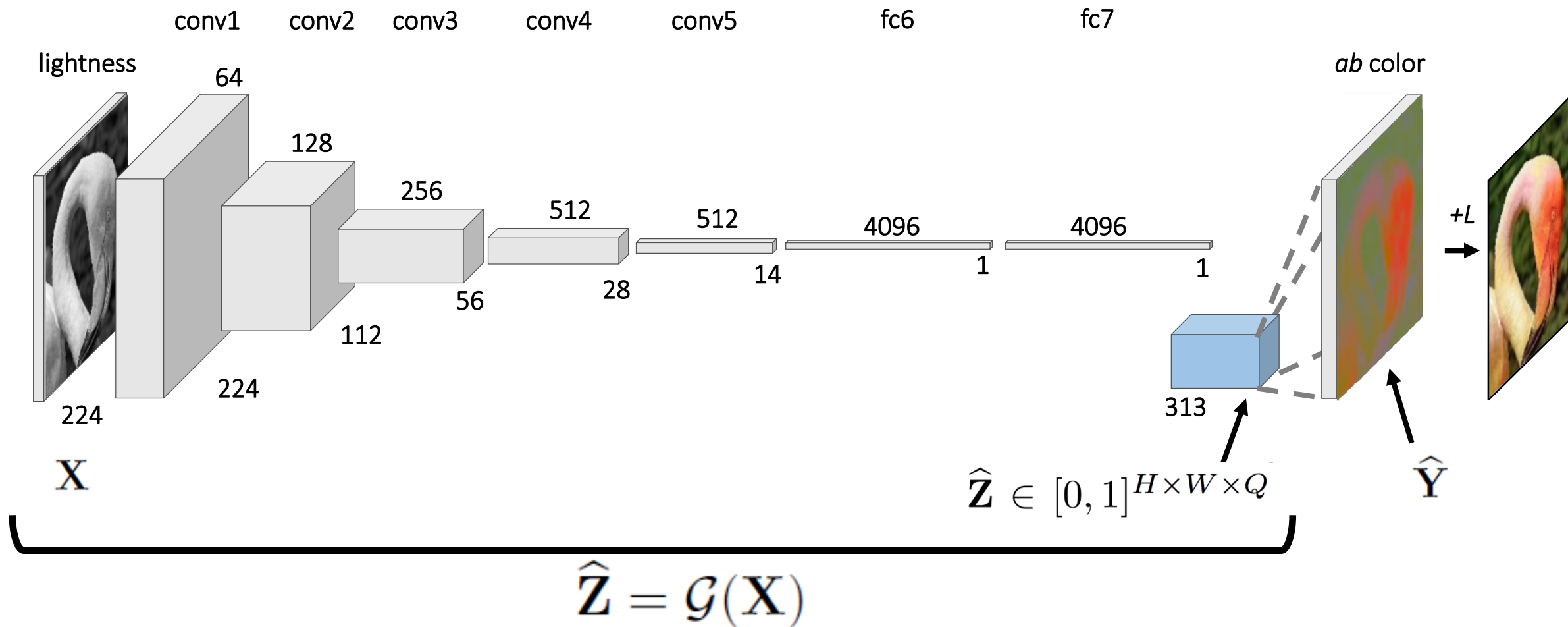
# Network Architecture



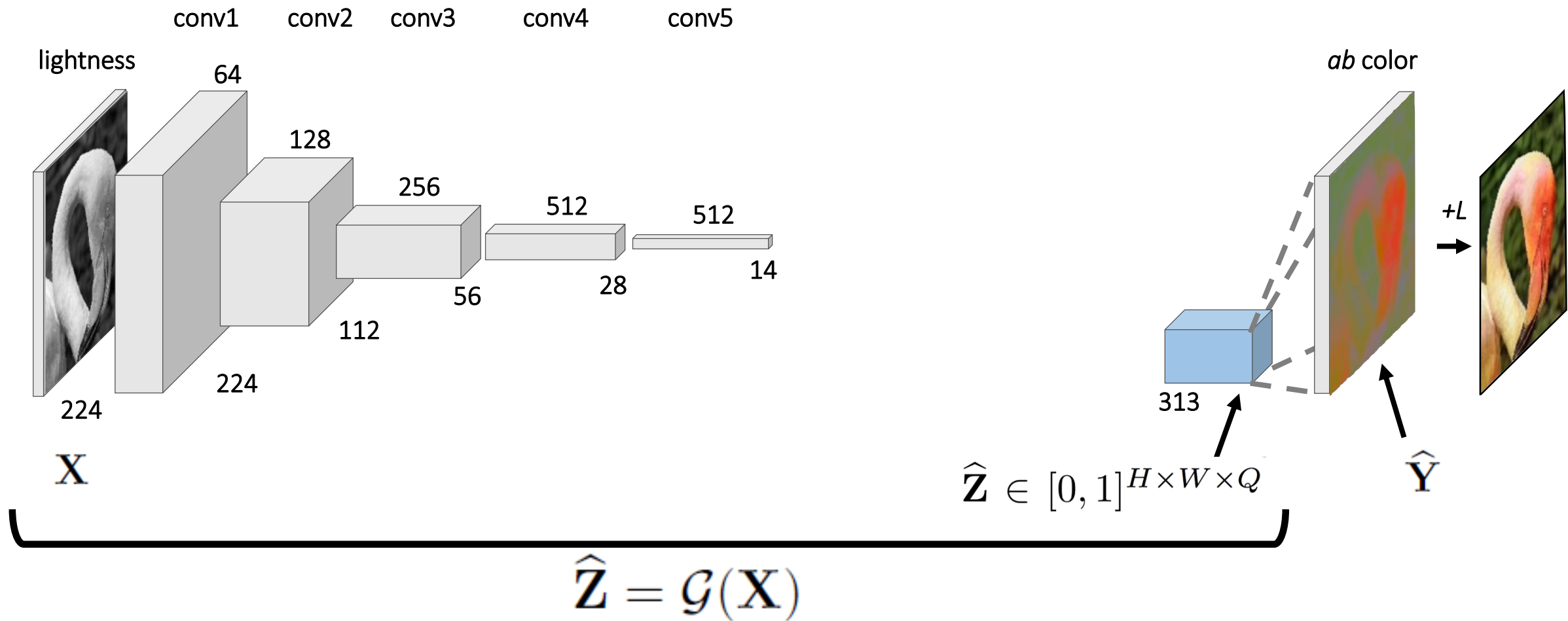
# Network Architecture



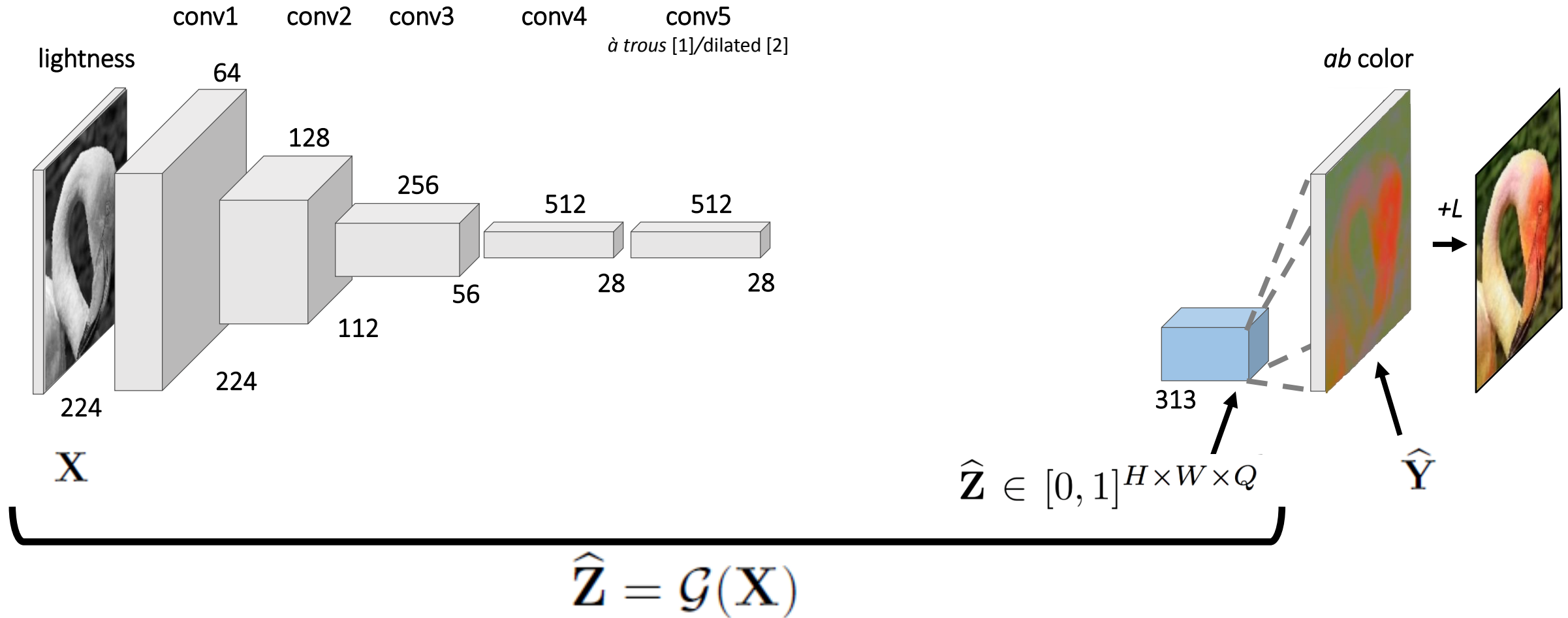
# Network Architecture



# Network Architecture



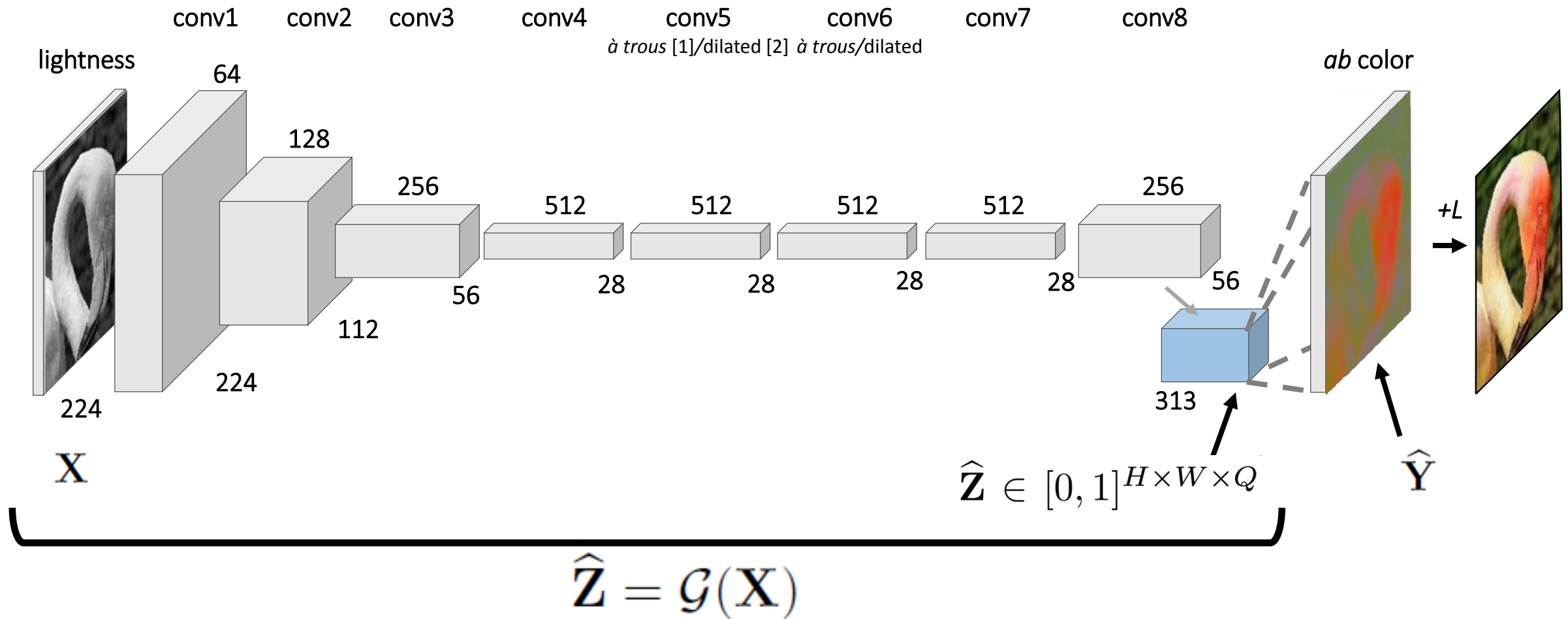
# Network Architecture



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

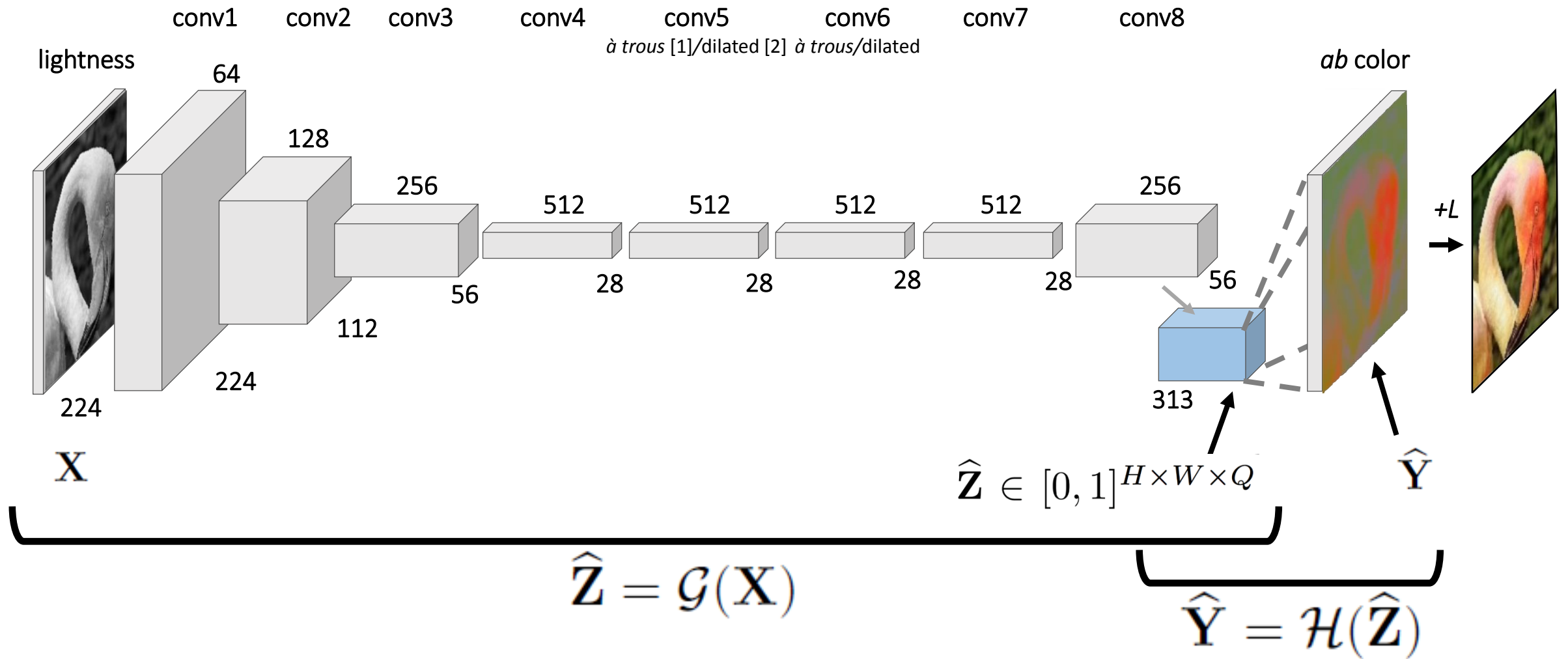


# Network Architecture



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

# Network Architecture



[1] Chen *et al.* In arXiv, 2016.

[2] Yu and Koltun. In ICLR, 2016

# Input



Input



L2 Regression



Input



L2 Regression



Class w/ Rebalancing



Input



L2 Regression



Class w/ Rebalancing



Input



L2 Regression



Class w/ Rebalancing



Input

L2 Regression

Class w/ Rebalancing





Ground Truth

L2 Regression

Class w/ Rebalancing



# Failure Cases



# Biases



# Evaluation

**Visual Quality**

Per-pixel accuracy

**Quantitative**

# Evaluation

|                     | <b>Visual Quality</b>   |
|---------------------|---|
| <b>Quantitative</b> | Per-pixel accuracy<br>Perceptual realism<br>Semantic interpretability |
| <b>Qualitative</b>  | Low-level stimuli<br>Legacy grayscale photos                          |

# Evaluation

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

# Evaluation

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

# Evaluation

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

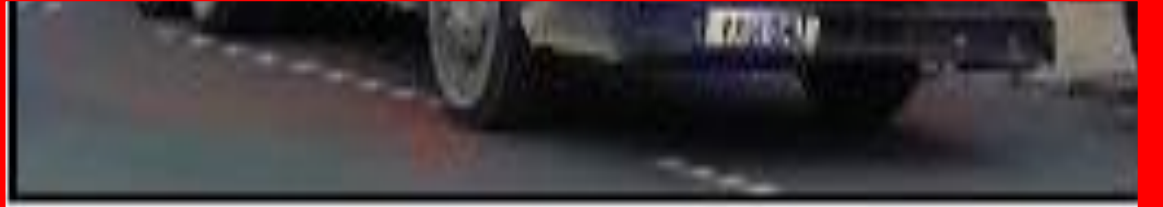
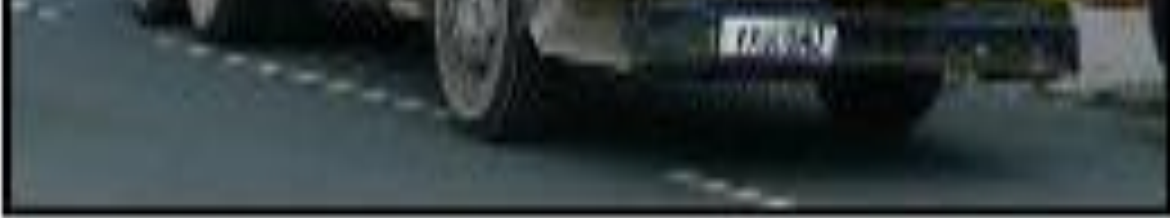


Perceptual Realism / Amazon Mechanical Turk Test

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

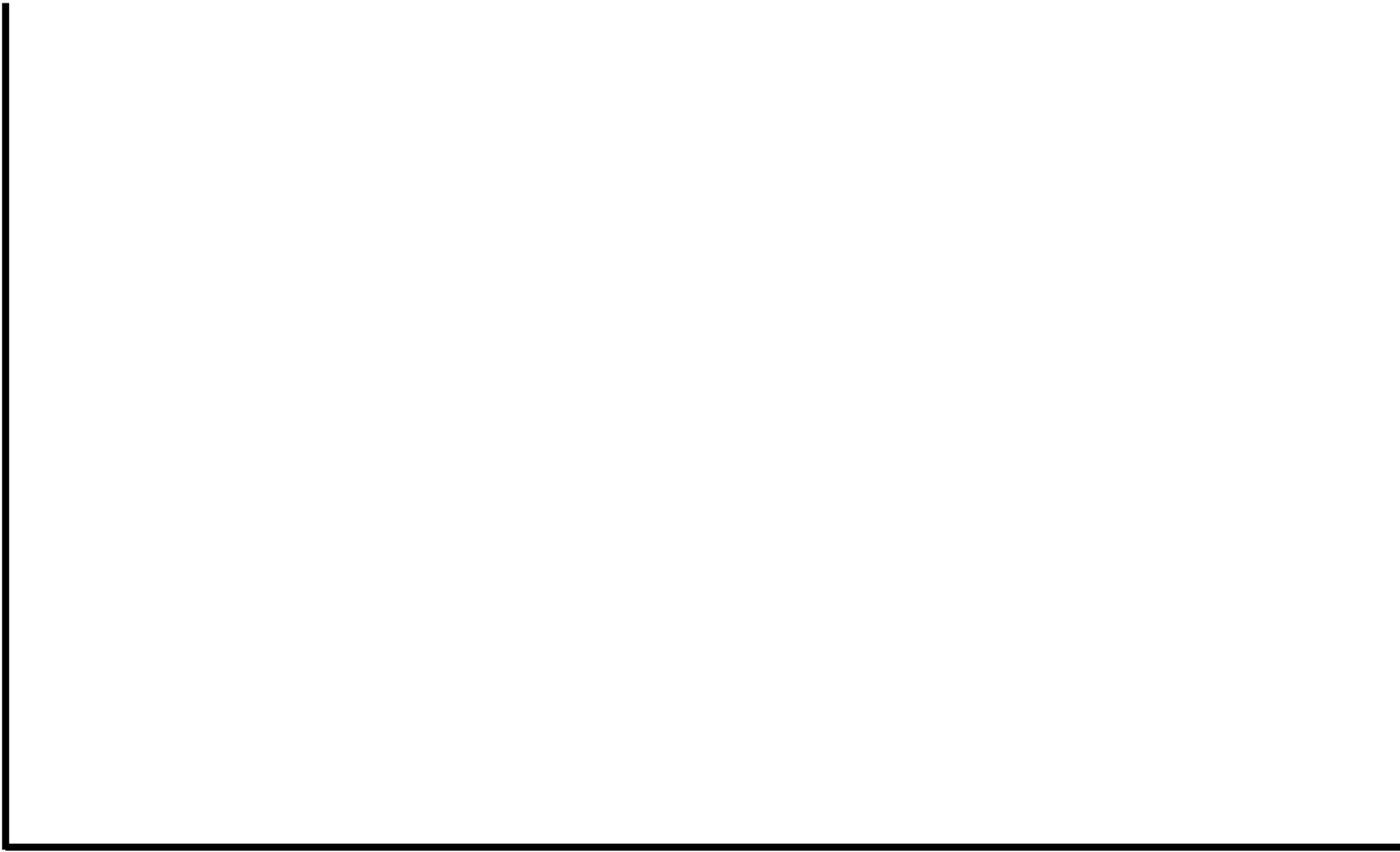
# Perceptual Realism Test

50%

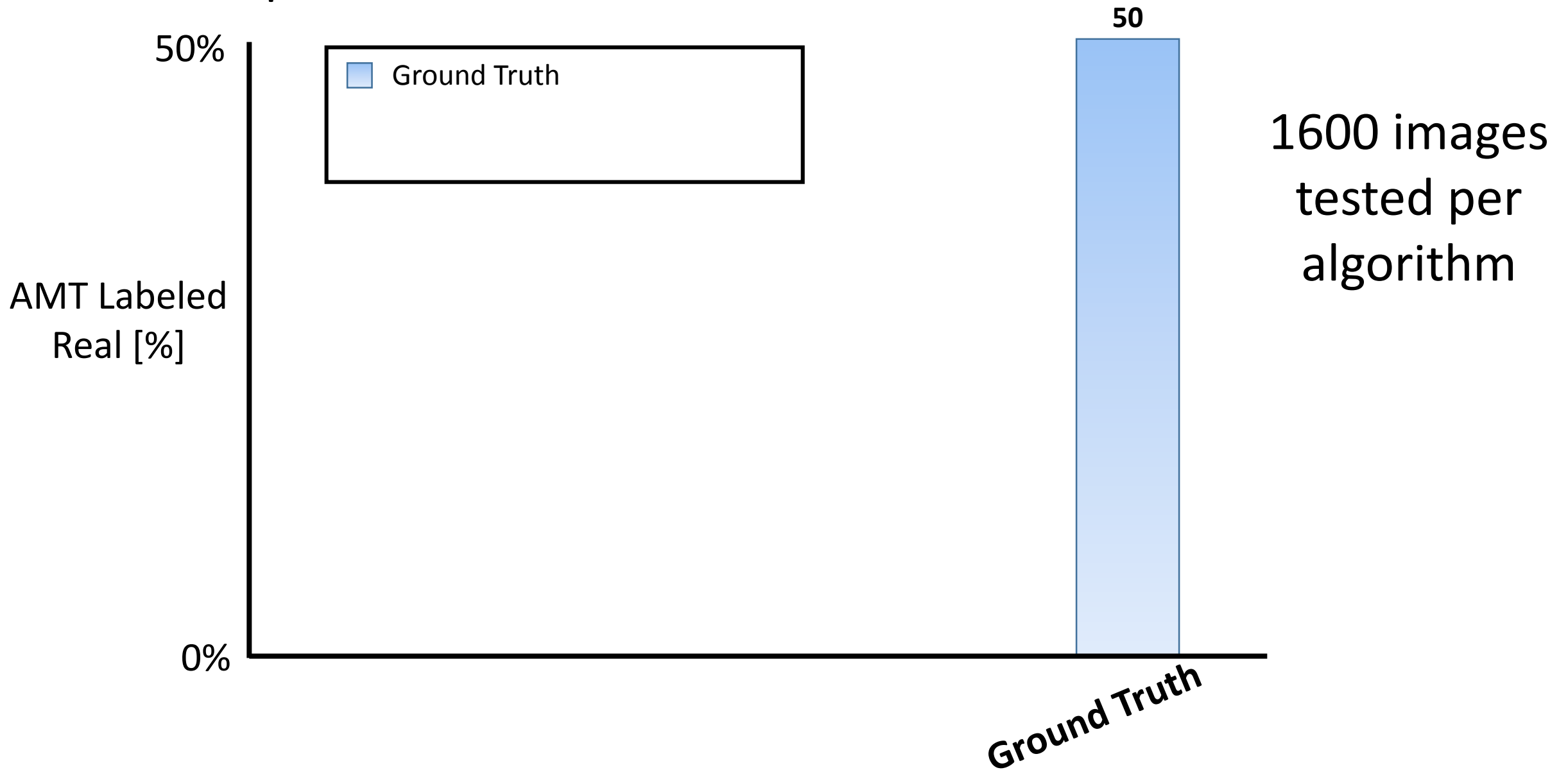
0%

AMT Labeled  
Real [%]

1600 images  
tested per  
algorithm

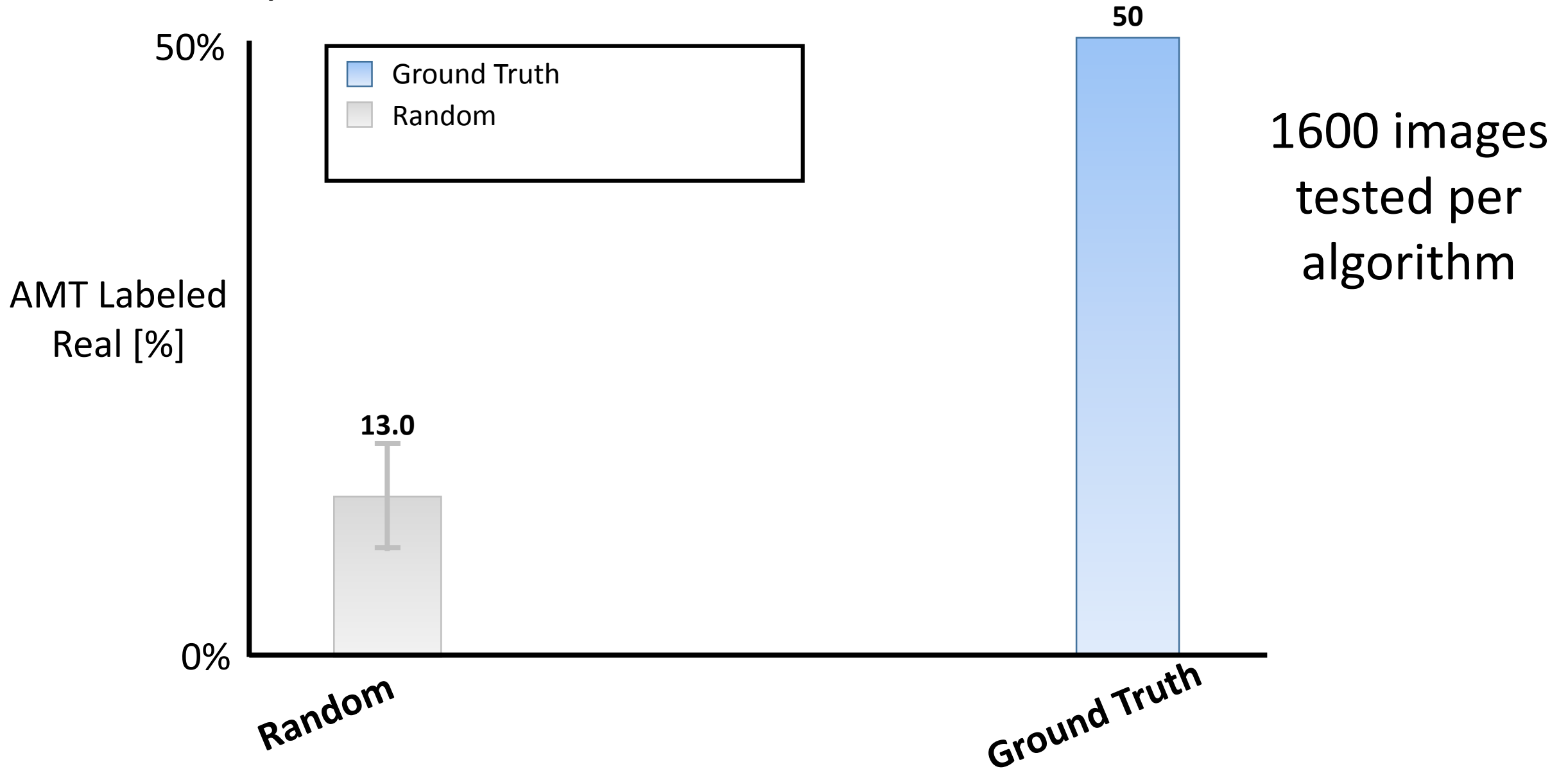


# Perceptual Realism Test

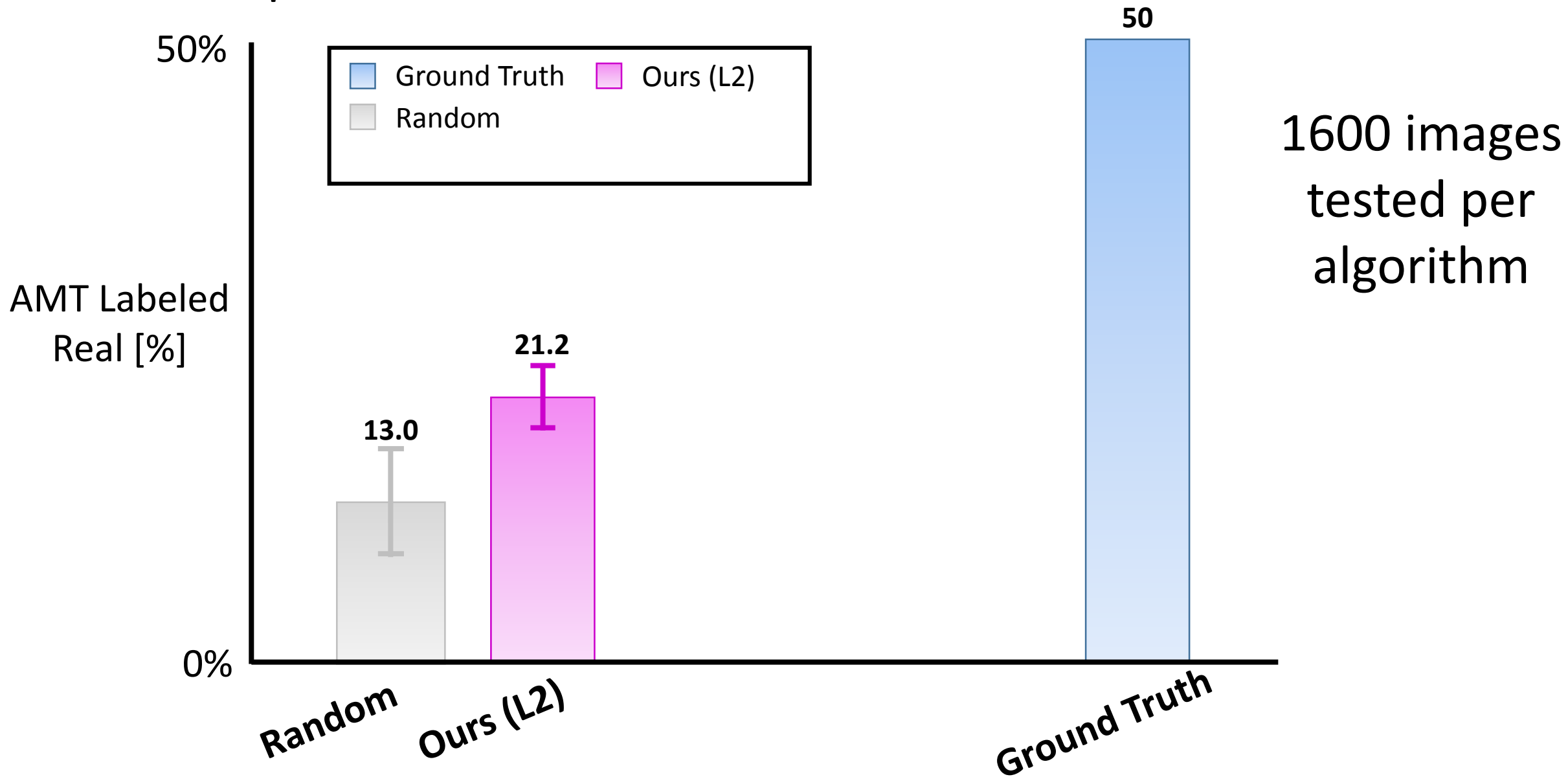




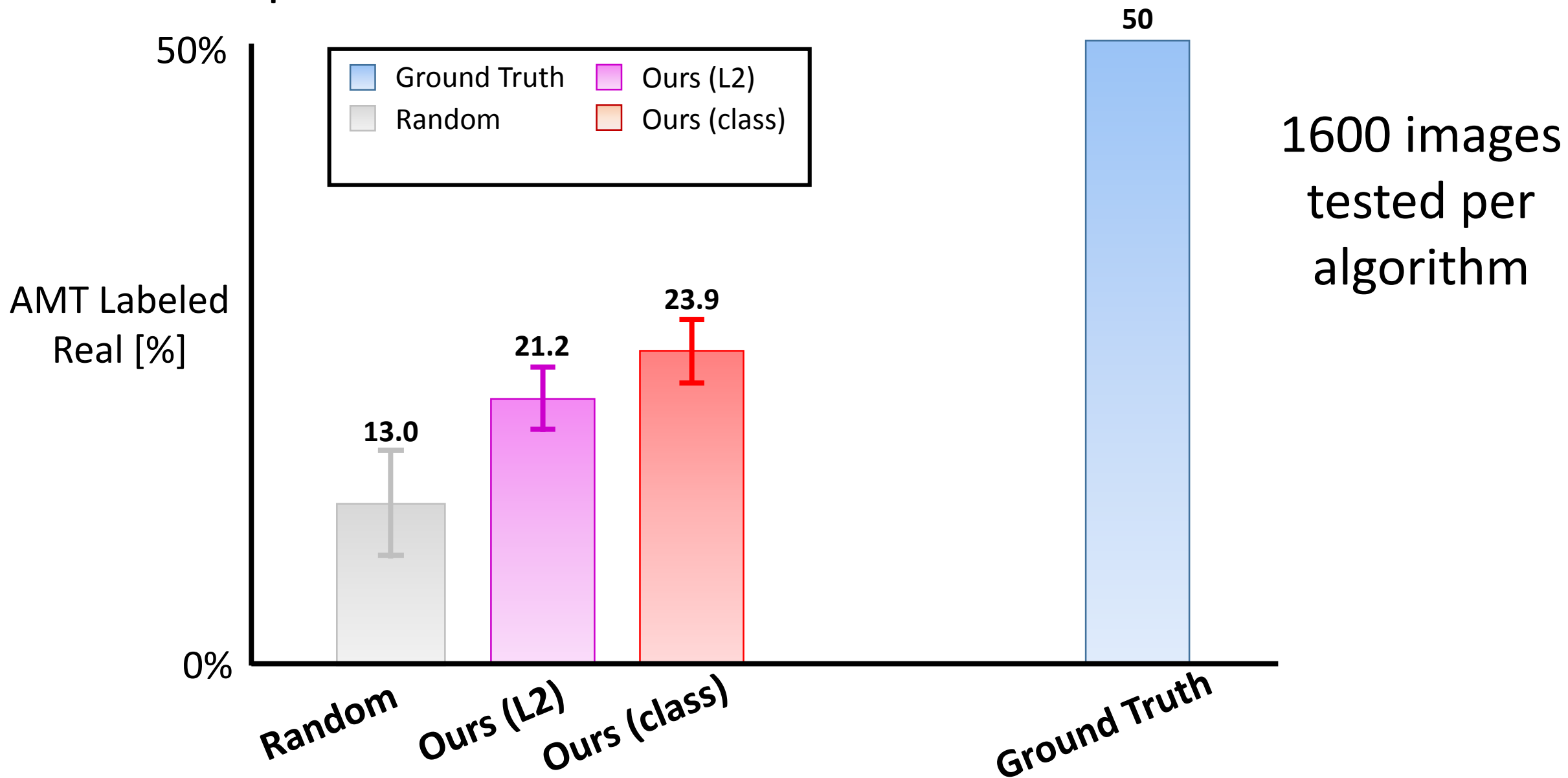
# Perceptual Realism Test



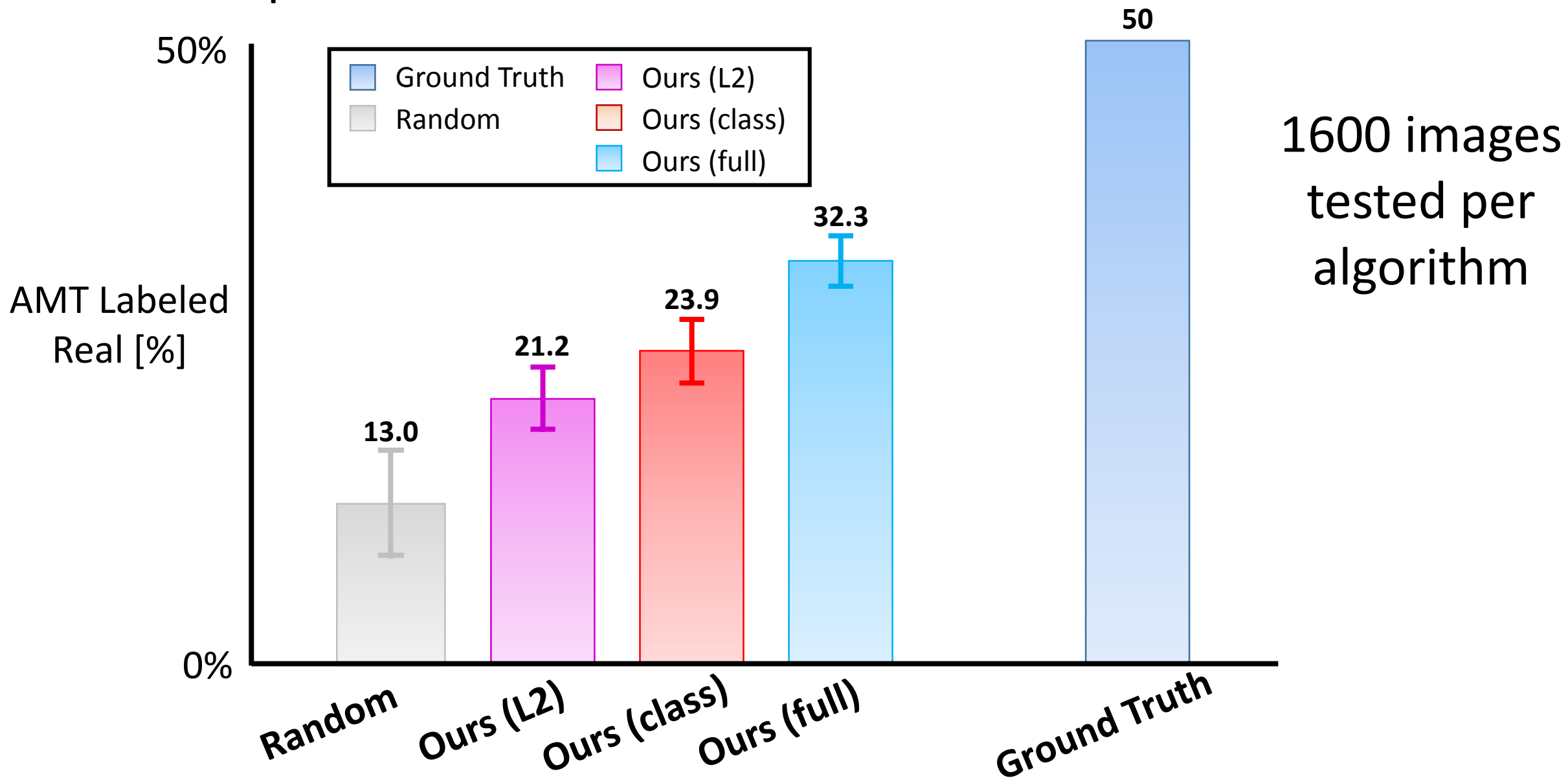
# Perceptual Realism Test



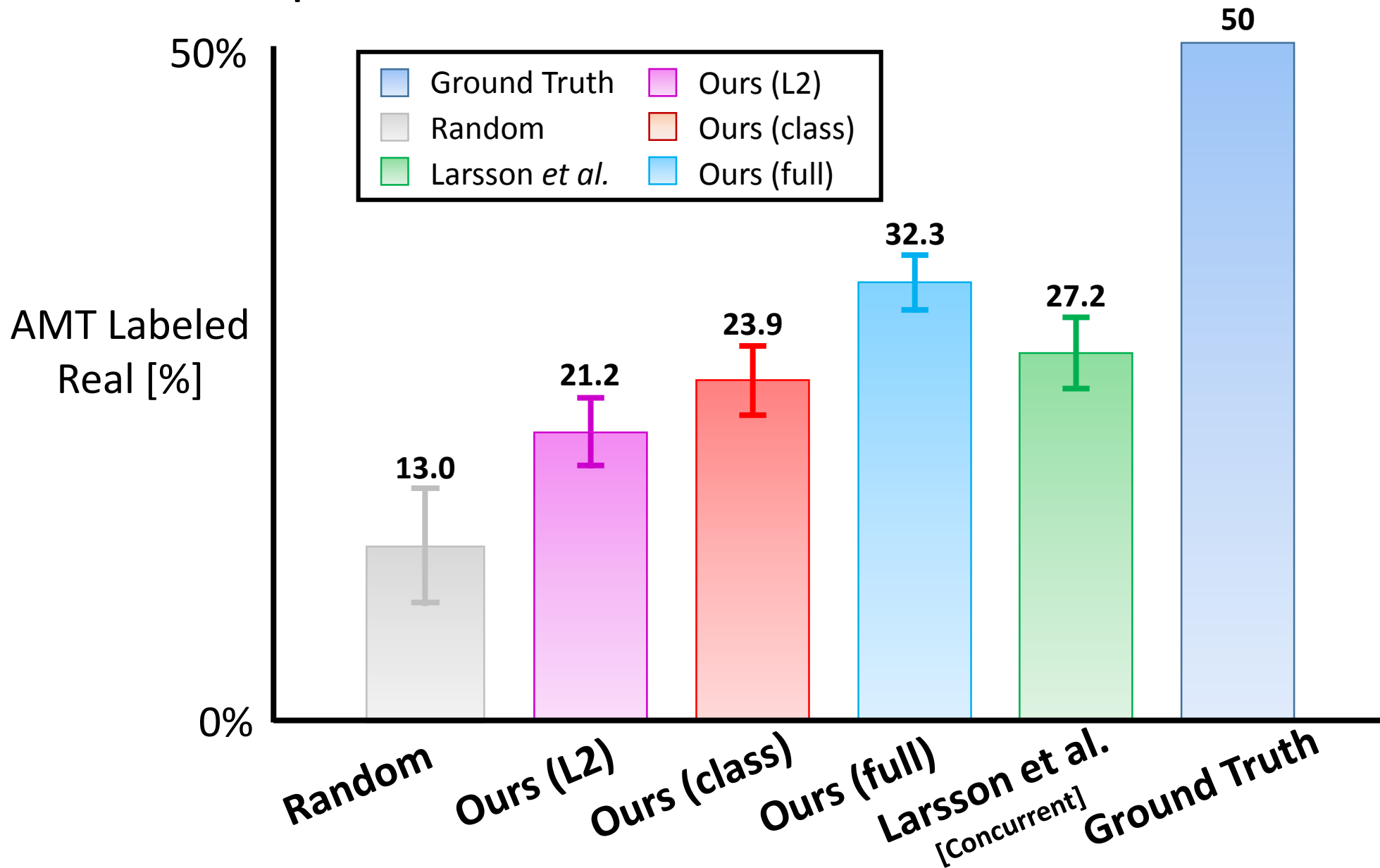
# Perceptual Realism Test



# Perceptual Realism Test



# Perceptual Realism Test



1600 images  
tested per  
algorithm

**Input**



**Ground Truth**



**Input**



**Ground Truth**



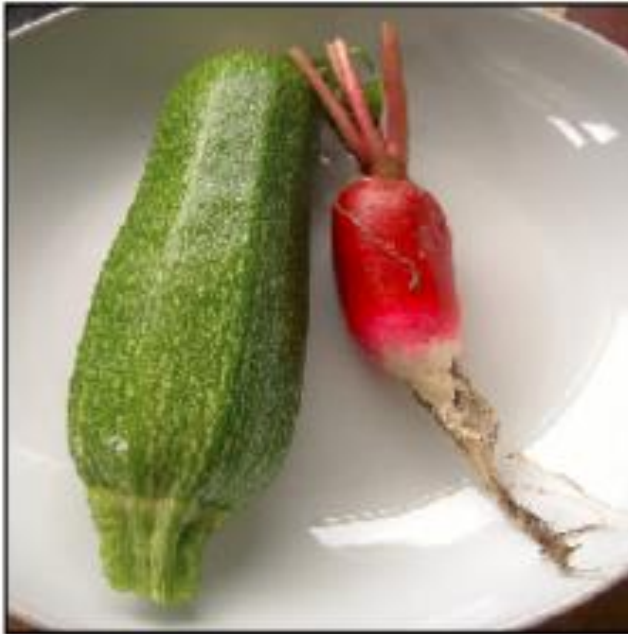
**Output**



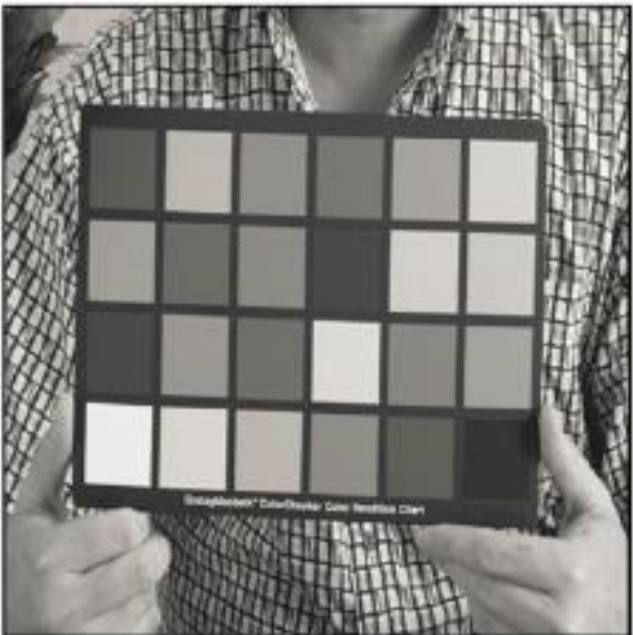
# Input



# Ground Truth



# Output

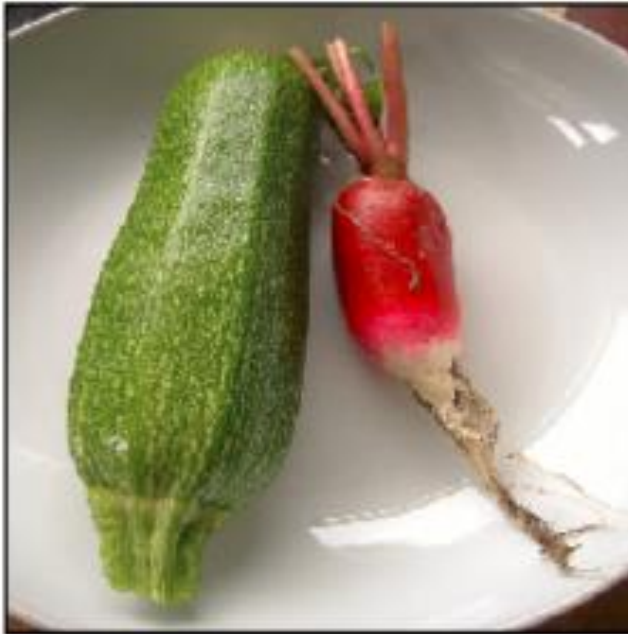




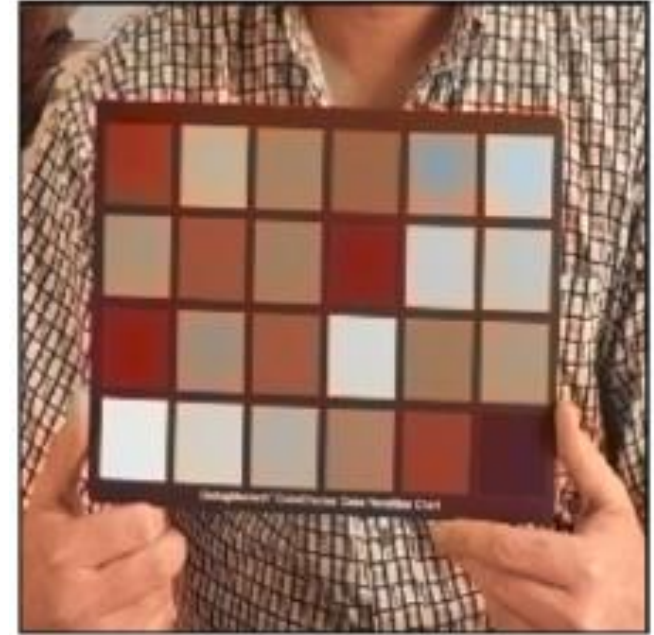
# Input



# Ground Truth



# Output



# Predicting Labels from Data

**Supervised  
training**

**Data  $x$**

**ImageNet  
images**



**Learned feature  
hierarchy**

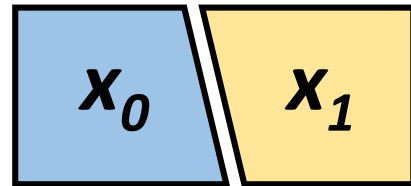


**Label  $y$**

**ImageNet  
labels**

# Predicting Data from Data

**Supervised  
training**



**ImageNet  
images**



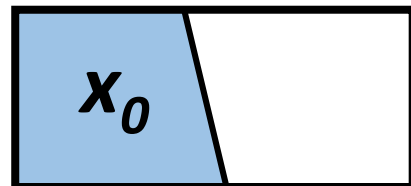
**Learned feature  
hierarchy**



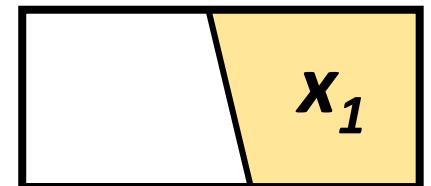
**Label  $y$**

**ImageNet  
labels**

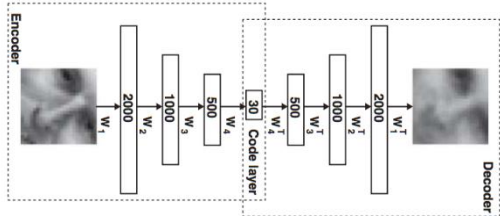
**Unsupervised/  
Self-supervised  
training**



**Learned feature  
hierarchy**

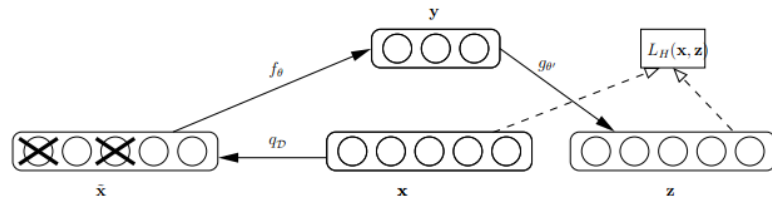


## Autoencoders



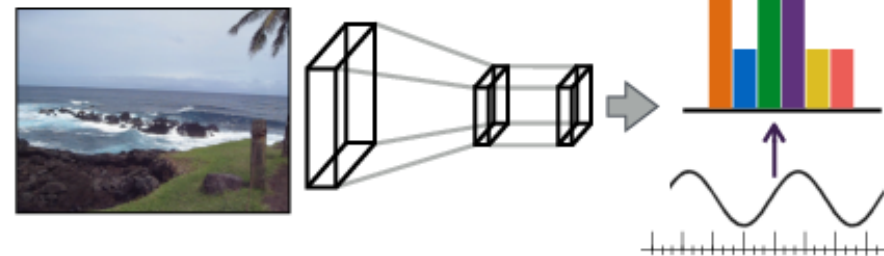
Hinton & Salakhutdinov.  
Science 2006.

## Denoising Autoencoders



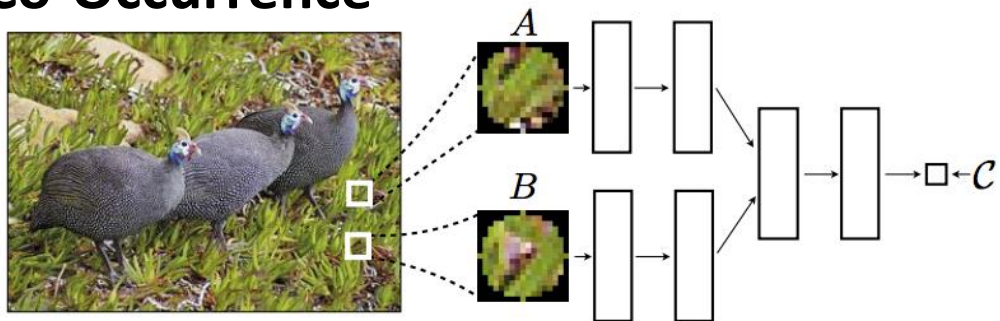
Vincent *et al.* ICML 2008.

## Audio



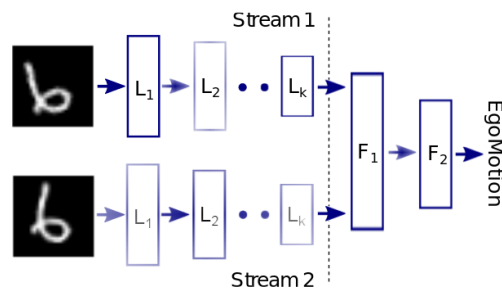
Owens *et al.* CVPR 2016, ECCV 2016

## Co-Occurrence



Isola *et al.* ICLR Workshop 2016.

## Egomotion

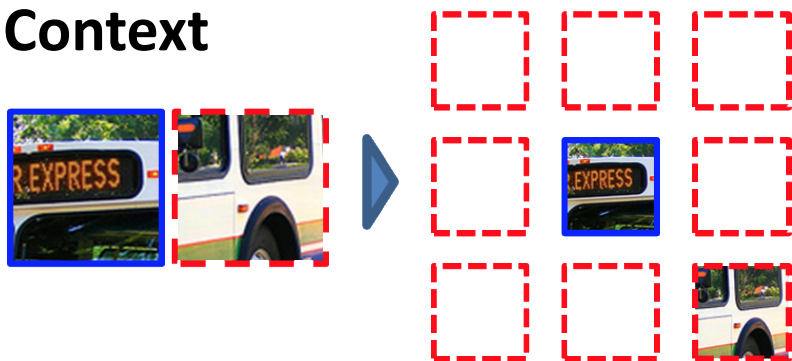


Agrawal *et al.* ICCV 2015



Jayaraman *et al.* ICCV 2015

## Context



Doersch *et al.* ICCV 2015



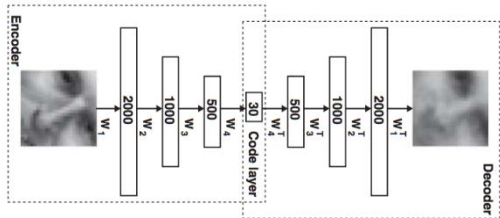
Pathak *et al.* CVPR 2016

## Video



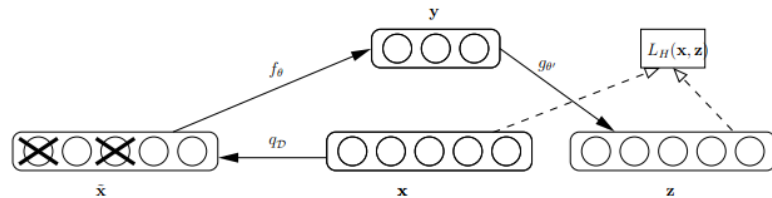
Wang *et al.* ICCV 2015

## Autoencoders



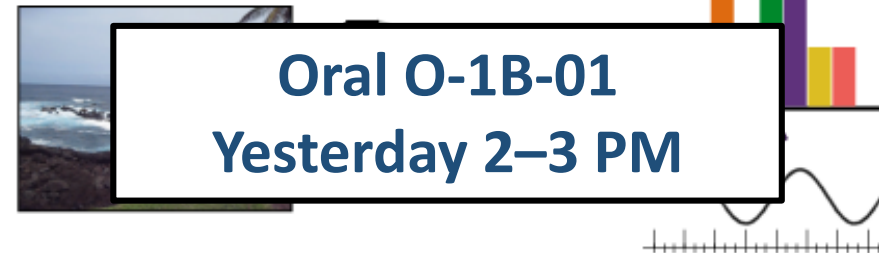
Hinton & Salakhutdinov.  
Science 2006.

## Denoising Autoencoders



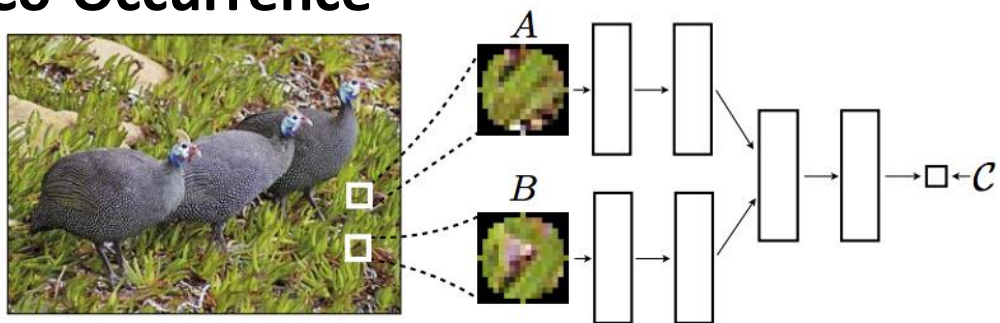
Vincent *et al.* ICML 2008.

## Audio



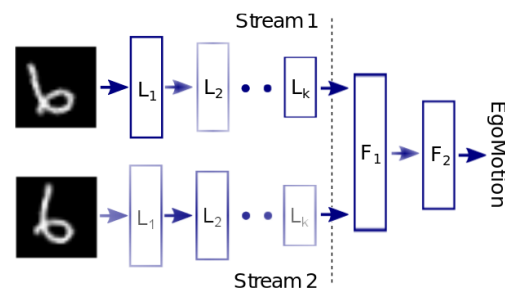
Owens *et al.* CVPR 2016, ECCV 2016

## Co-Occurrence



Isola *et al.* ICLR Workshop 2016.

## Egomotion

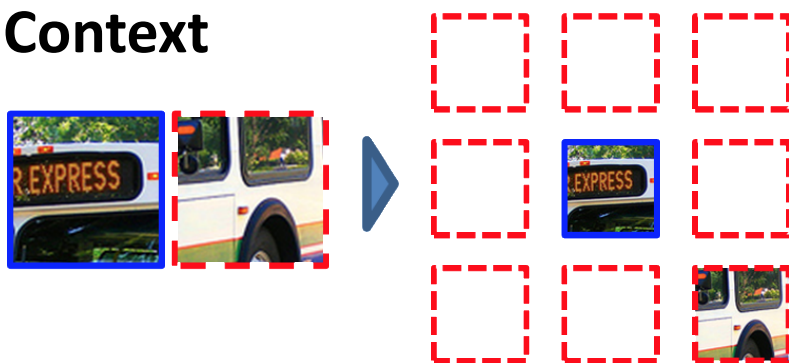


Agrawal *et al.* ICCV 2015



Jayaraman *et al.* ICCV 2015

## Context

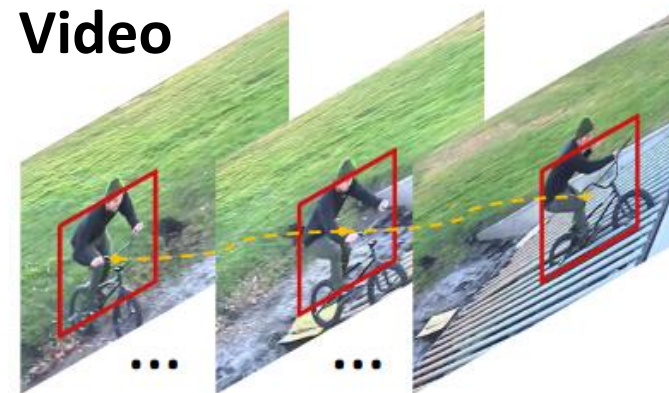


Doersch *et al.* ICCV 2015



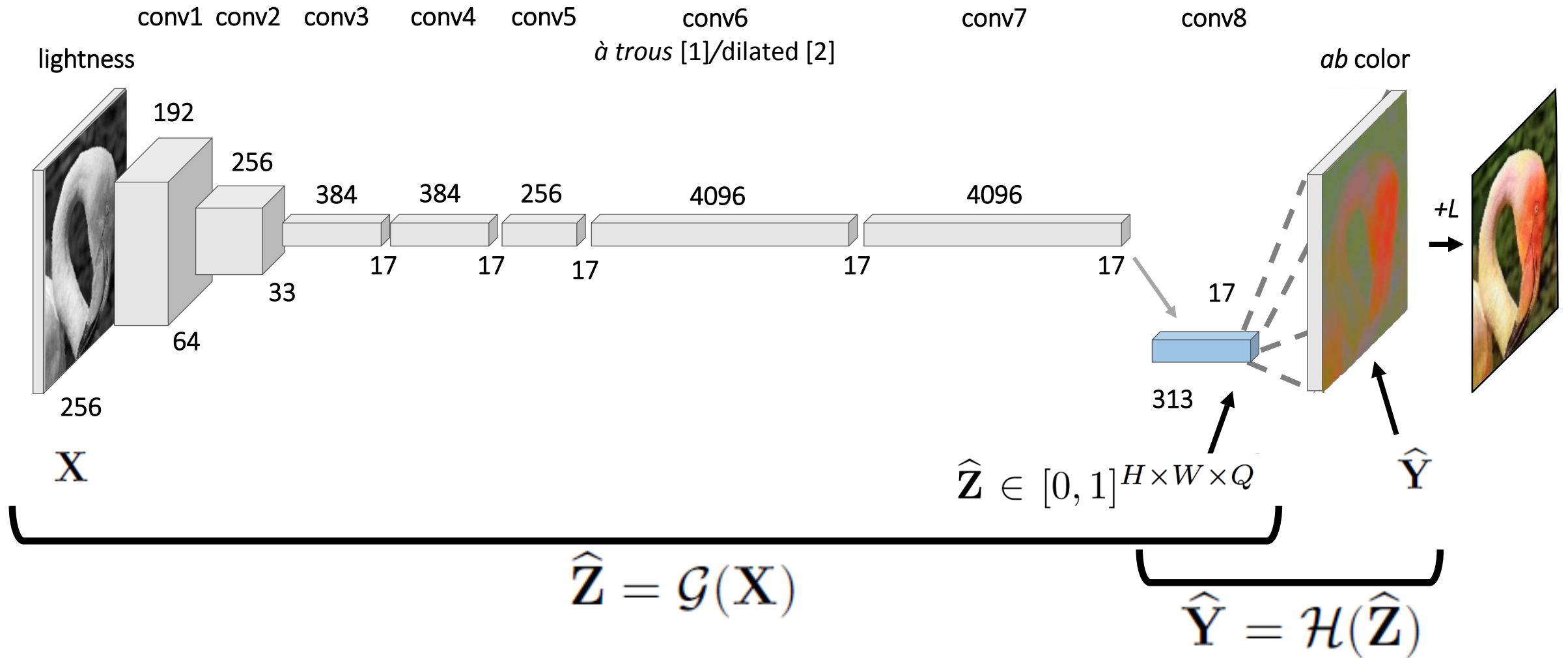
Pathak *et al.* CVPR 2016

## Video



Wang *et al.* ICCV 2015

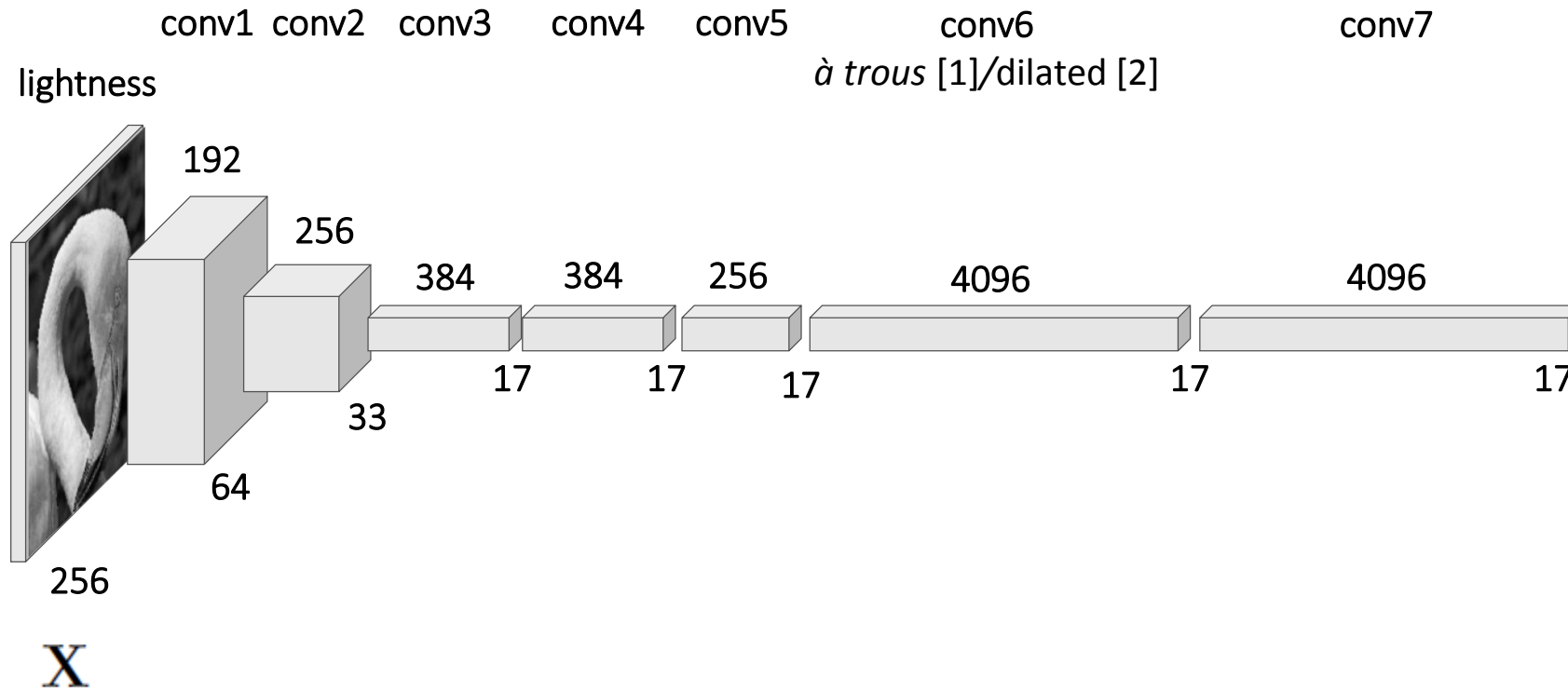
# Cross-Channel Encoder



[1] Chen *et al.* In arXiv, 2016.

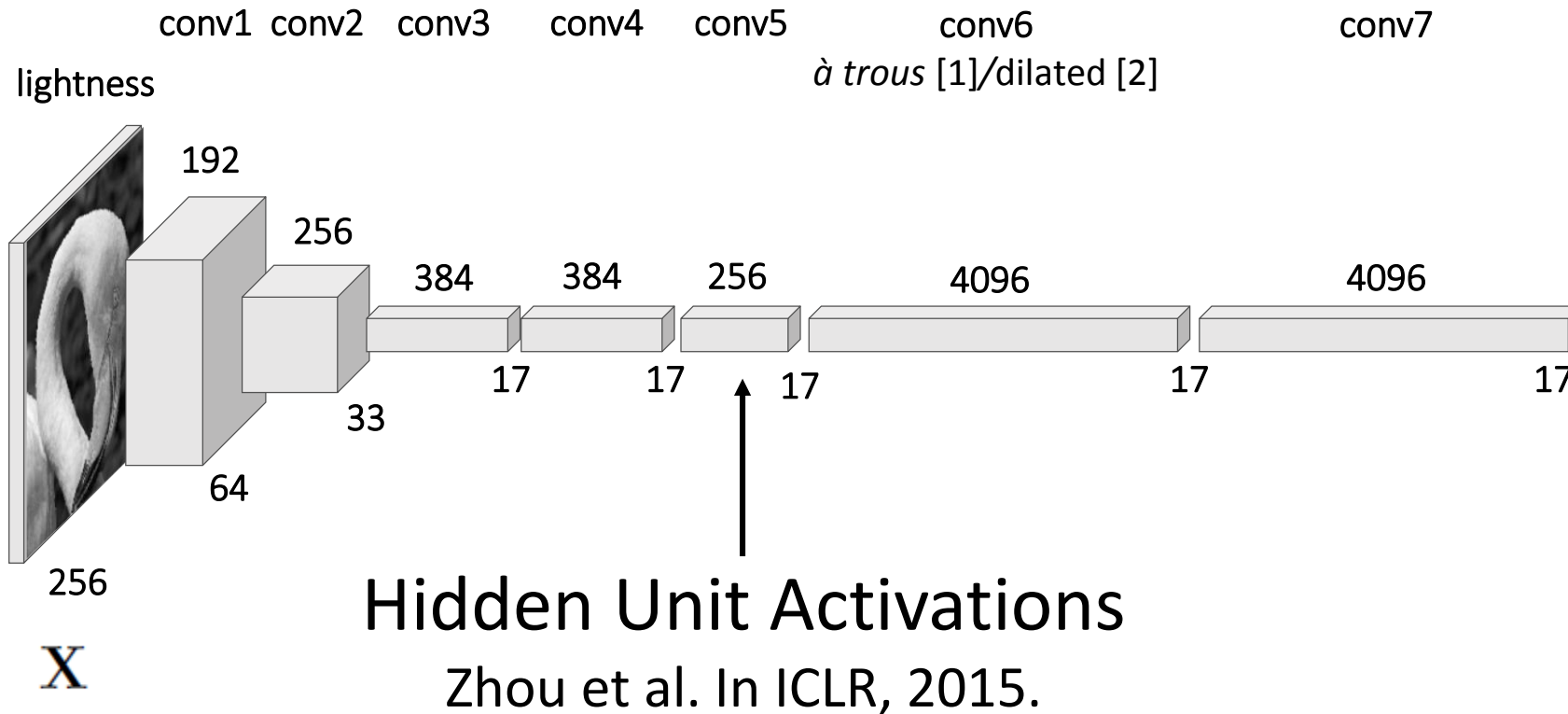
[2] Yu and Koltun. In ICLR, 2016

# Cross-Channel Encoder



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

# Cross-Channel Encoder



[1] Chen *et al.* In arXiv, 2016.

[2] Yu and Koltun. In ICLR, 2016



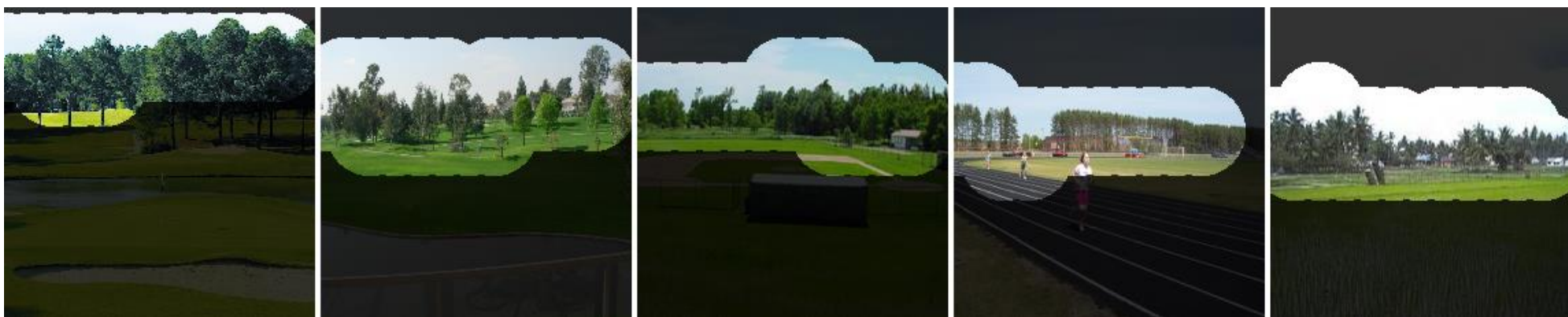
# Hidden Unit (conv5) Activations

# Hidden Unit (conv5) Activations

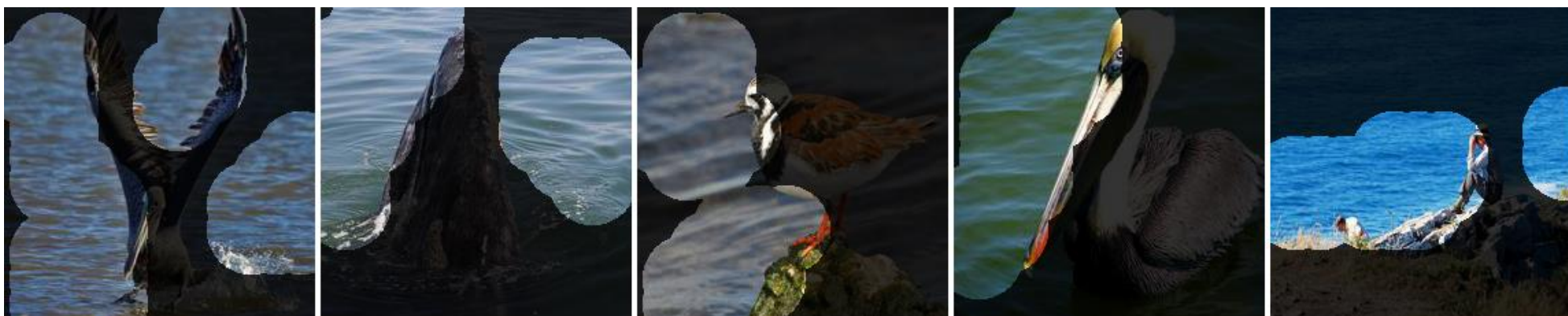
sky



trees



water



# Hidden Unit (conv5) Activations

# Hidden Unit (conv5) Activations

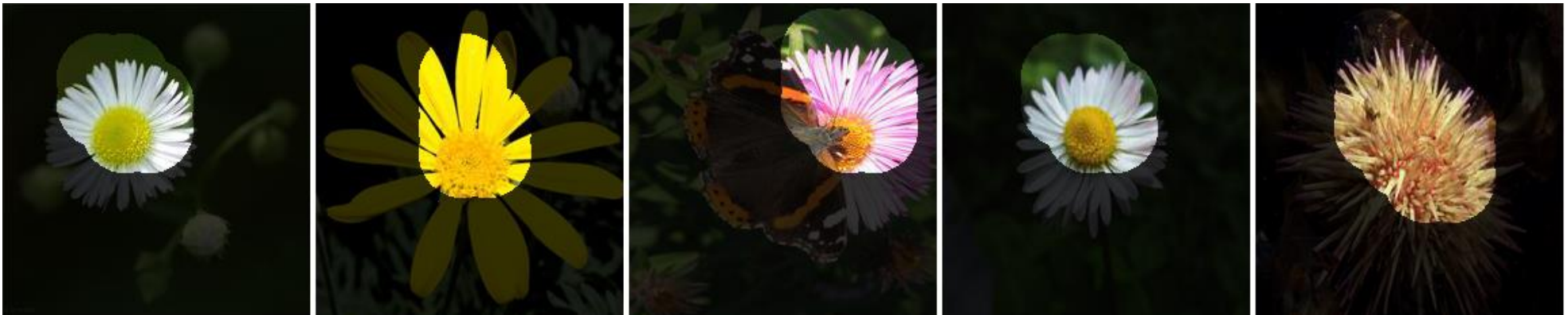
faces



dog  
faces



flowers



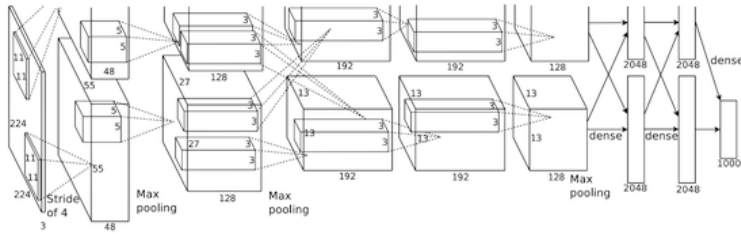
# Dataset & Task Generalization on PASCAL VOC

# Dataset & Task Generalization on PASCAL VOC

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

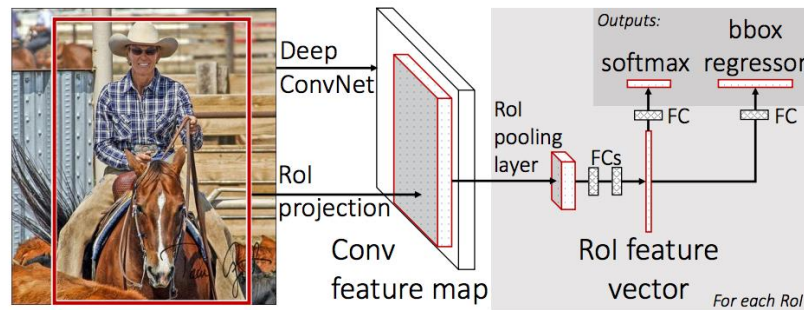
# Dataset & Task Generalization on PASCAL VOC

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



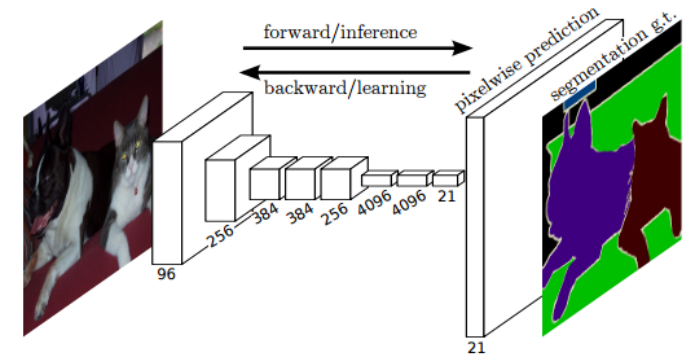
## Classification

Krähenbühl et al. In ICLR, 2016.



## Detection

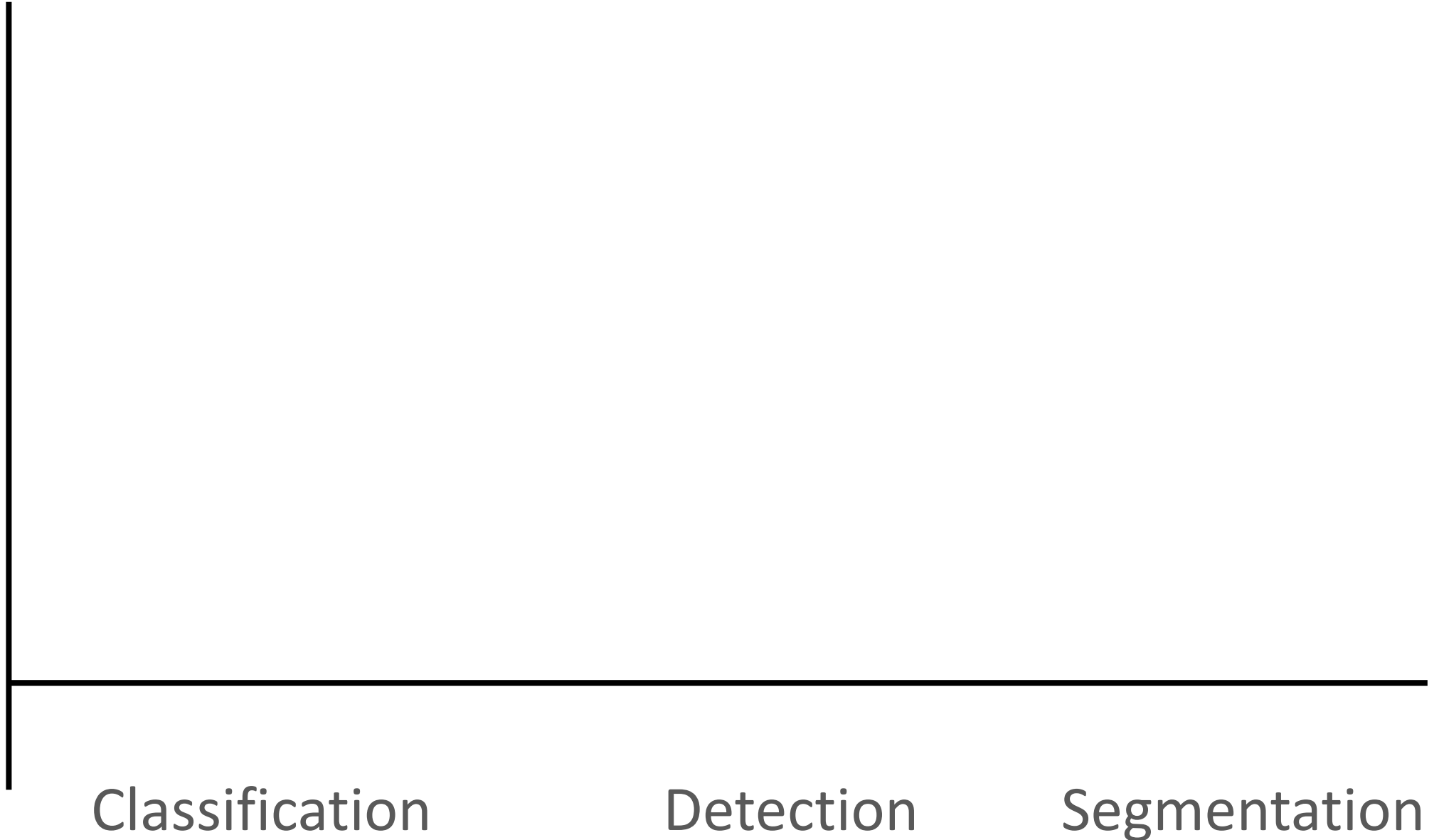
Fast R-CNN. Girshick. In ICCV, 2015.



## Segmentation

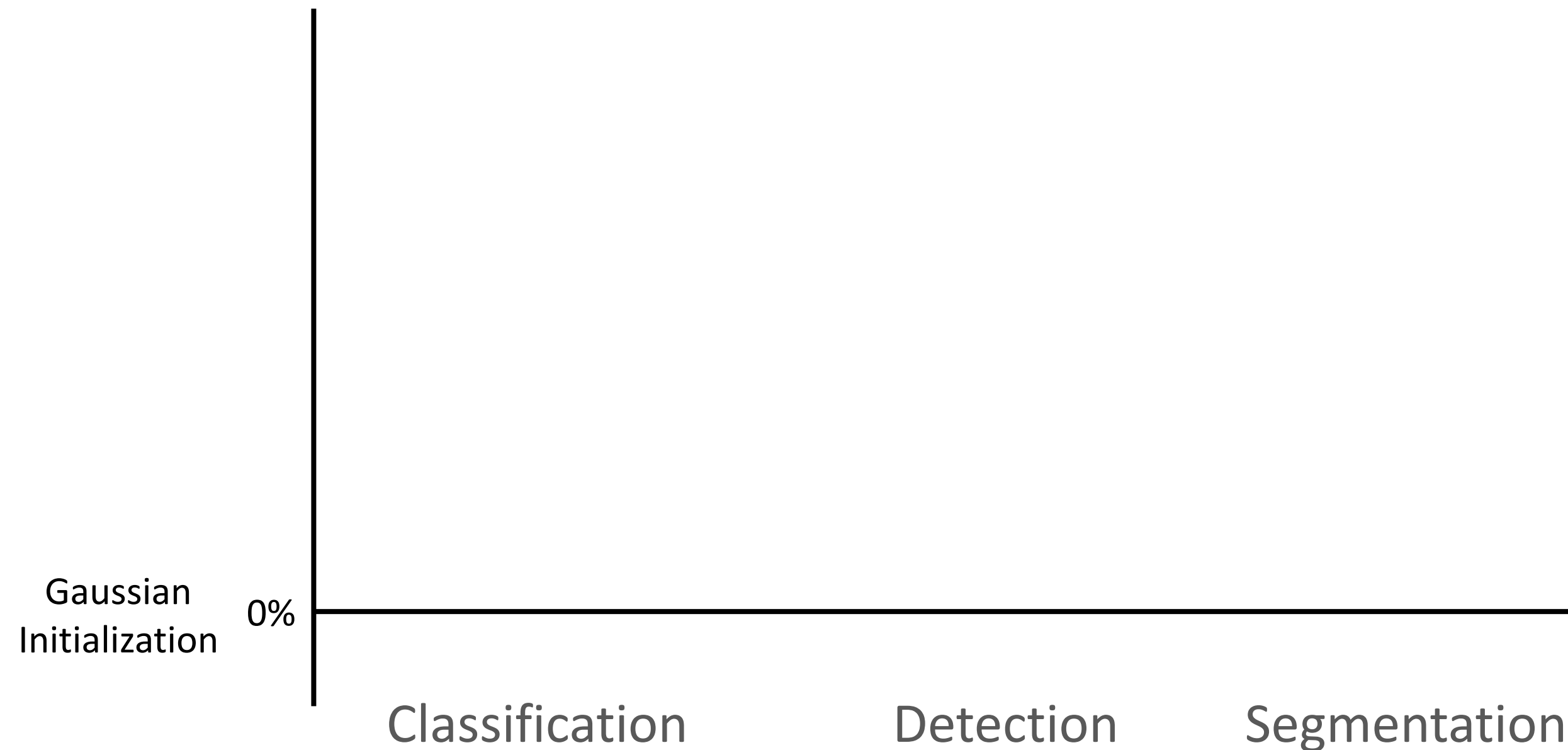
FCNs. Long et al. In CVPR, 2015.

# Dataset & Task Generalization on PASCAL VOC





# Dataset & Task Generalization on PASCAL VOC



# Dataset & Task Generalization on PASCAL VOC

ImageNet  
Labels

100%

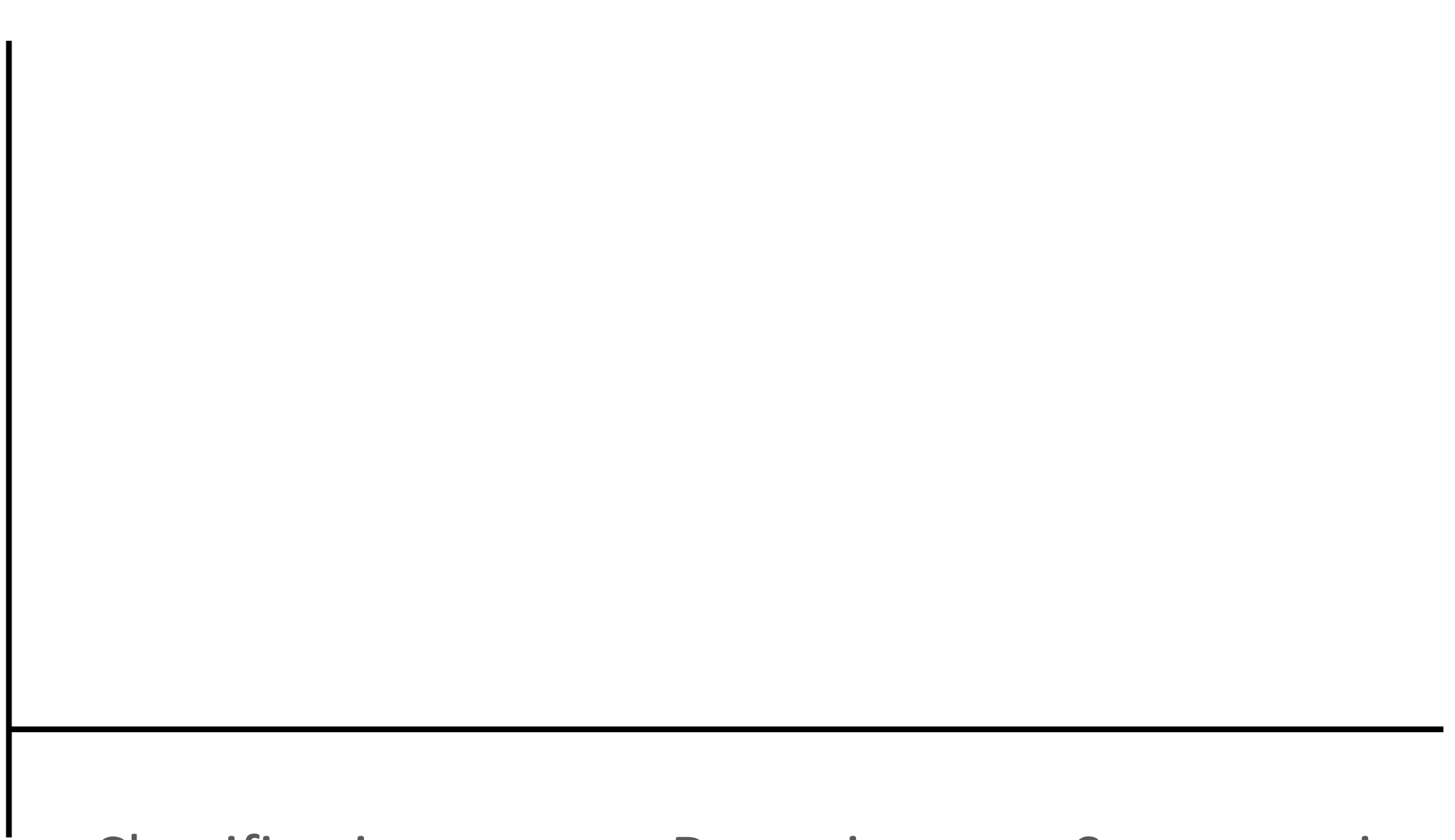
Gaussian  
Initialization

0%

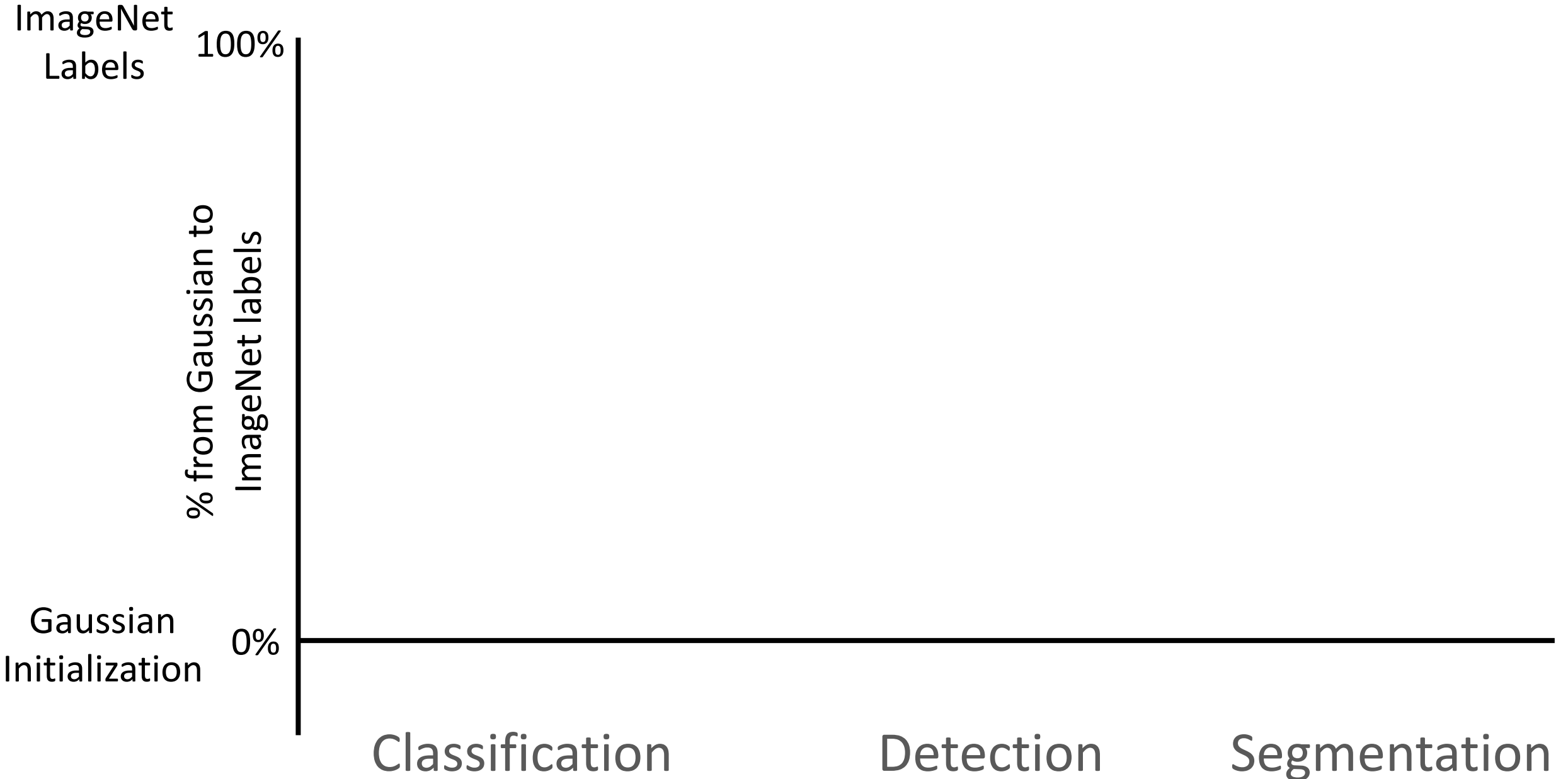
Classification

Detection

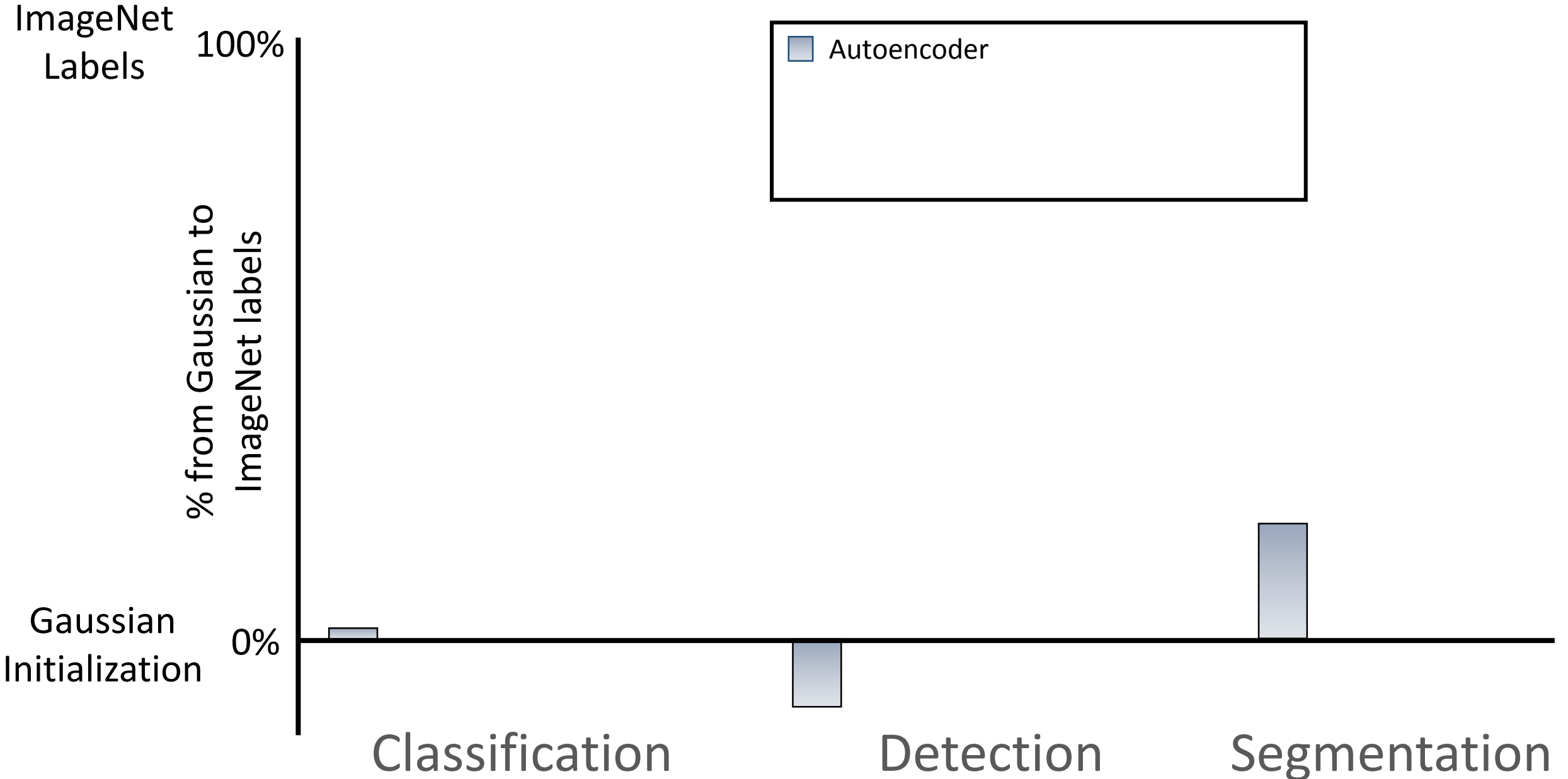
Segmentation



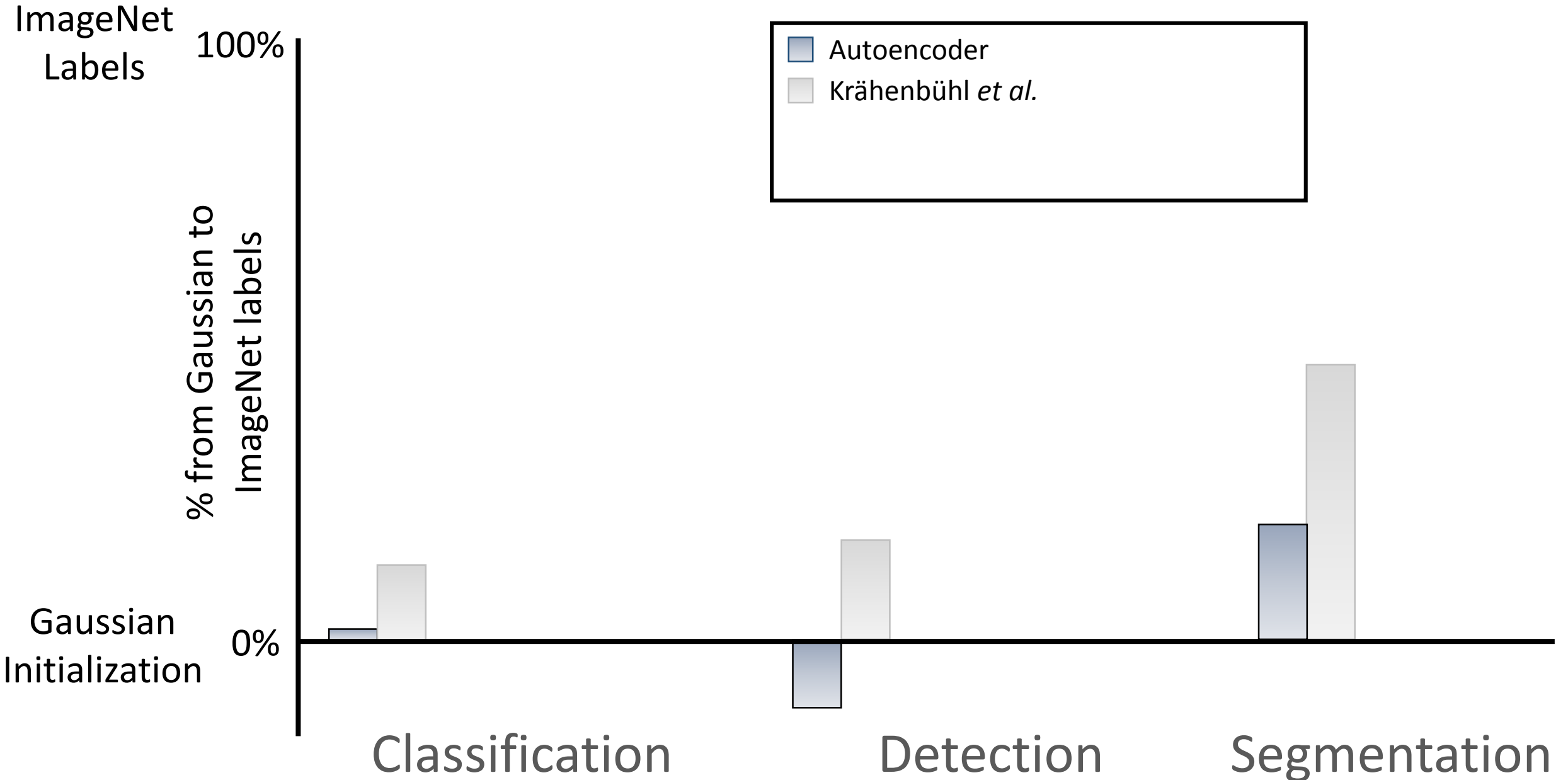
# Dataset & Task Generalization on PASCAL VOC



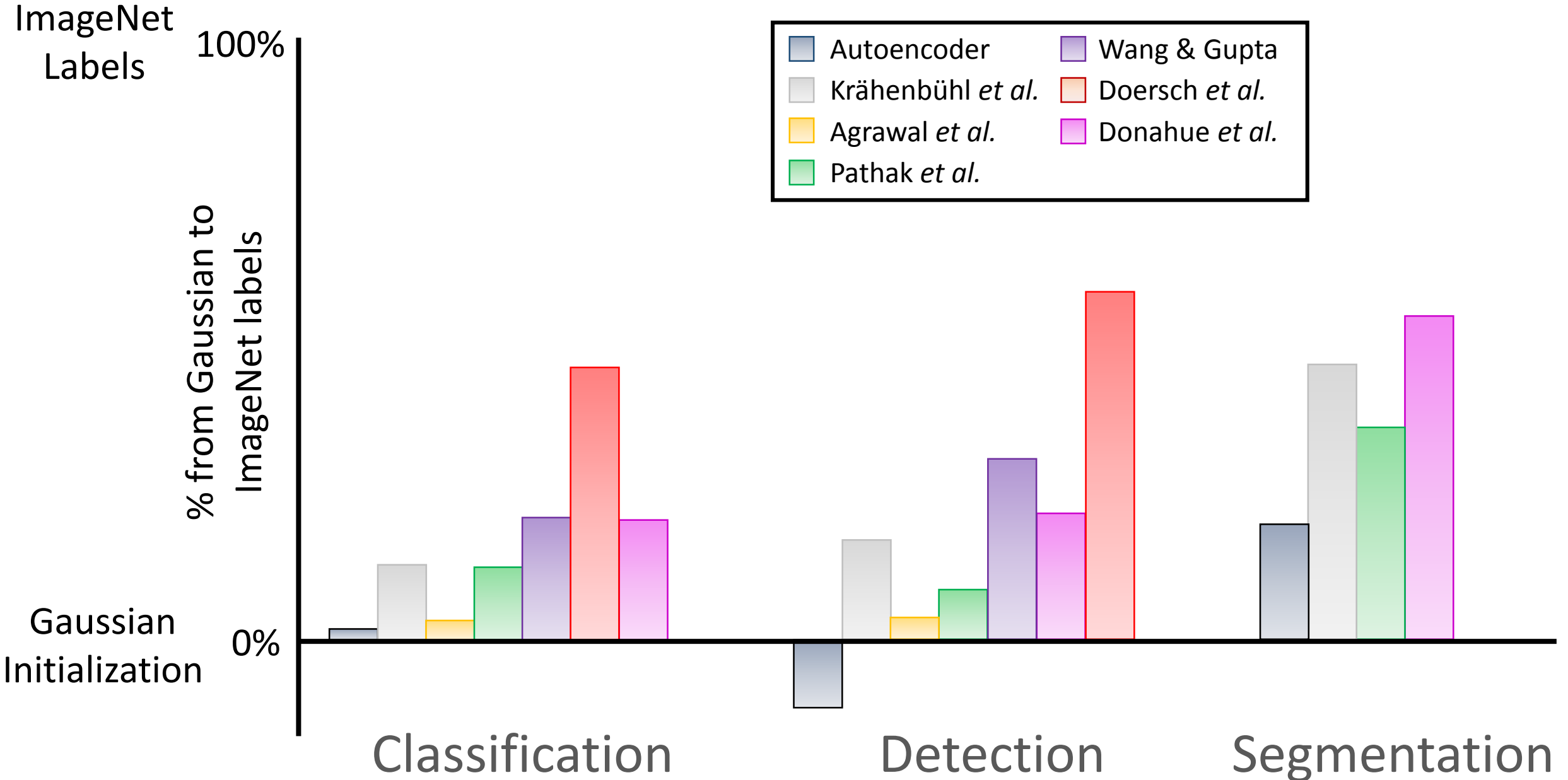
# Dataset & Task Generalization on PASCAL VOC



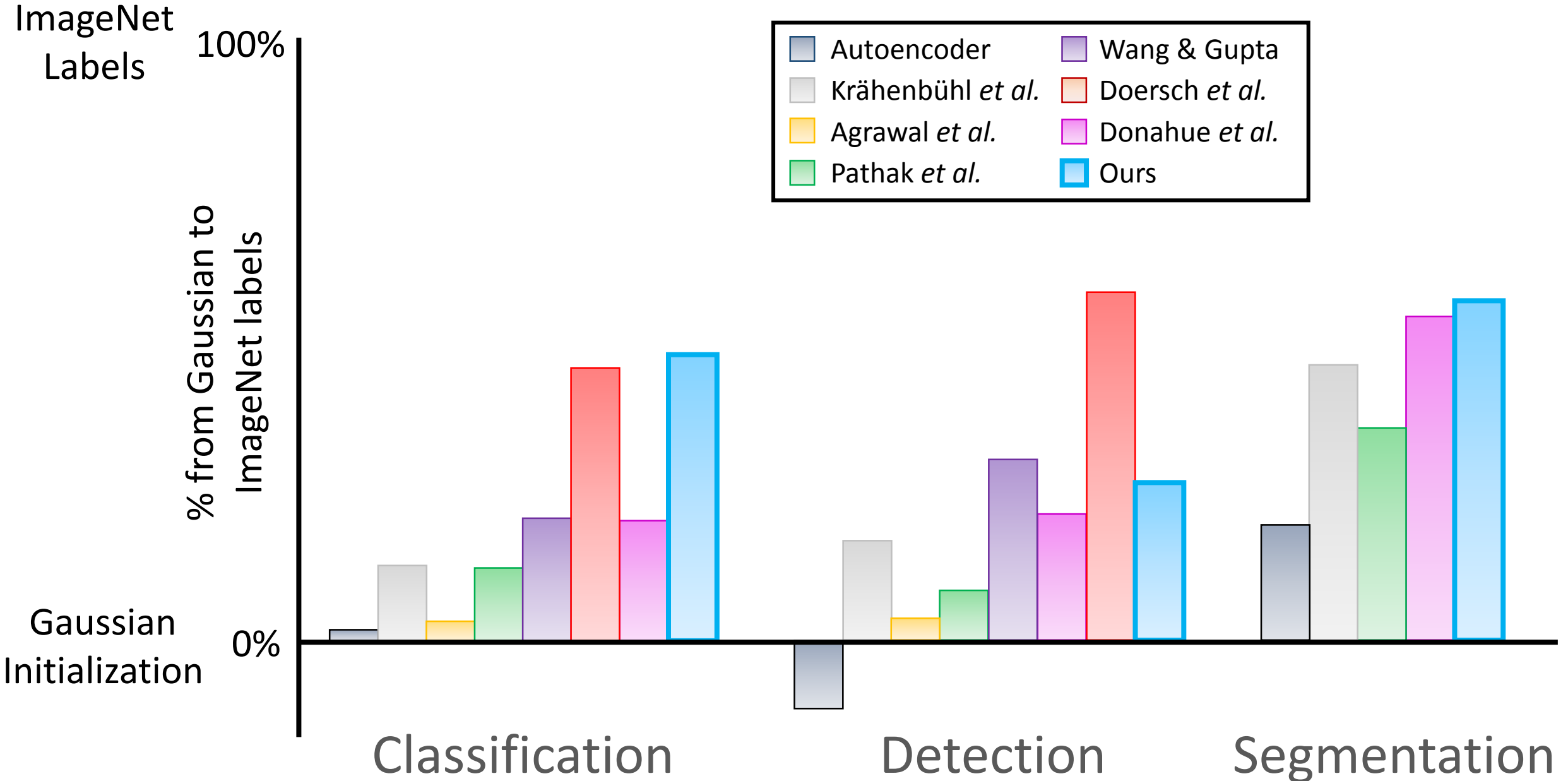
# Dataset & Task Generalization on PASCAL VOC



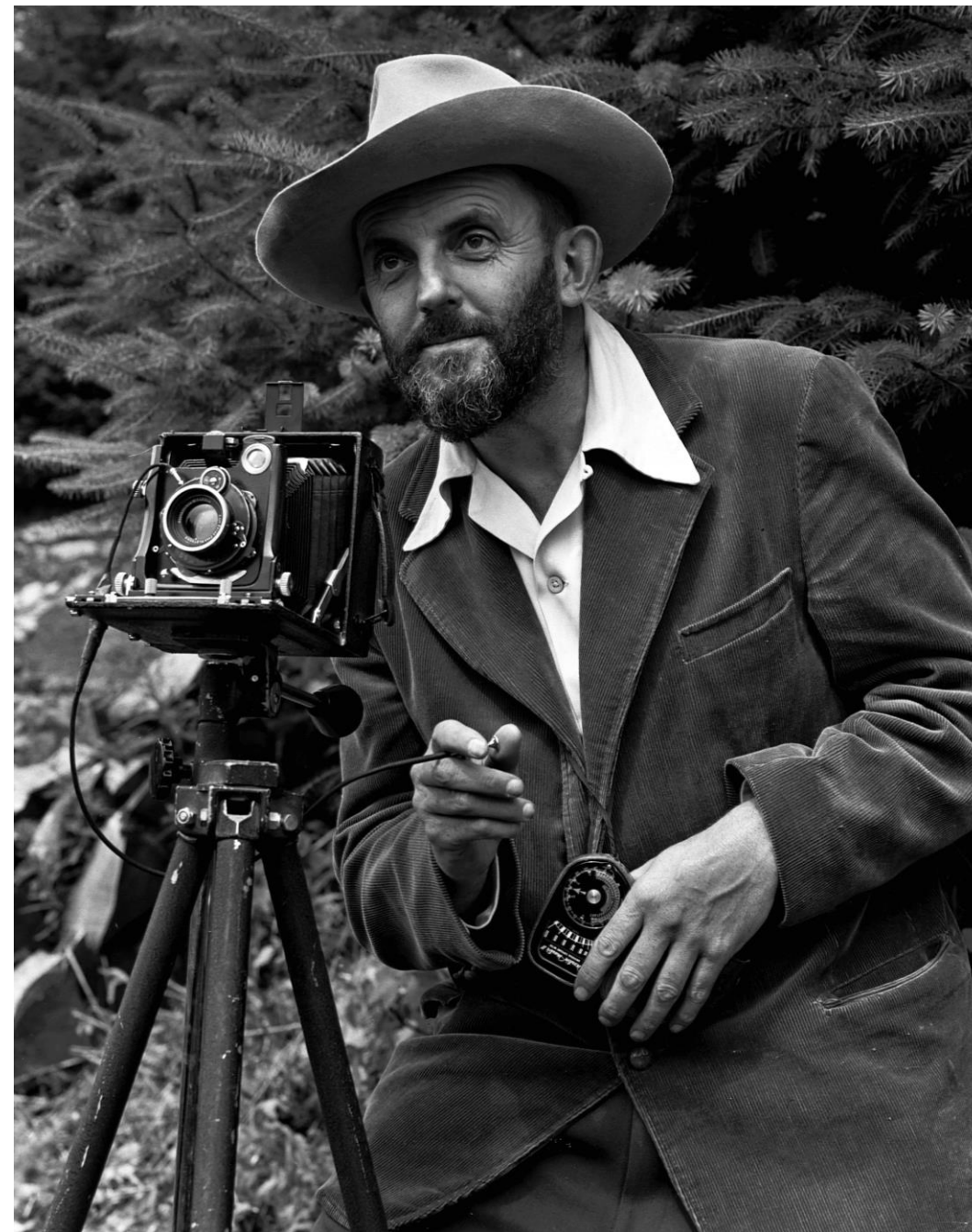
# Dataset & Task Generalization on PASCAL VOC



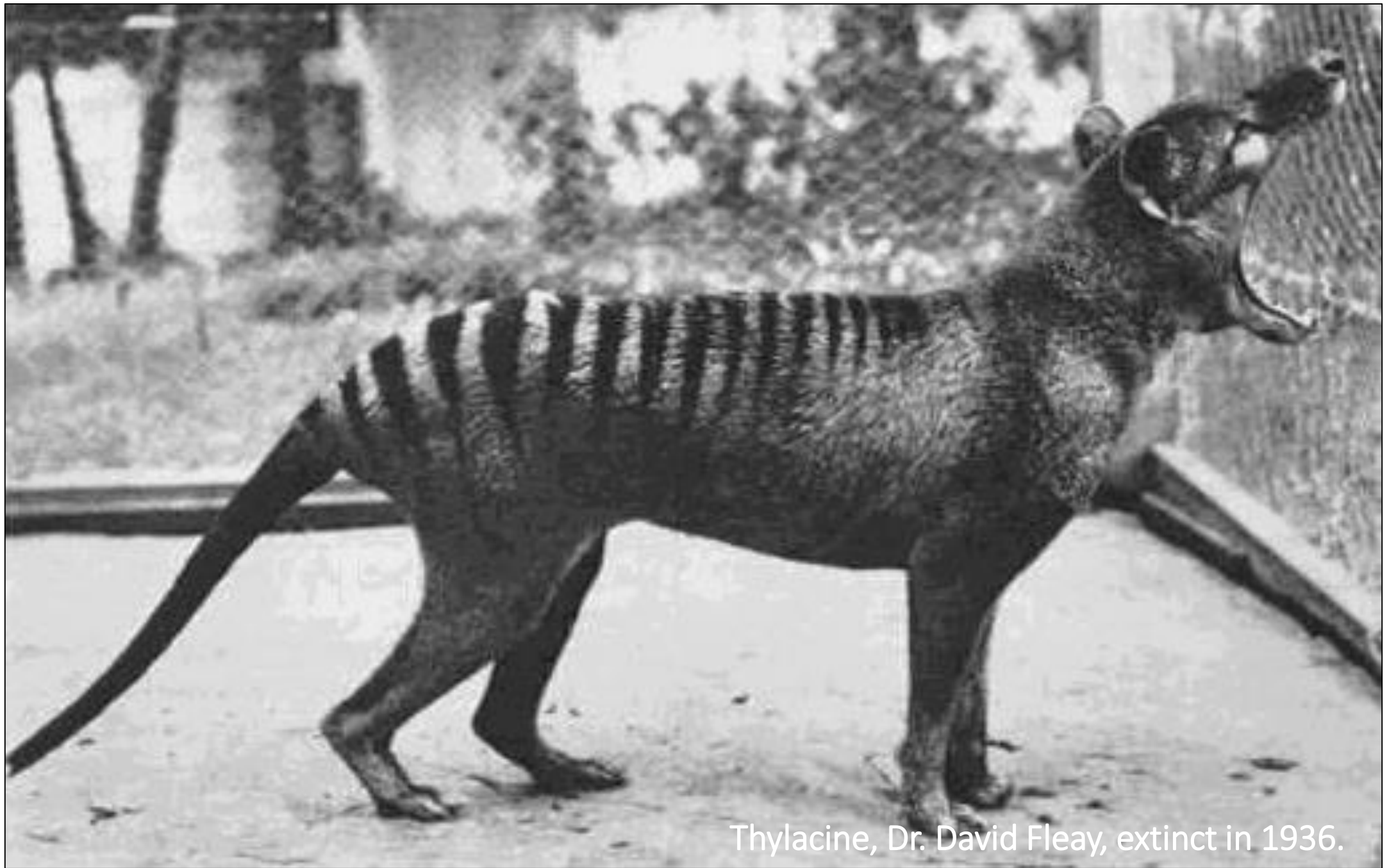
# Dataset & Task Generalization on PASCAL VOC



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







Thylacine, Dr. David Fleay, extinct in 1936.



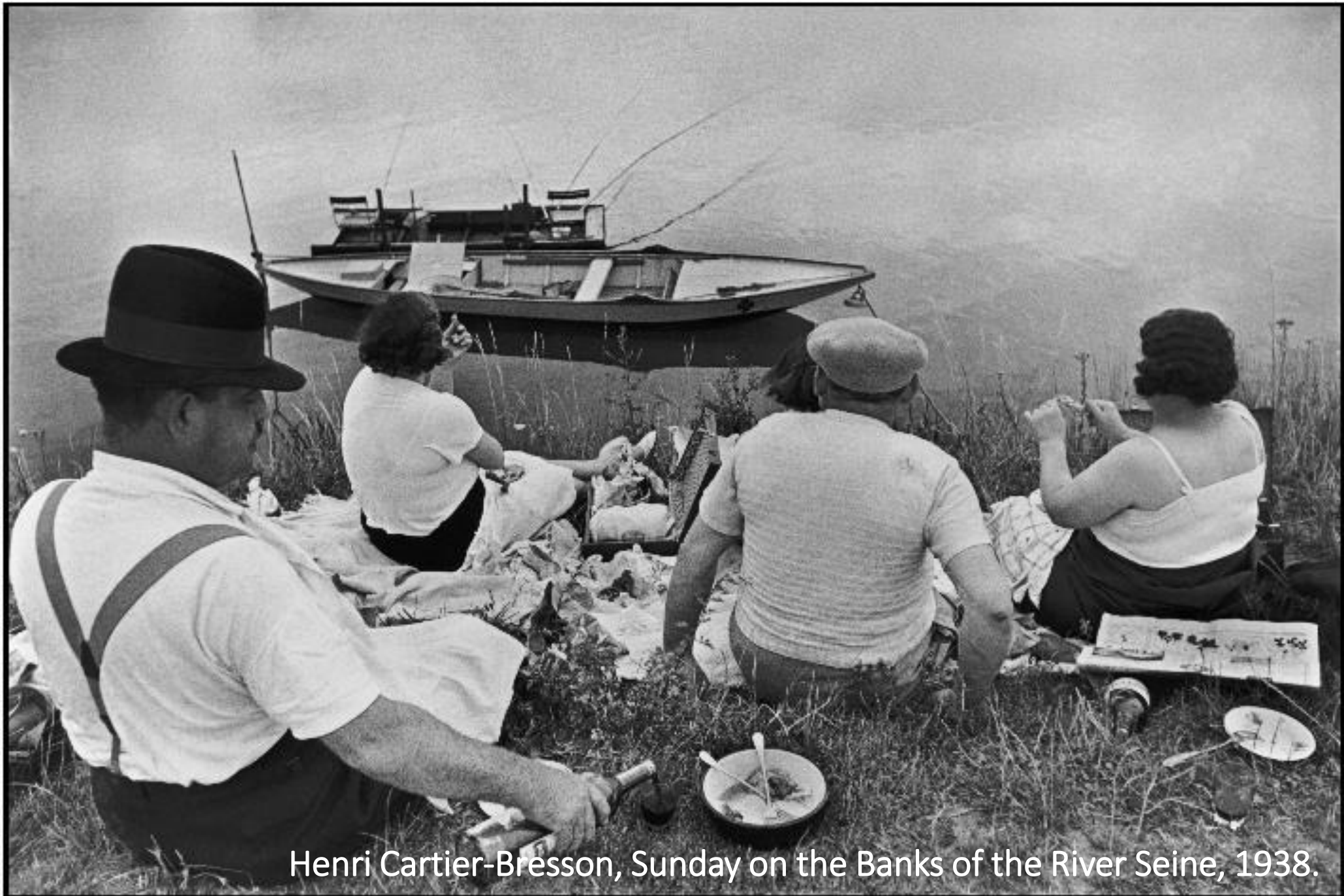
Thylacine, Dr. David Fleay, extinct in 1936.



Amateur Family Photo, 1956.



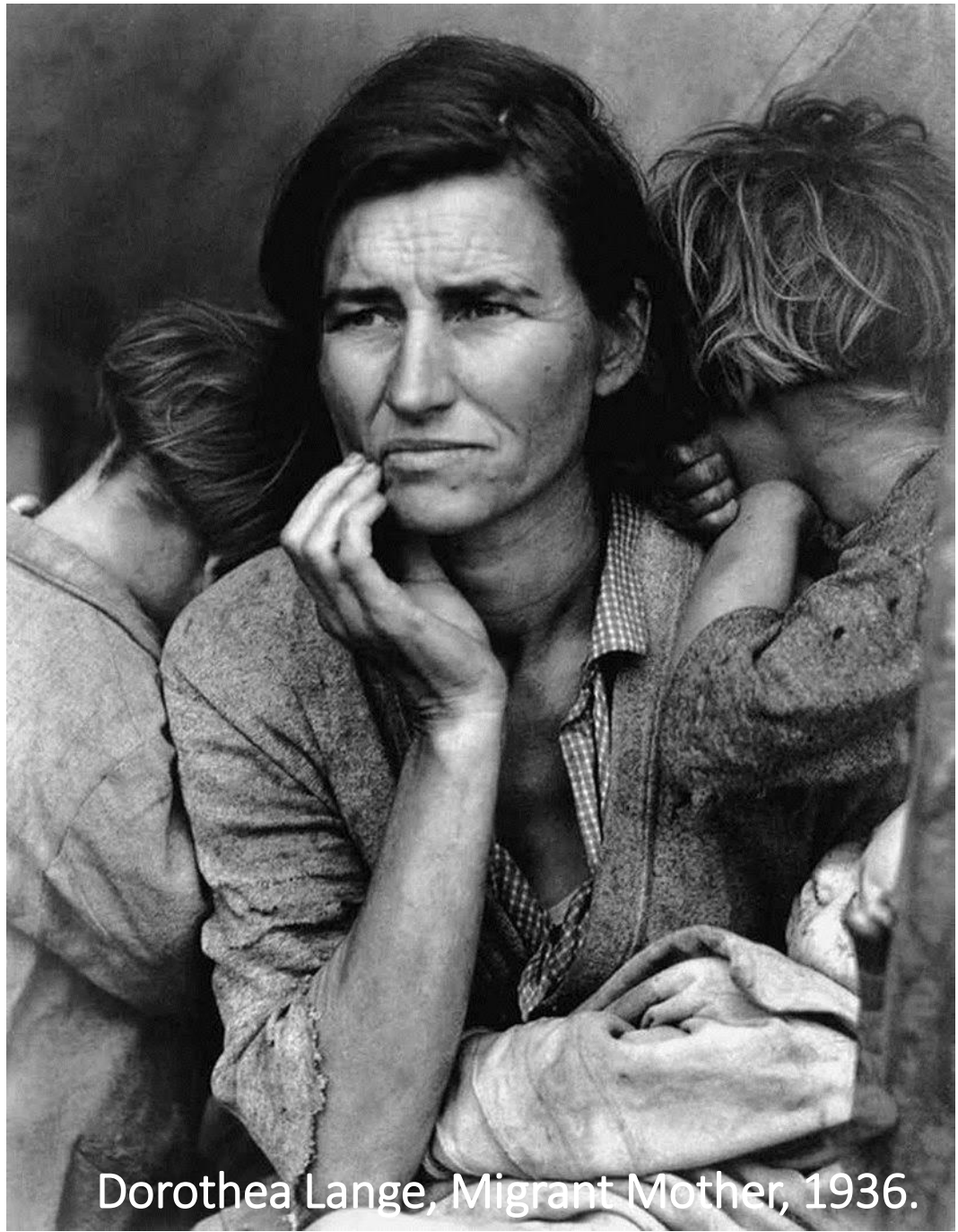
Amateur Family Photo, 1956.



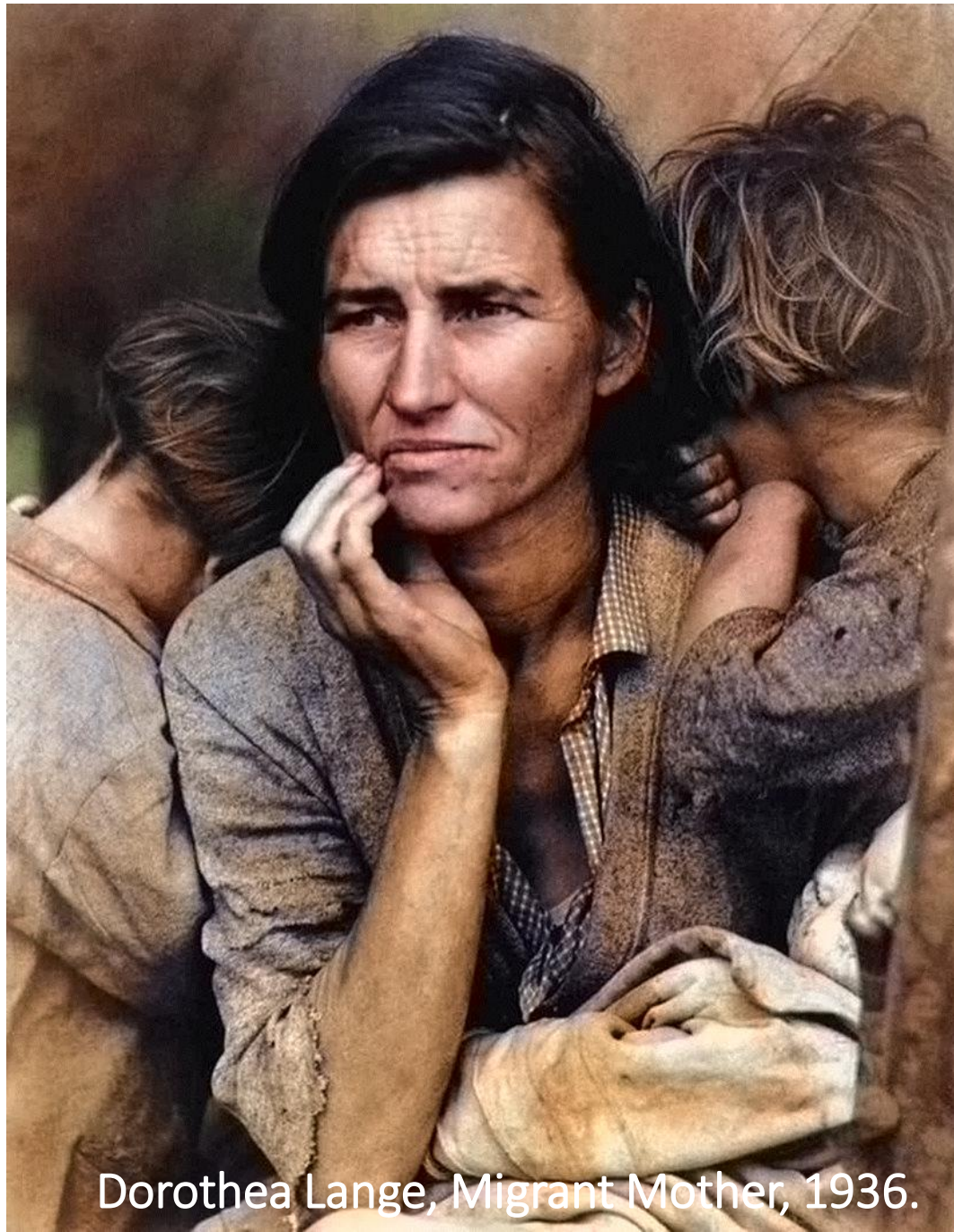
Henri Cartier-Bresson, Sunday on the Banks of the River Seine, 1938.



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



Dorothea Lange, Migrant Mother, 1936.



Dorothea Lange, Migrant Mother, 1936.



# Additional Information

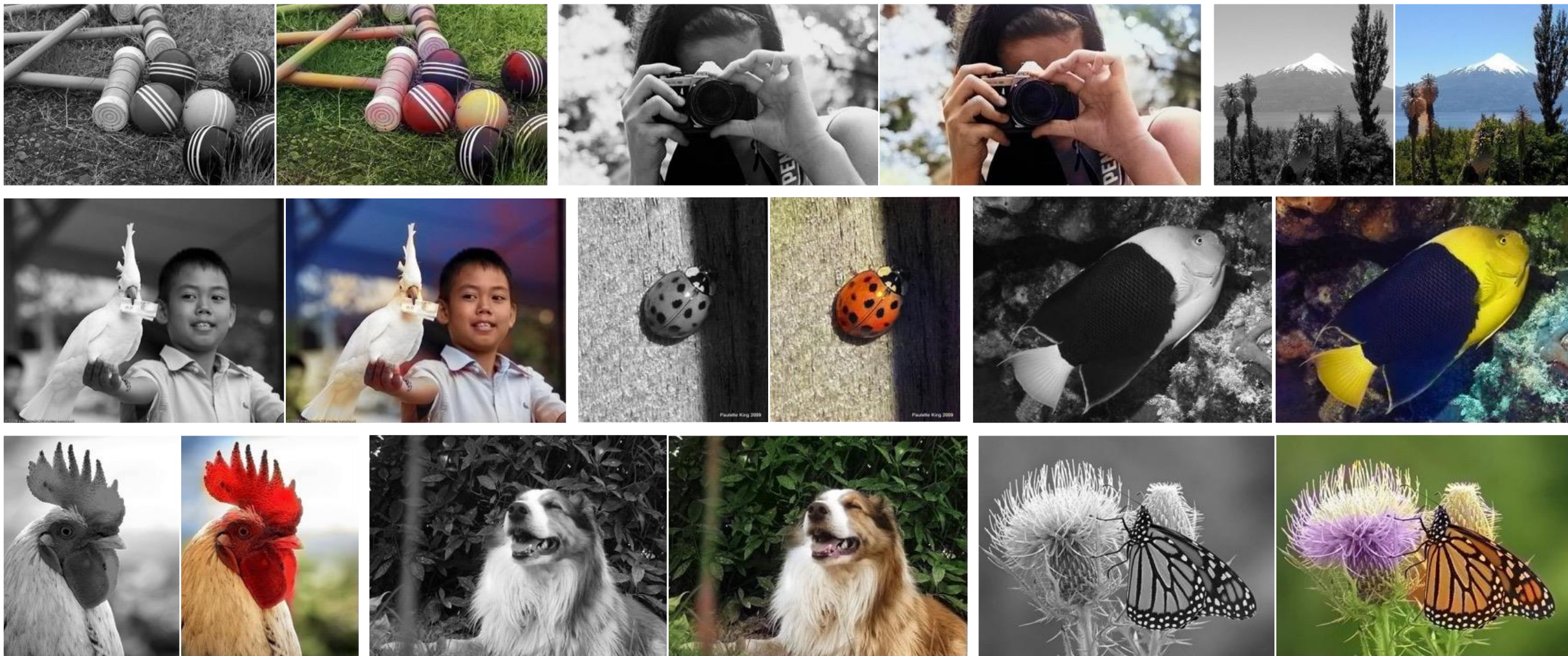
- Demo
  - <http://demos.algorithmia.com/colorize-photos/>
- Reddit ColorizeBot
  - Type “colorizebot” under any image post
- Code
  - <https://github.com/richzhang/colorization>
- Website – full paper, user examples, visualizations
  - <http://richzhang.github.io/colorization>

# Lukas Graham – 7 Years

Submitted by Ron Zohar

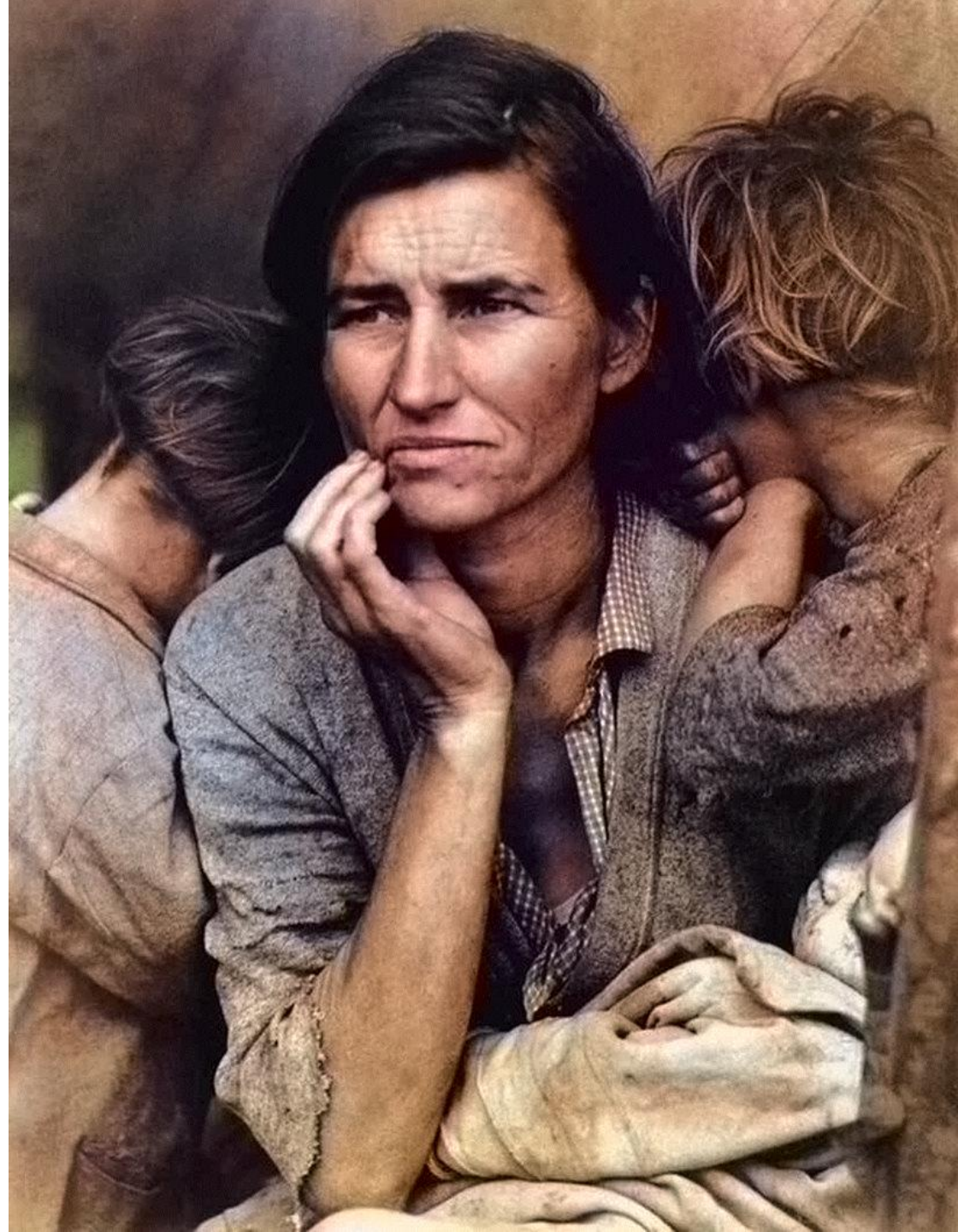


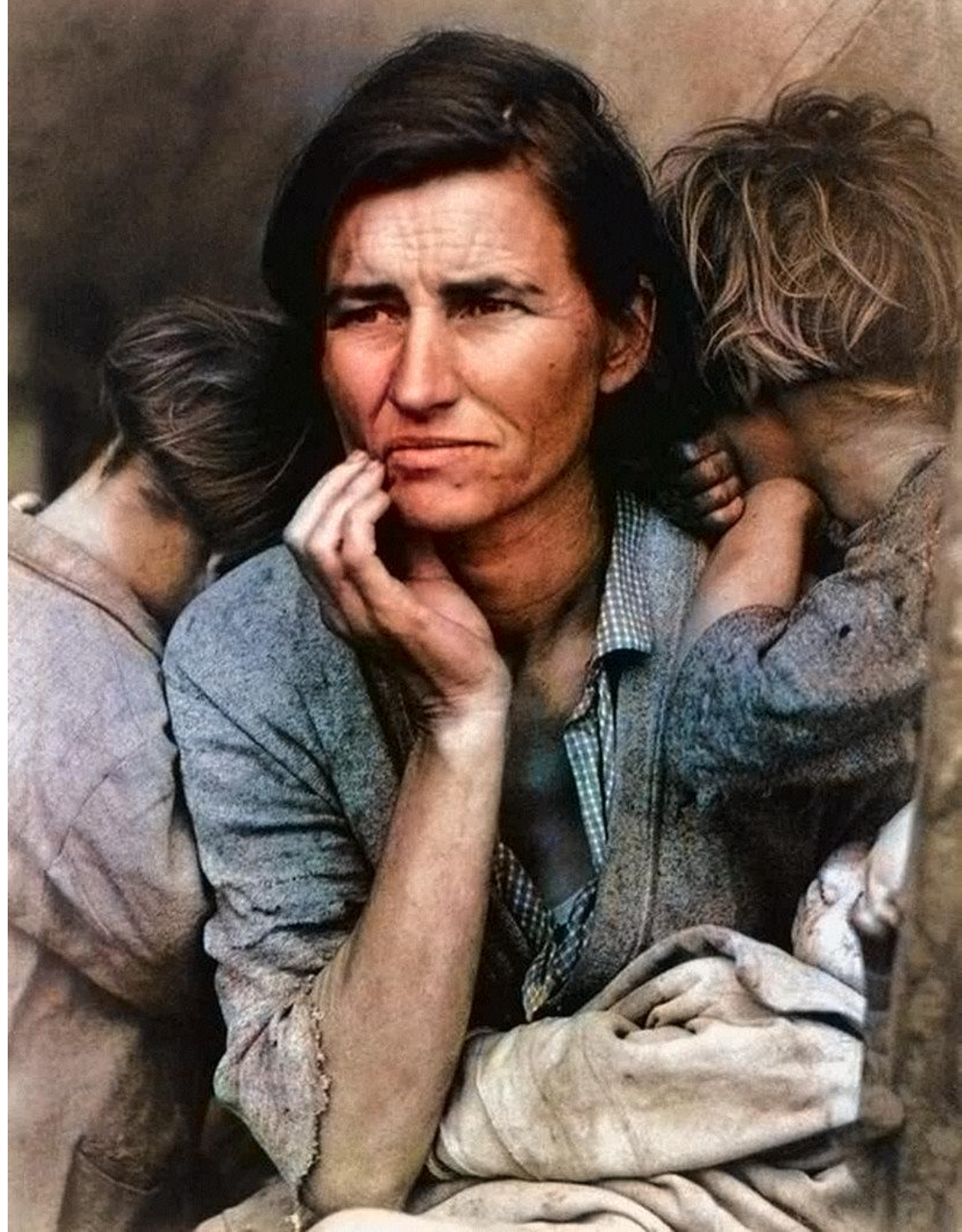
**For the full paper, code, and live demo:**  
[richzhang.github.io/colorization](https://richzhang.github.io/colorization)

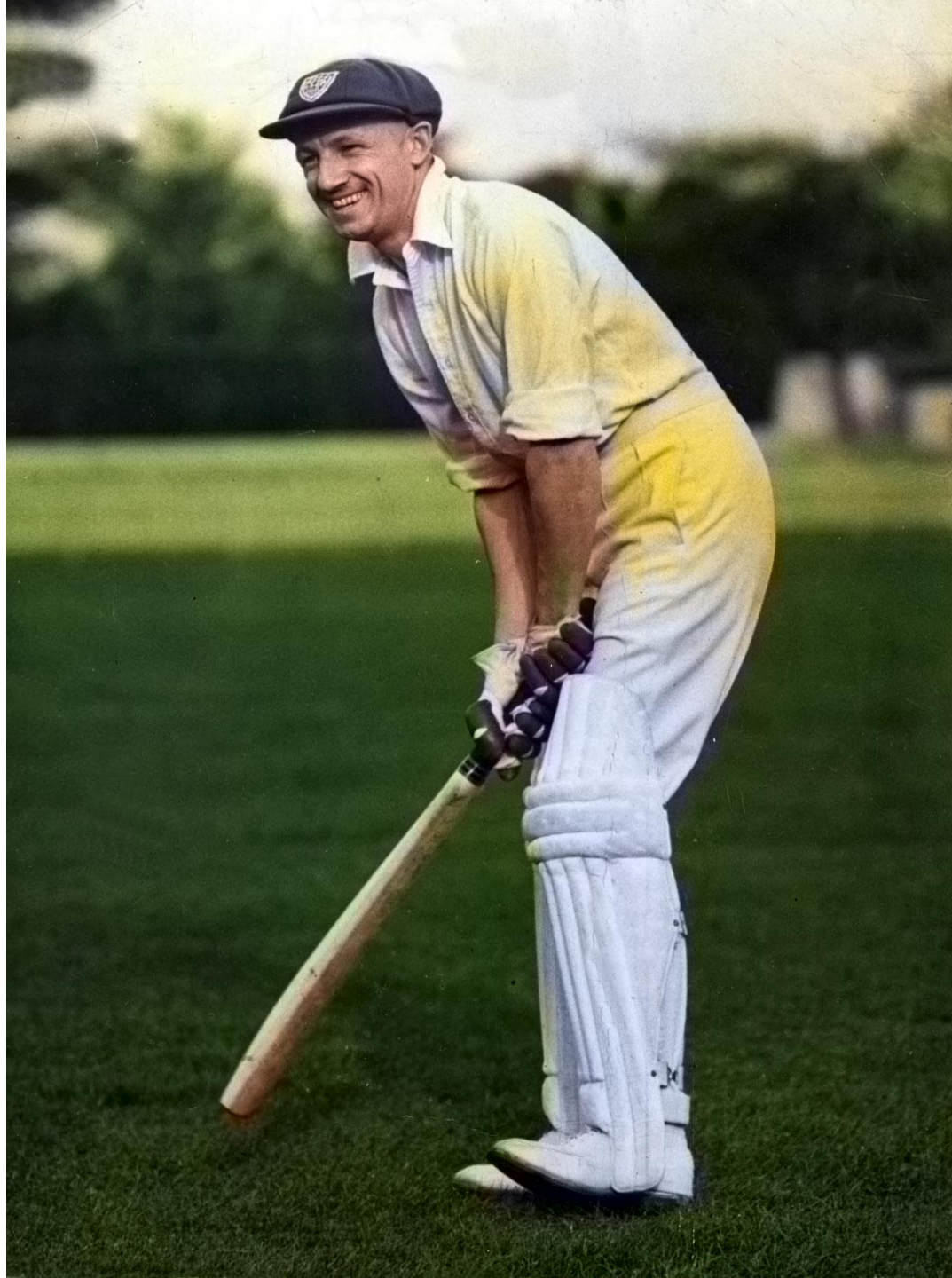


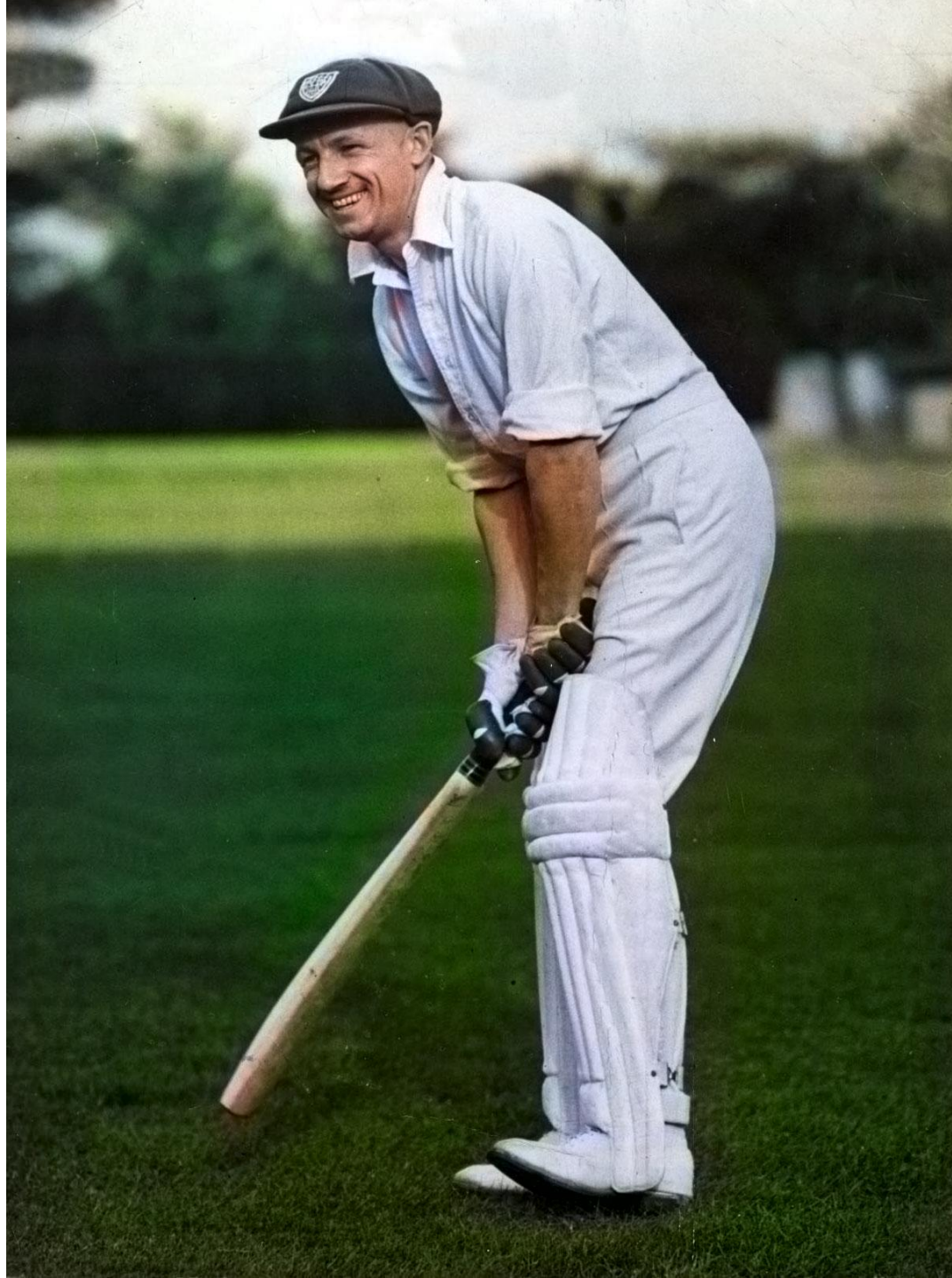
For the full paper, additional examples and our model:  
[richzhang.github.io/colorization](https://richzhang.github.io/colorization)

Backup









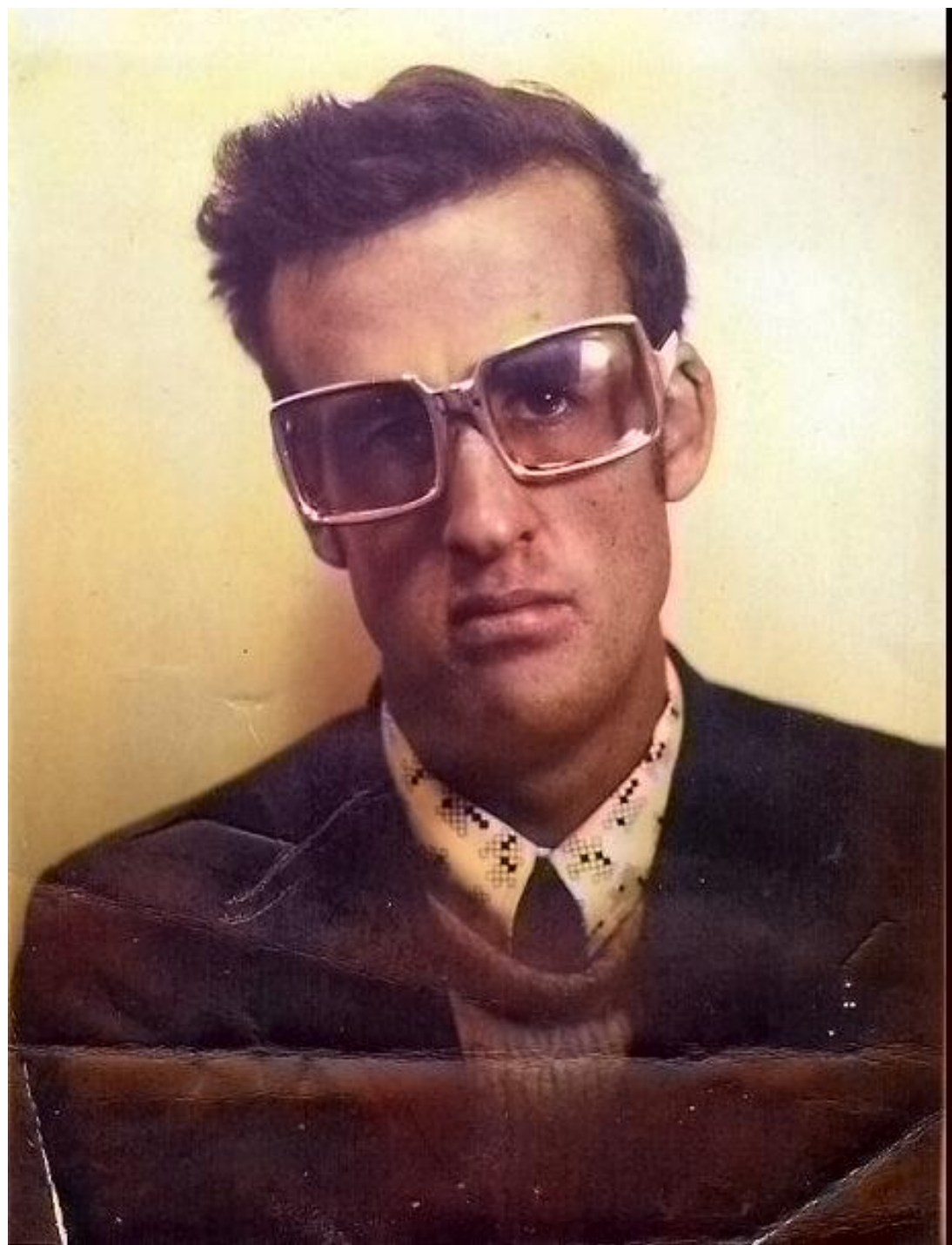


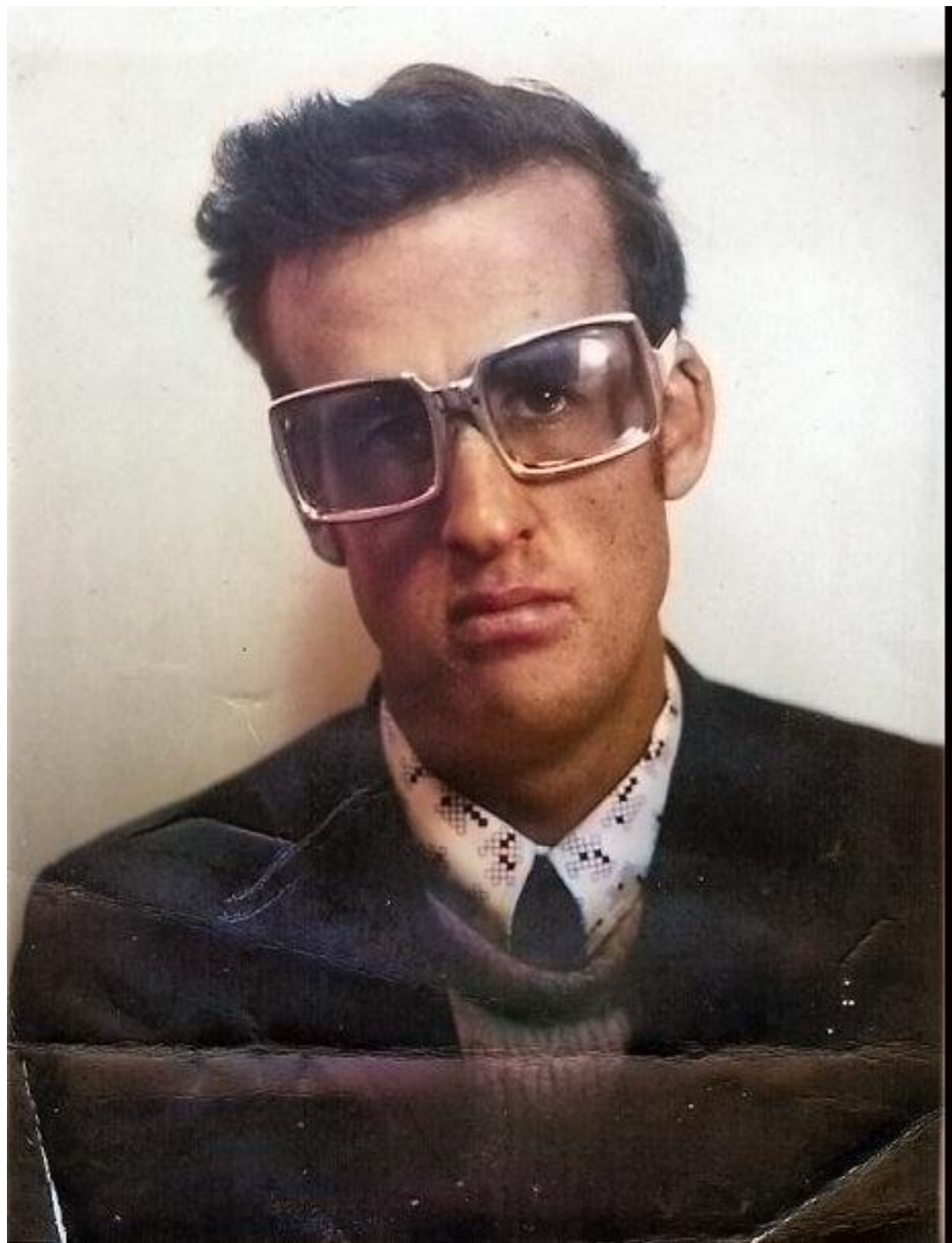










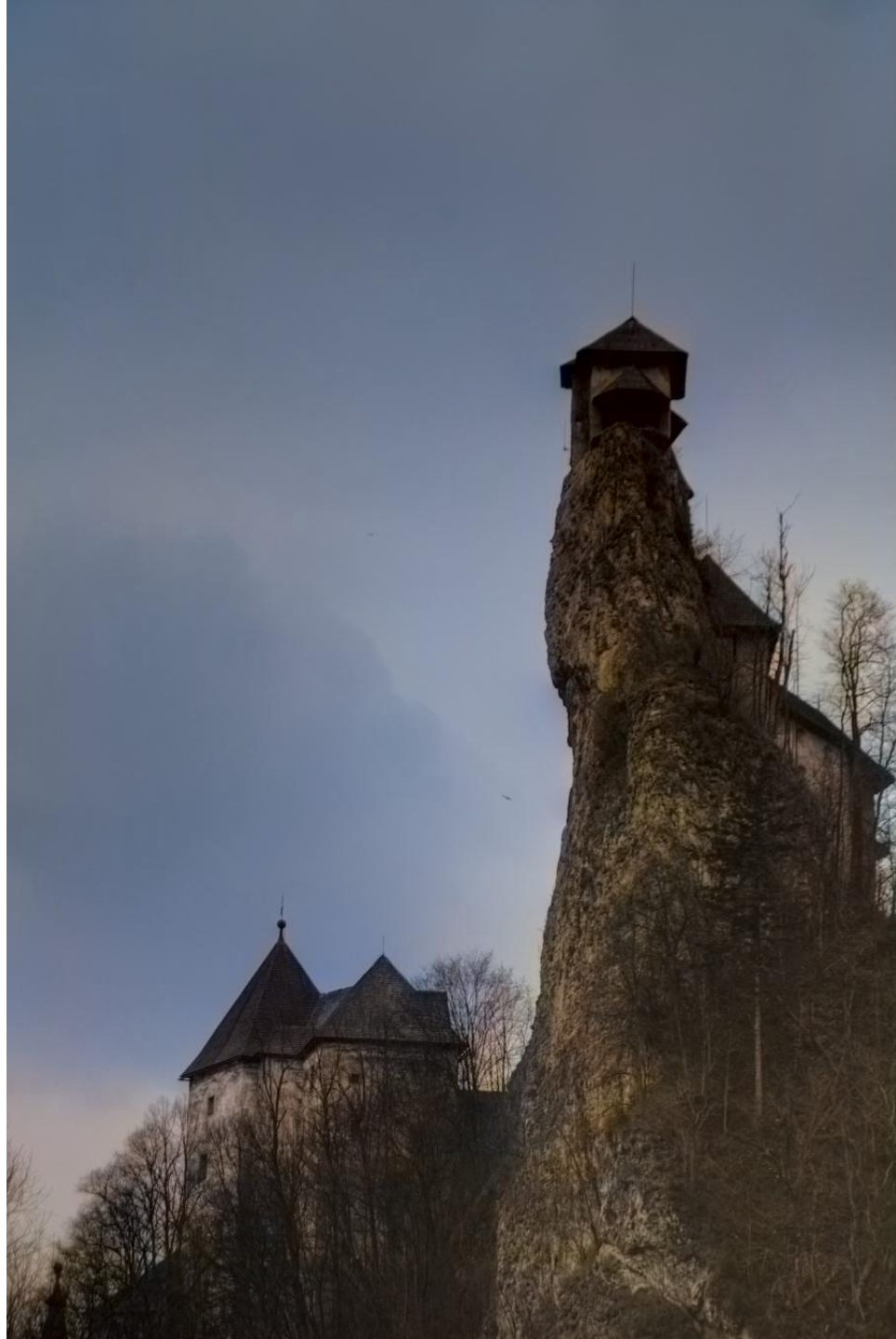






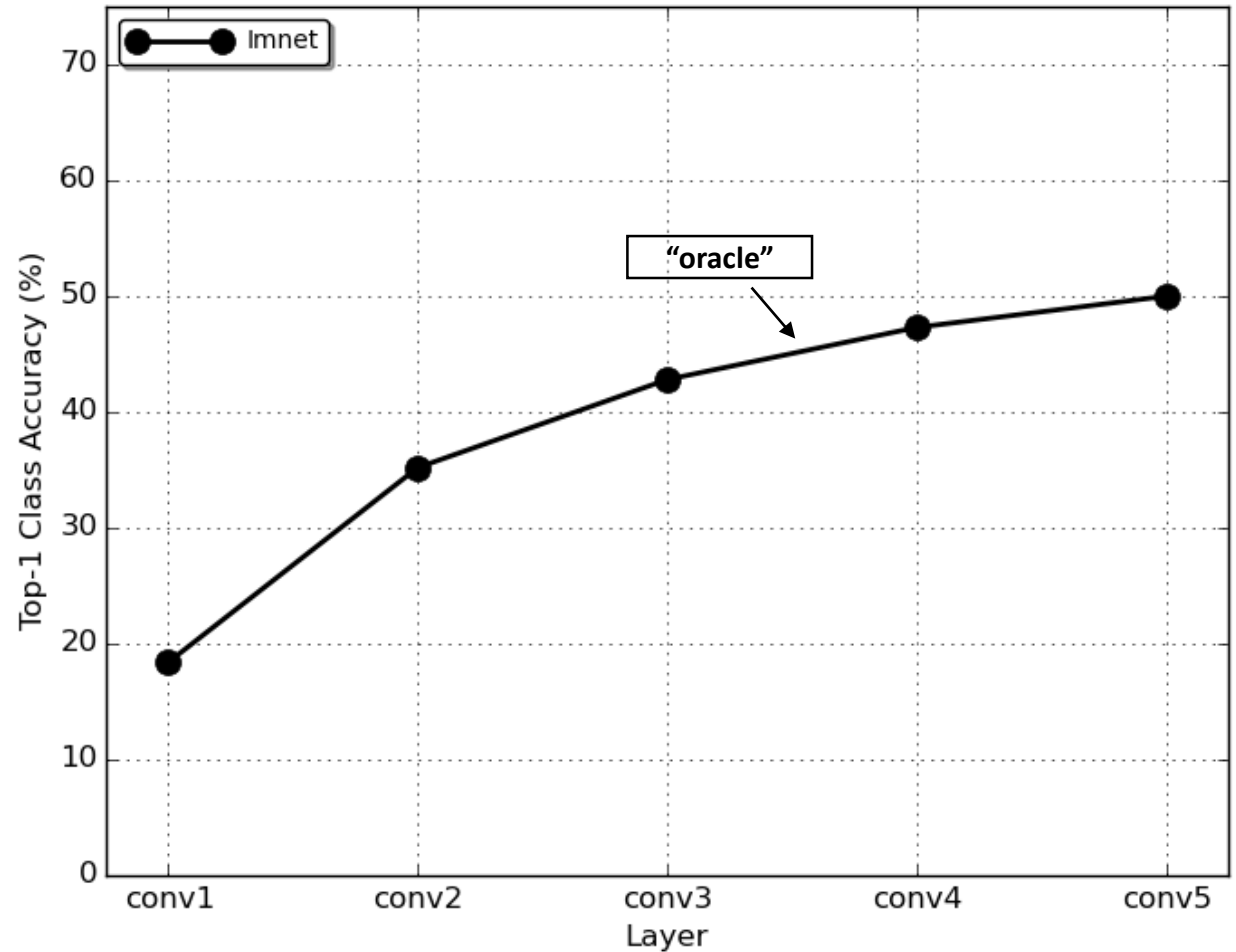






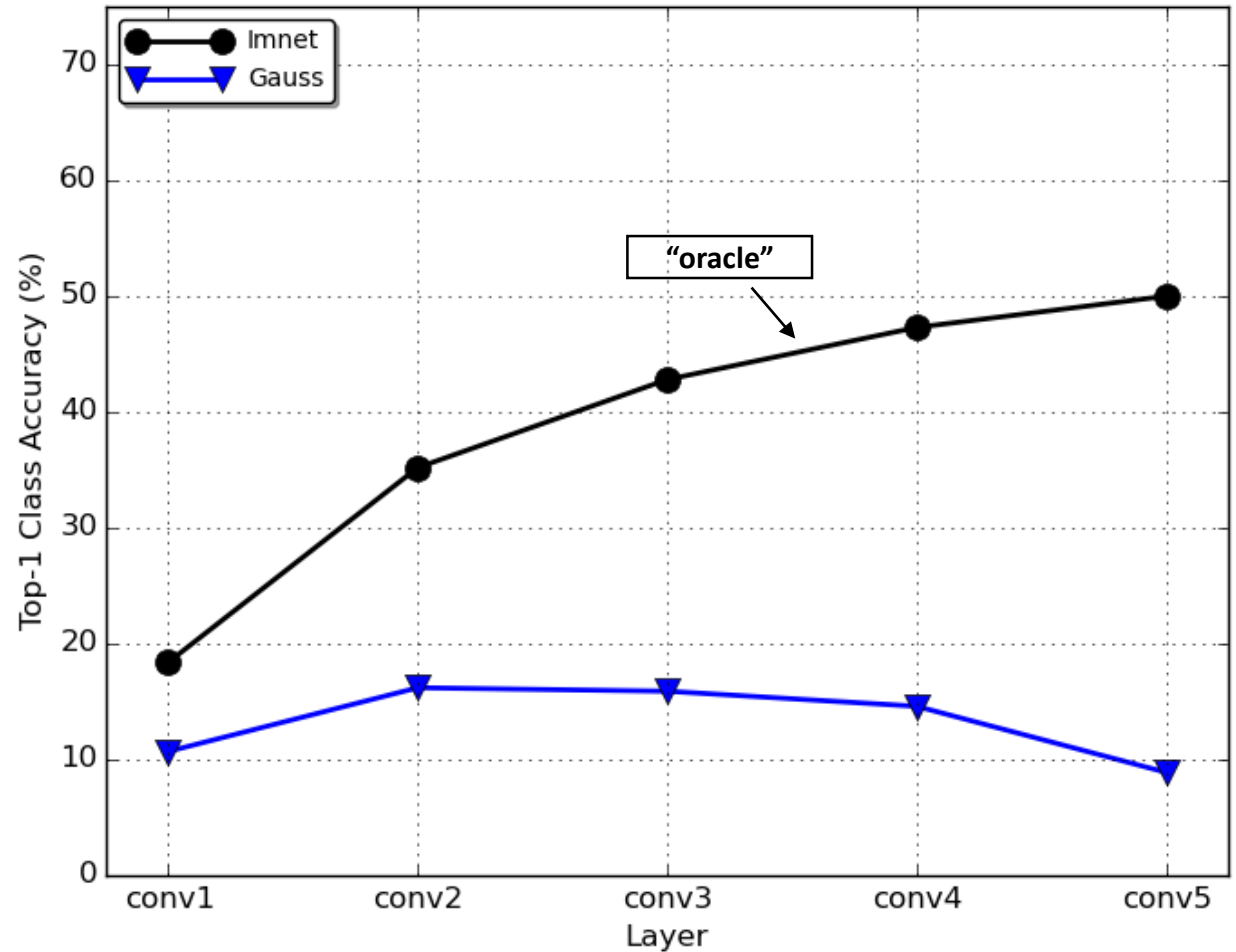
# Task Generalization - ILSVRC linear classification

- Directly training on labels provides “oracle”



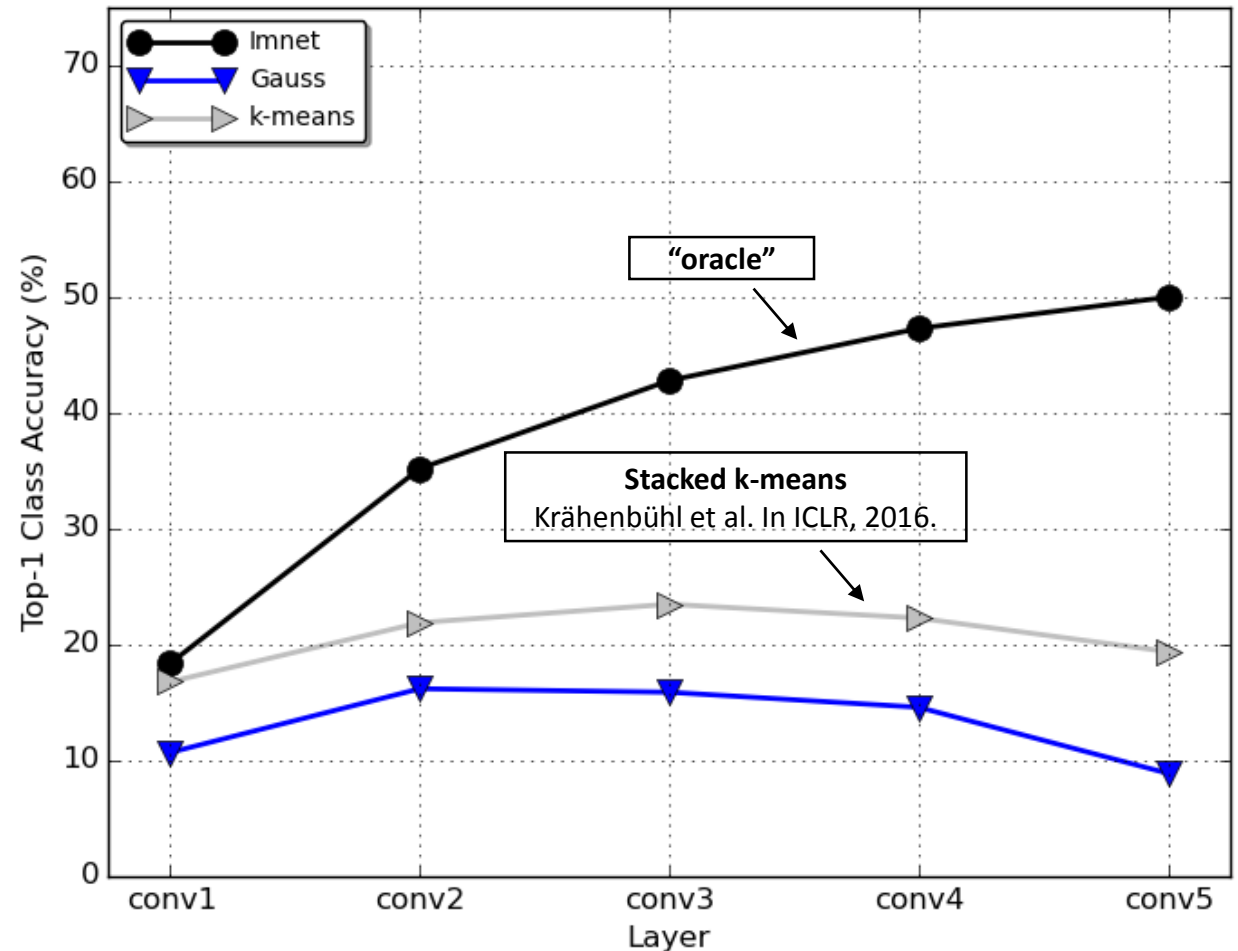
# Task Generalization - ILSVRC linear classification

- Directly training on labels provides “oracle”



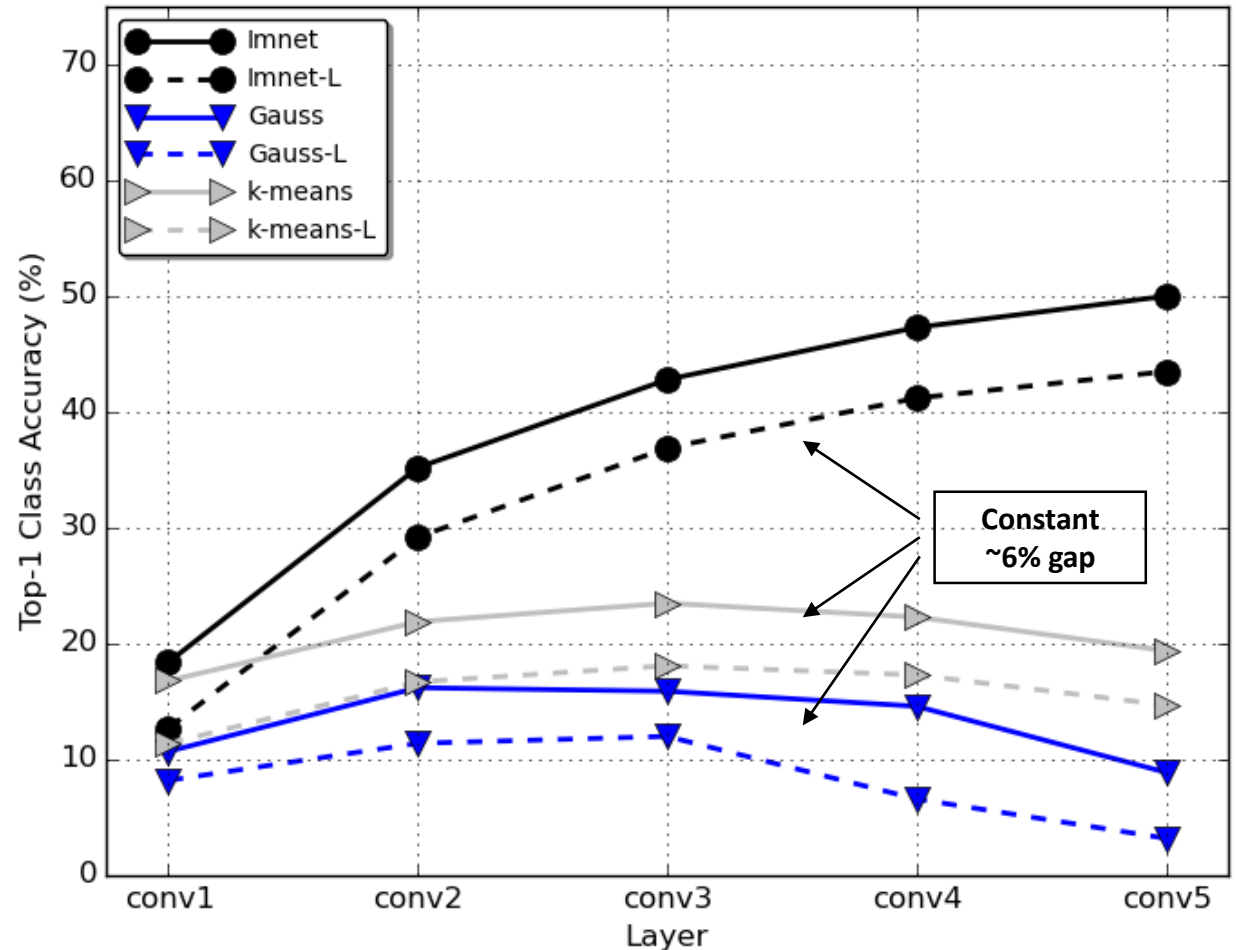
# Task Generalization - ILSVRC linear classification

- Directly training on labels provides “oracle”
- Stacked k-means provides strong baseline



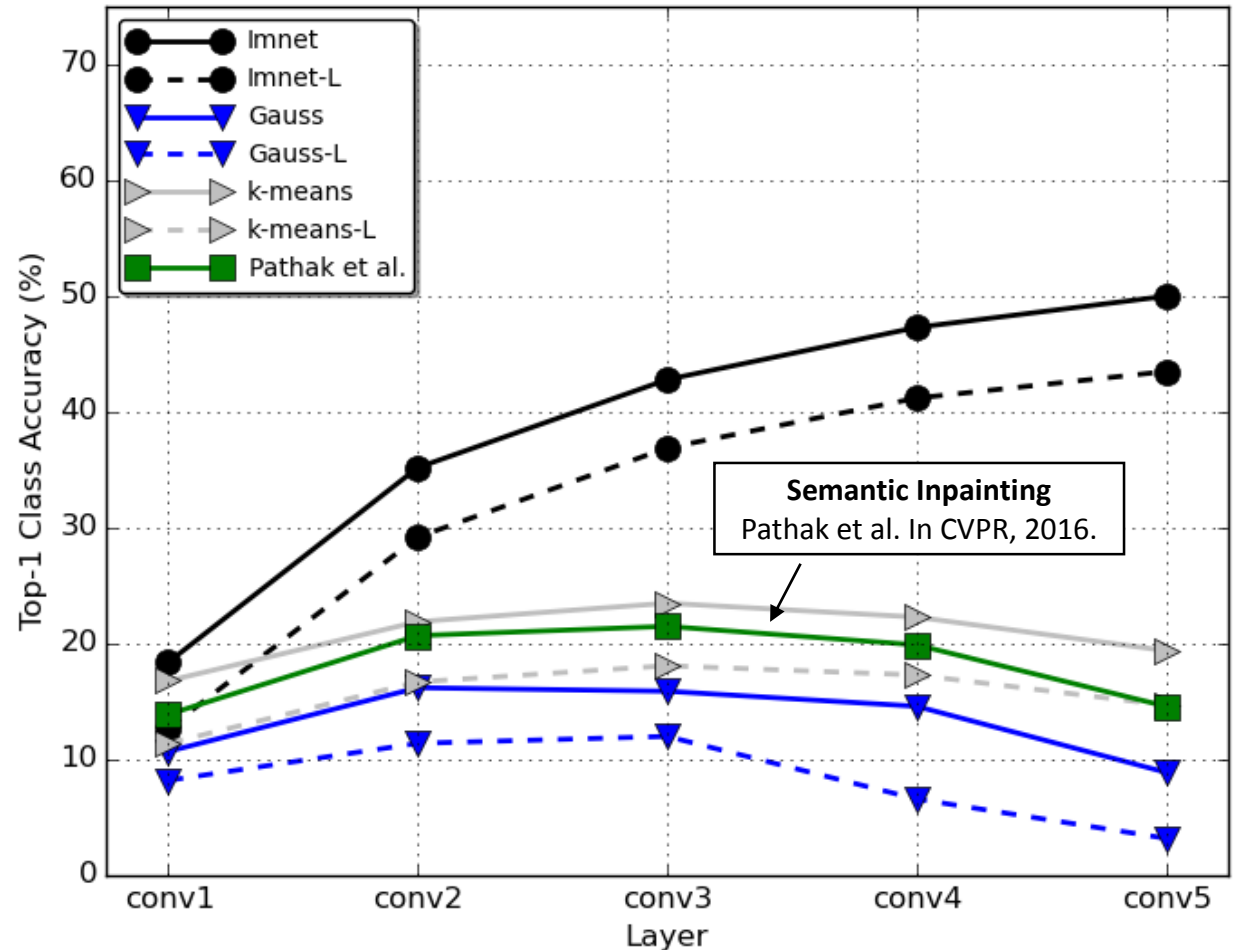
# Task Generalization - ILSVRC linear classification

- Directly training on labels provides “oracle”
- Stacked k-means provides strong baseline
- Constant 6% gap from grayscale handicap



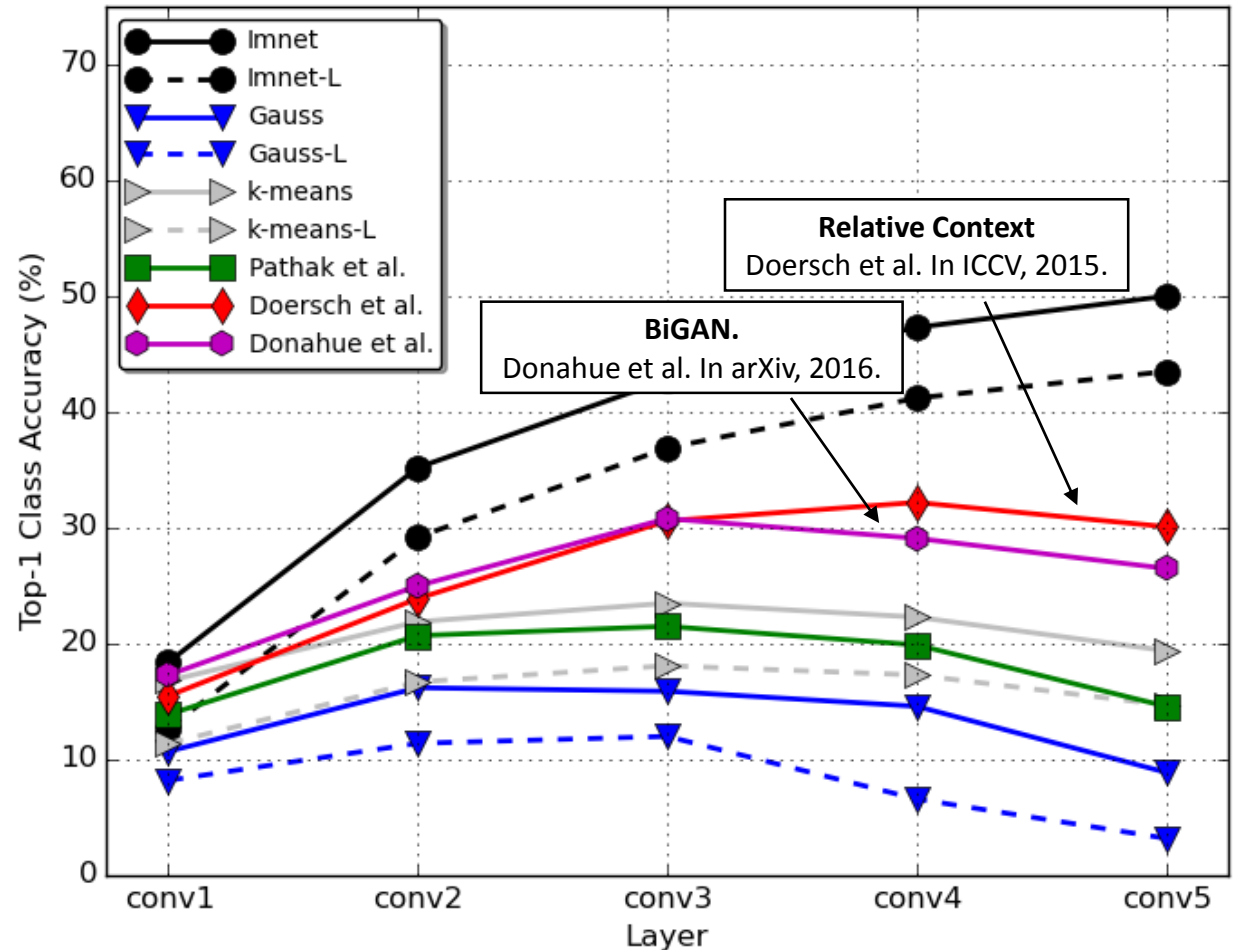
# Task Generalization - ILSVRC linear classification

- Directly training on labels provides “oracle”
- Stacked k-means provides strong baseline
- Constant 6% gap from grayscale handicap



# Task Generalization - ILSVRC linear classification

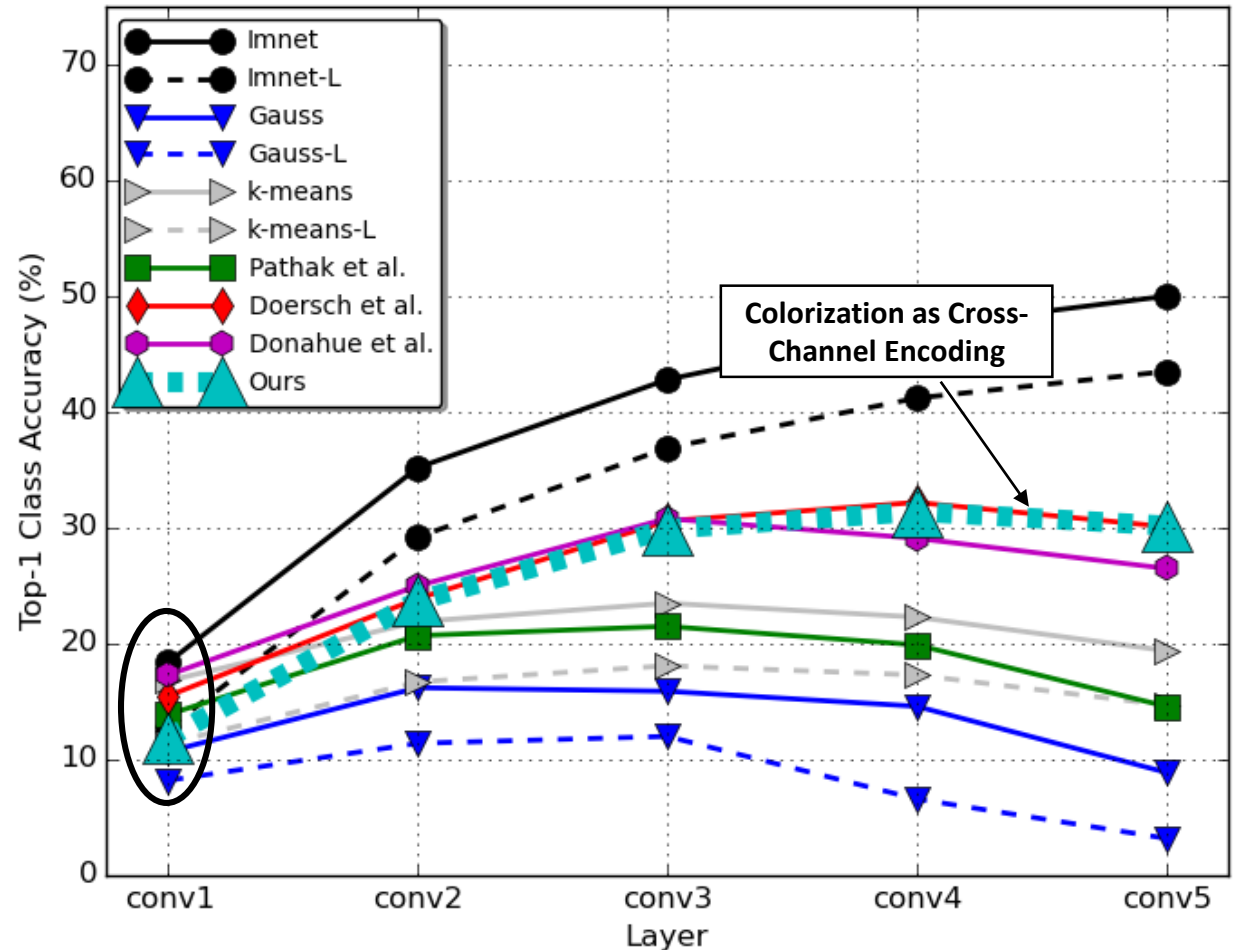
- Directly training on labels provides “oracle”
- Stacked k-means provides strong baseline
- Constant 6% gap from grayscale handicap





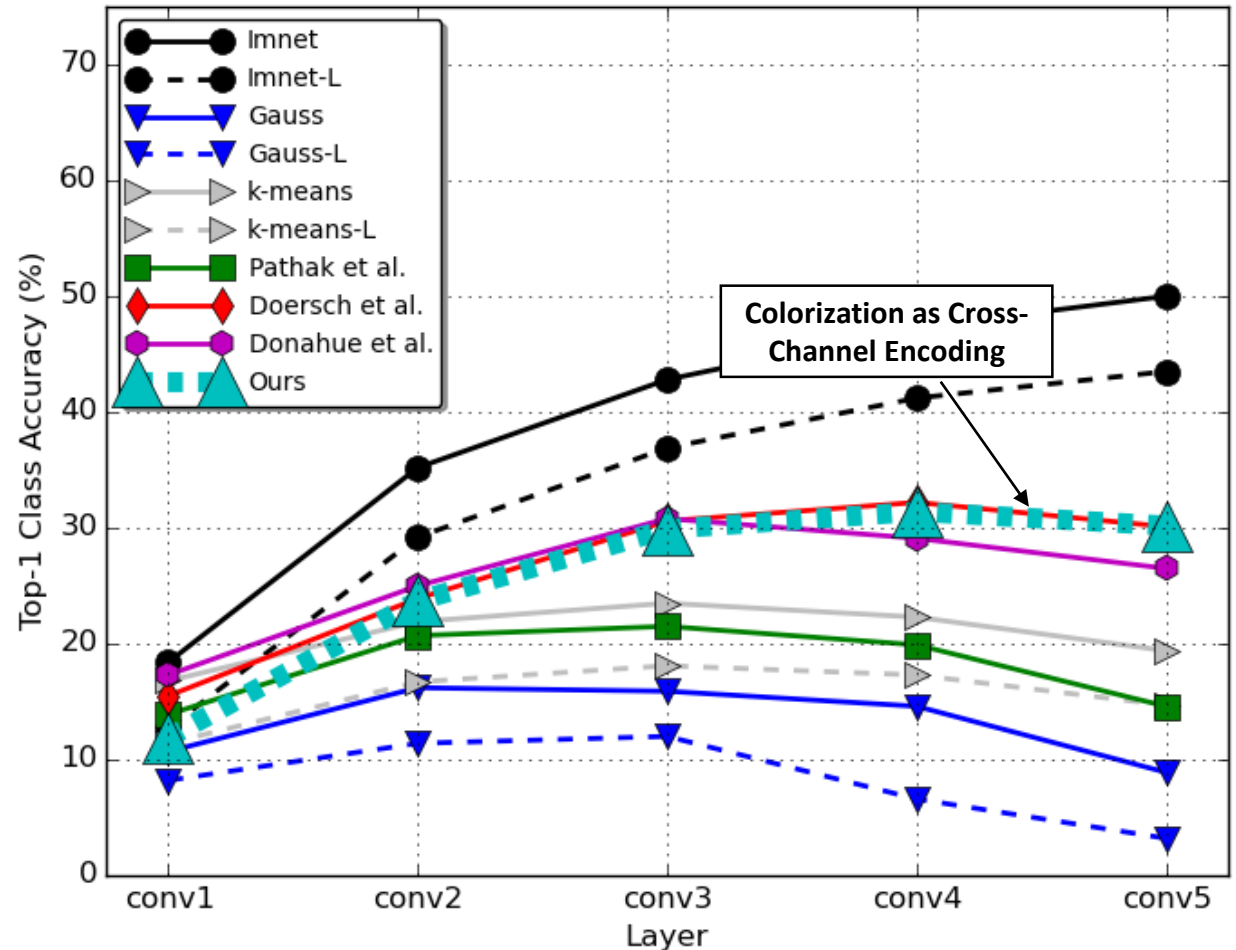
# Task Generalization - ILSVRC linear classification

- Directly training on labels provides “oracle”
- Stacked k-means provides strong baseline
- Constant 6% gap from grayscale handicap
  - Our *conv1* suffers from input handicap



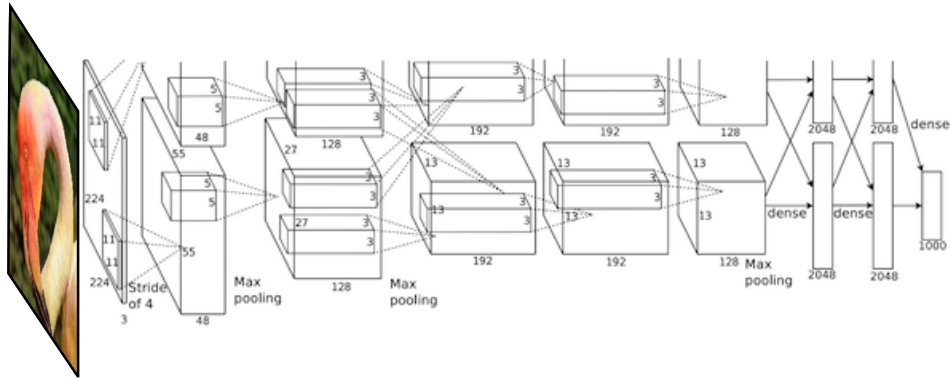
# Task Generalization - ILSVRC linear classification

- Directly training on labels provides “oracle”
- Stacked k-means provides strong baseline
- Constant 6% gap from grayscale handicap
  - Our *conv1* suffers from input handicap
- Our *conv2-5* performs competitively throughout



# Predicting Labels from Data

Image



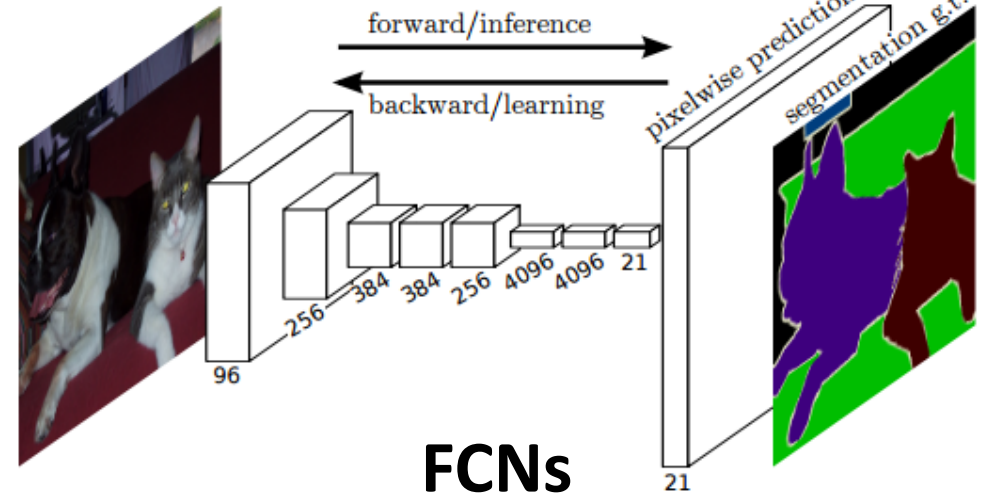
Label

“flamingo”

**Alexnet**

Krizhevsky et al. In *NIPS*, 2012.

Image



Dense Labels

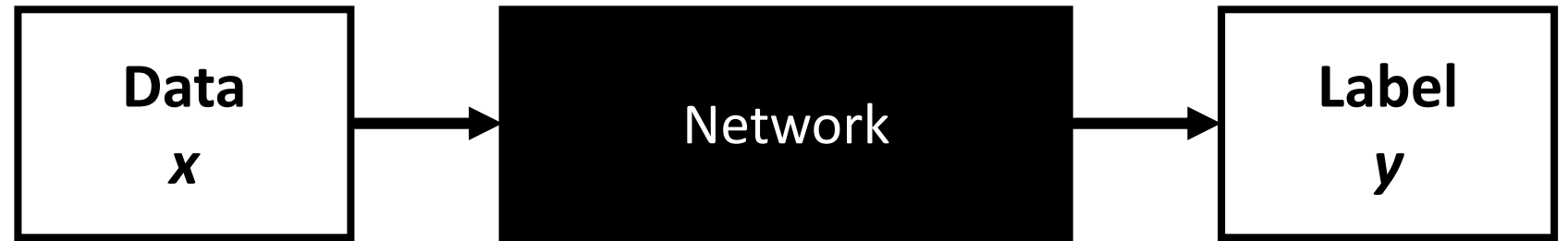
**FCNs**

Long et al. In *CVPR*, 2015.



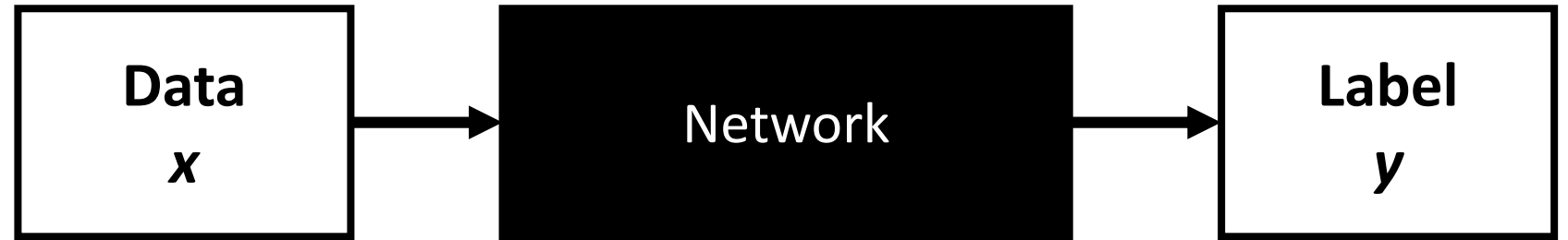
# Predicting Labels from Data

**Supervised  
training**



# Predicting Data from Data

**Supervised  
training**

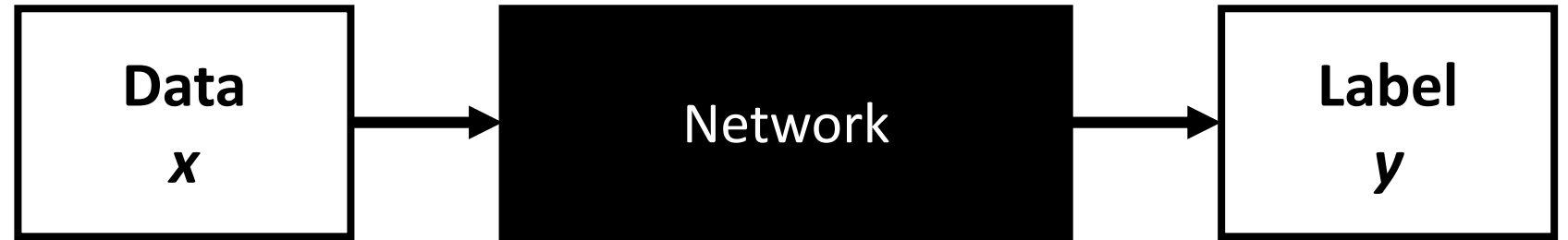


**Self-supervised  
training**

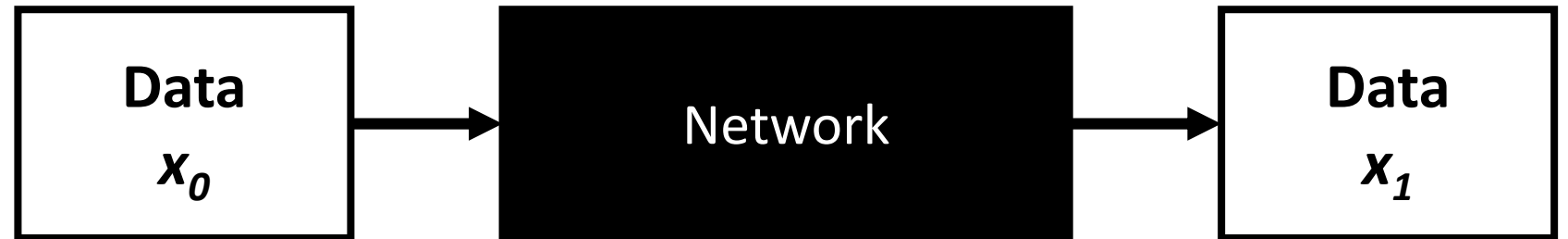


# Predicting Data from Data

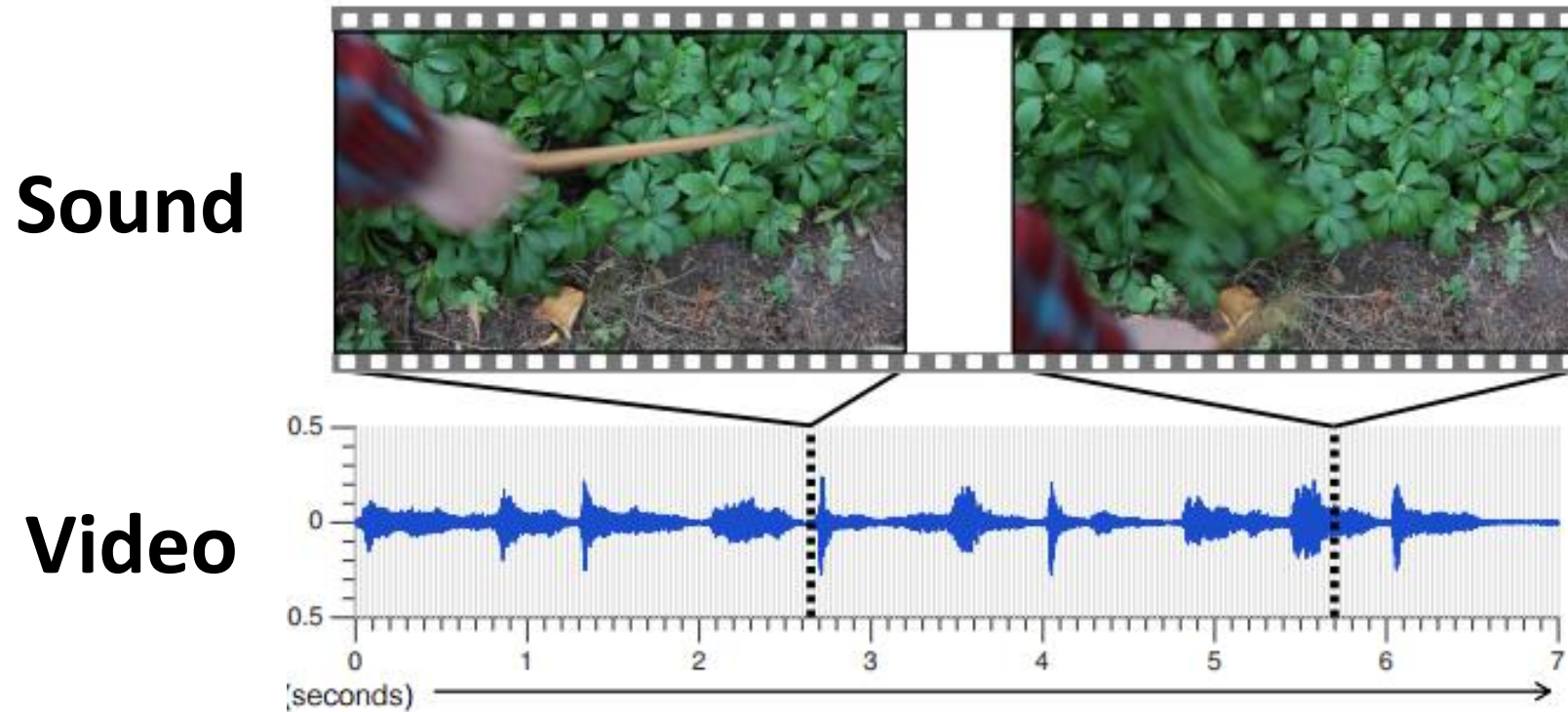
**Supervised  
training**



**Self-supervised  
training**



# Visually Indicated Sounds



Owens et al. **Visually Indicated Sounds**. In *CVPR*, 2016.



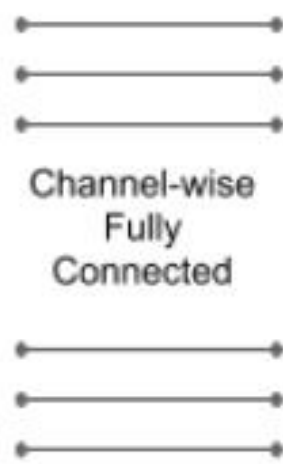
# Context Encoders

**Context pixels**



Encoder

Encoder Features



Decoder Features

Decoder

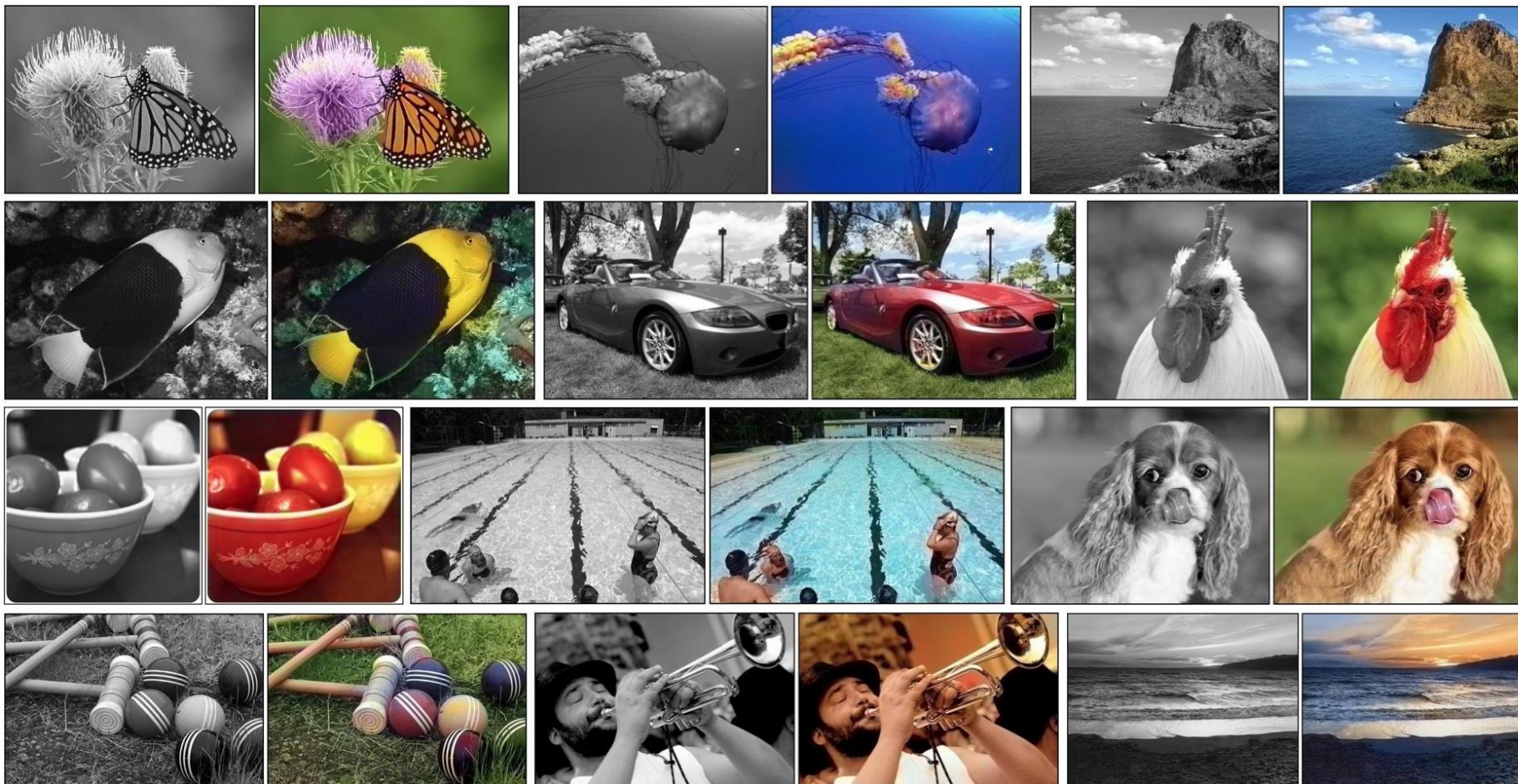
**Center pixels**



Pathak et al. **Context Encoders: Feature Learning by Inpainting**. In *CVPR*, 2016.







# Colourful Image Colourization

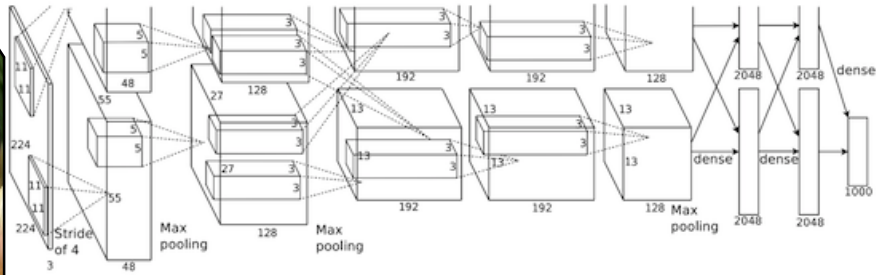
Richard Zhang, Phillip Isola, Alexei A. Efros

In ArXiv, March 2016.

[richzhang.github.io/colorization](http://richzhang.github.io/colorization)

# Predicting Labels from Data

**Image**



**Label**

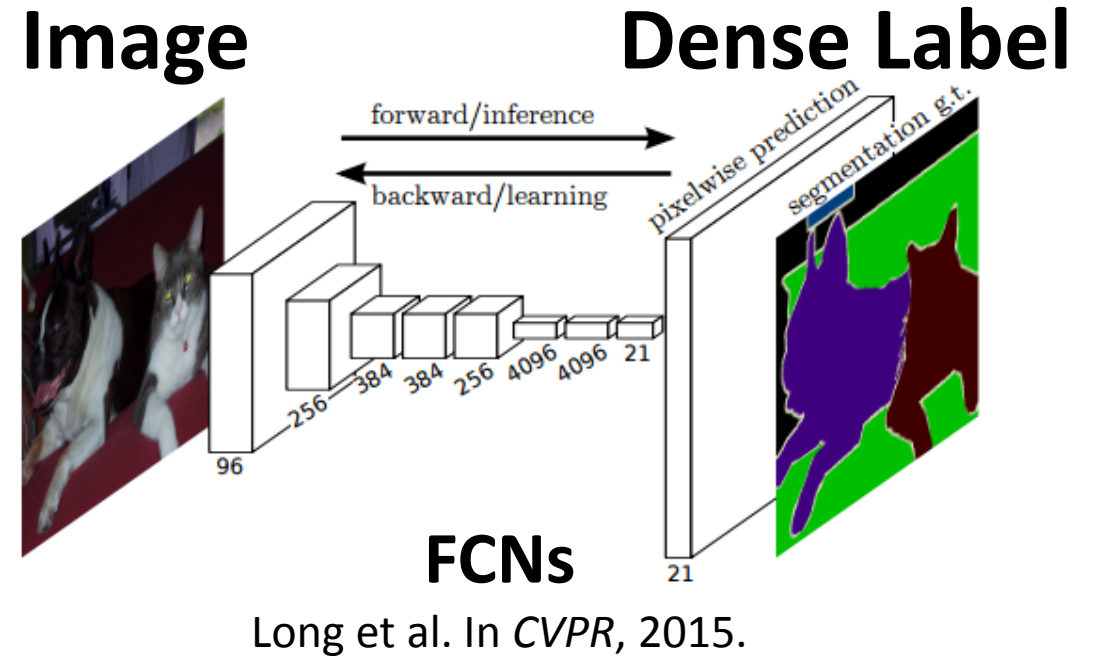
“flamingo”

**Alexnet**

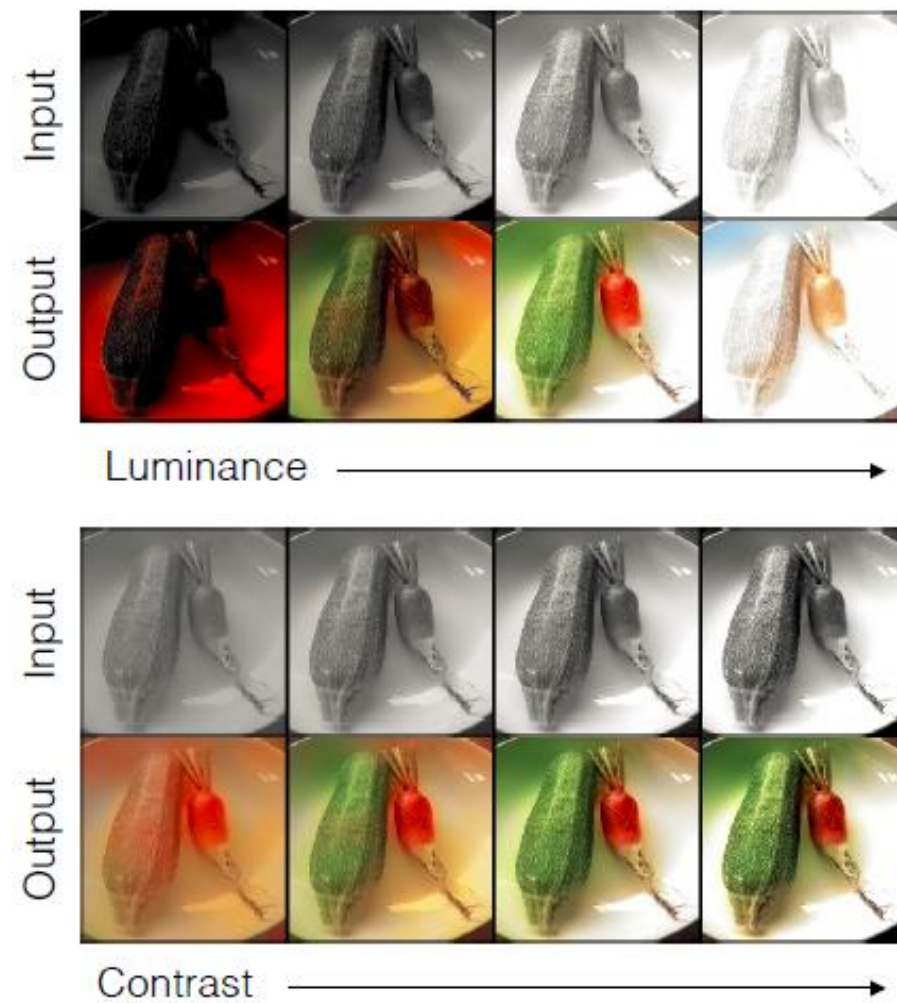
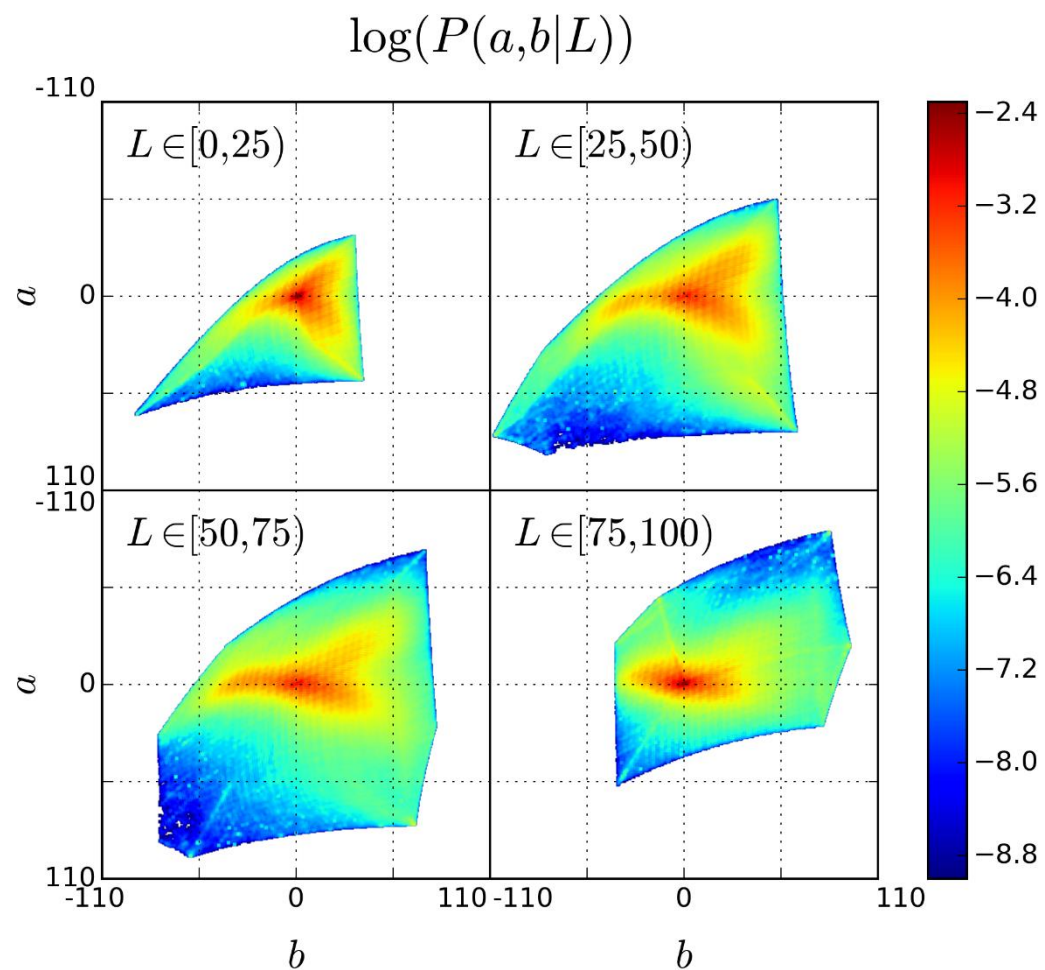
Krizhevsky et al. In *NIPS*, 2012.



# Predicting Labels from Data



# Low-level Perturbations



# Common Confusions

jacamar

Ground truth



standard schnauzer

Ground truth



# Common Confusions

jacamar  
bulbul

Ground  
truth



Recolored



standard schnauzer  
irish terrier

Ground  
truth



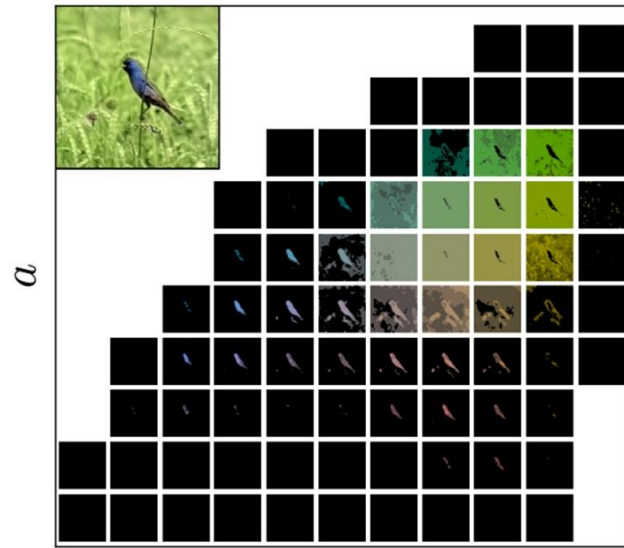
Recolored



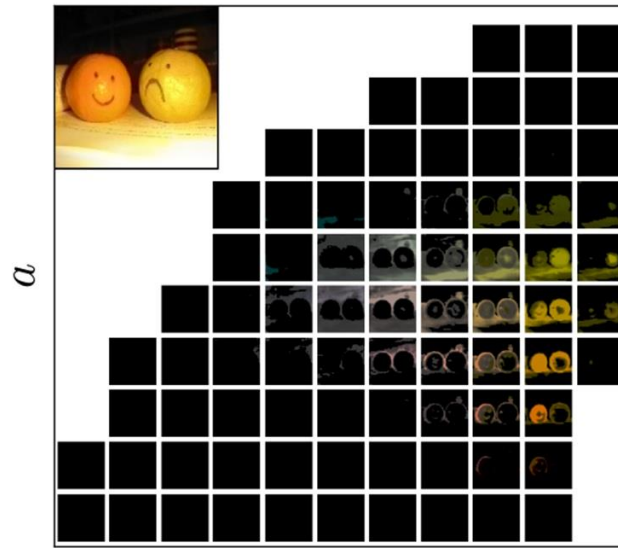
# Future steps

- Perceptual Losses
- Back-propagate end-to-end
- Train on “infinite” data
- Domain gap to legacy black and white images

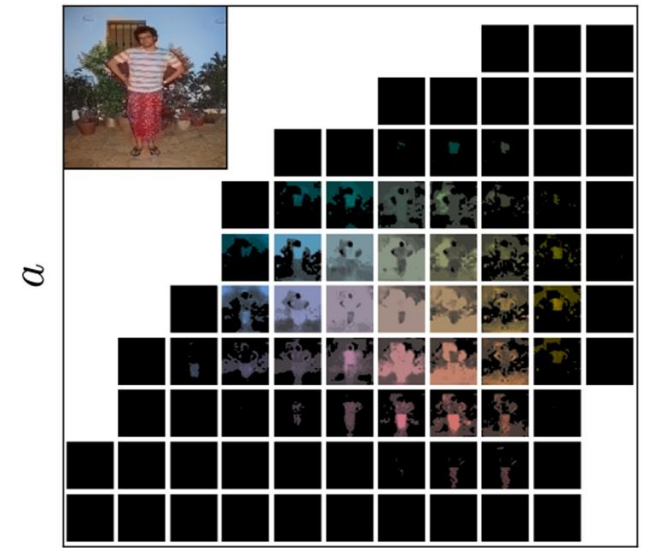
# Example Output distribution



(a)

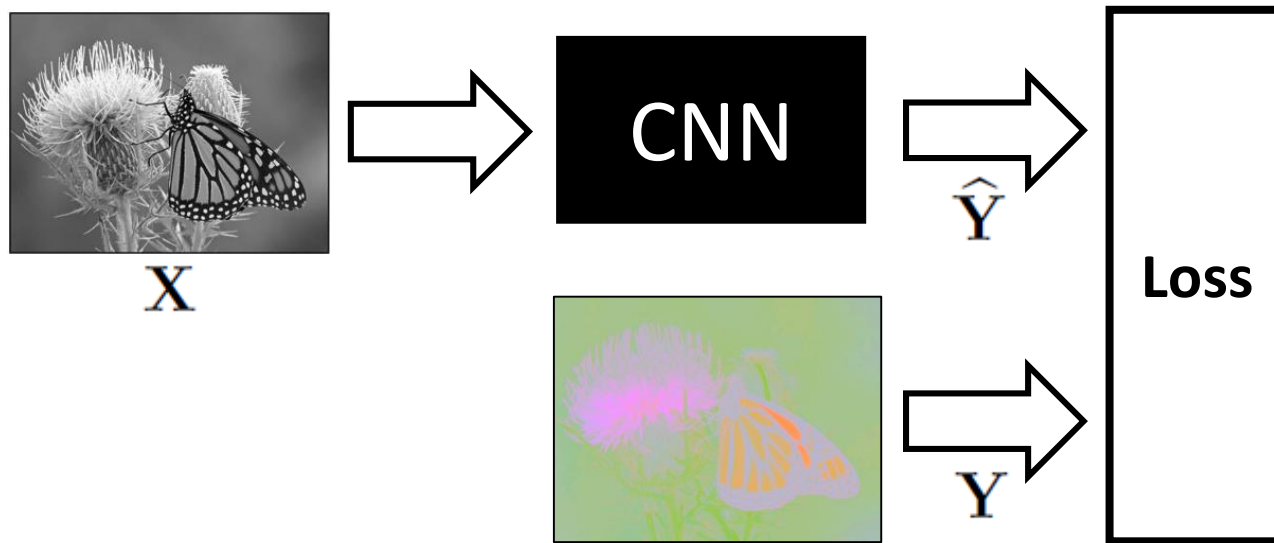


(b)



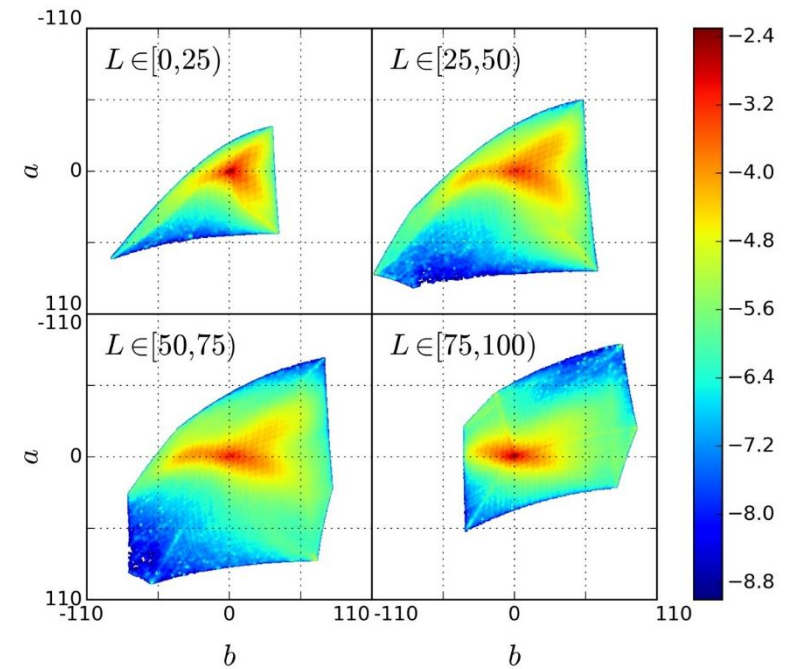
(c)





# Color statistics

Histogram over  $ab$  space  
Conditioned on  $L$



# LOSS

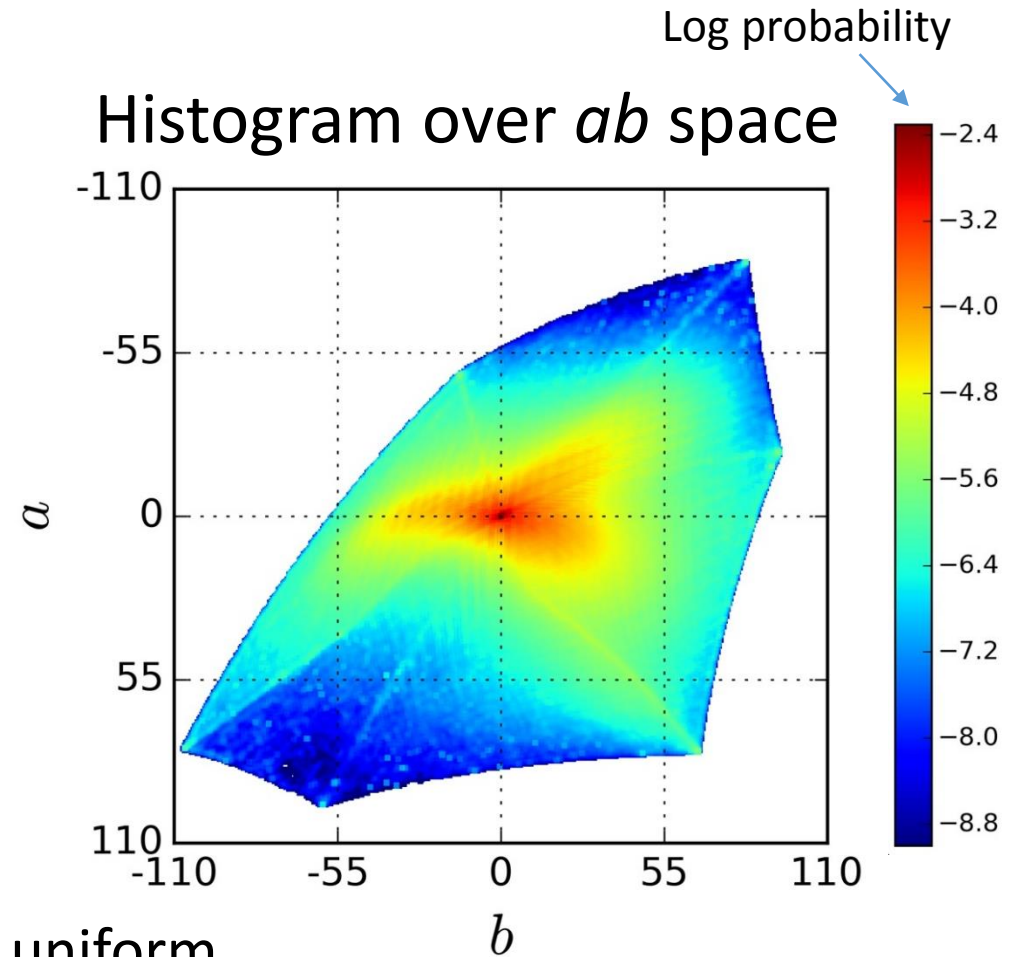
- Regression with L2 loss will not address inherent ambiguity
- Use ***multinomial classification***
  - quantize  $ab$  space into grid size 10
  - cross entropy loss
- ***Class rebalancing*** at train time to encourage learning of *rare* colors

$$w \propto ((1 - \lambda)(\mathbf{G}_\sigma \circ \mathbf{p}) + \lambda)^{-1}$$

Reweighting factor

empirical distribution

combine with uniform



# Probability Distribution to Point Estimate

$$\mathcal{H}(\mathbf{Z}_{h,w}) = \mathbb{E}(f_T(\log \mathbf{Z}_{h,w})), \quad f_T(\mathbf{z}) = \frac{\exp(\mathbf{z}/T)}{\sum_q \exp(\mathbf{z}_q/T)}$$

Mean  
T=1

T=.77

T=.58

T=.38

T=.29

T=.14

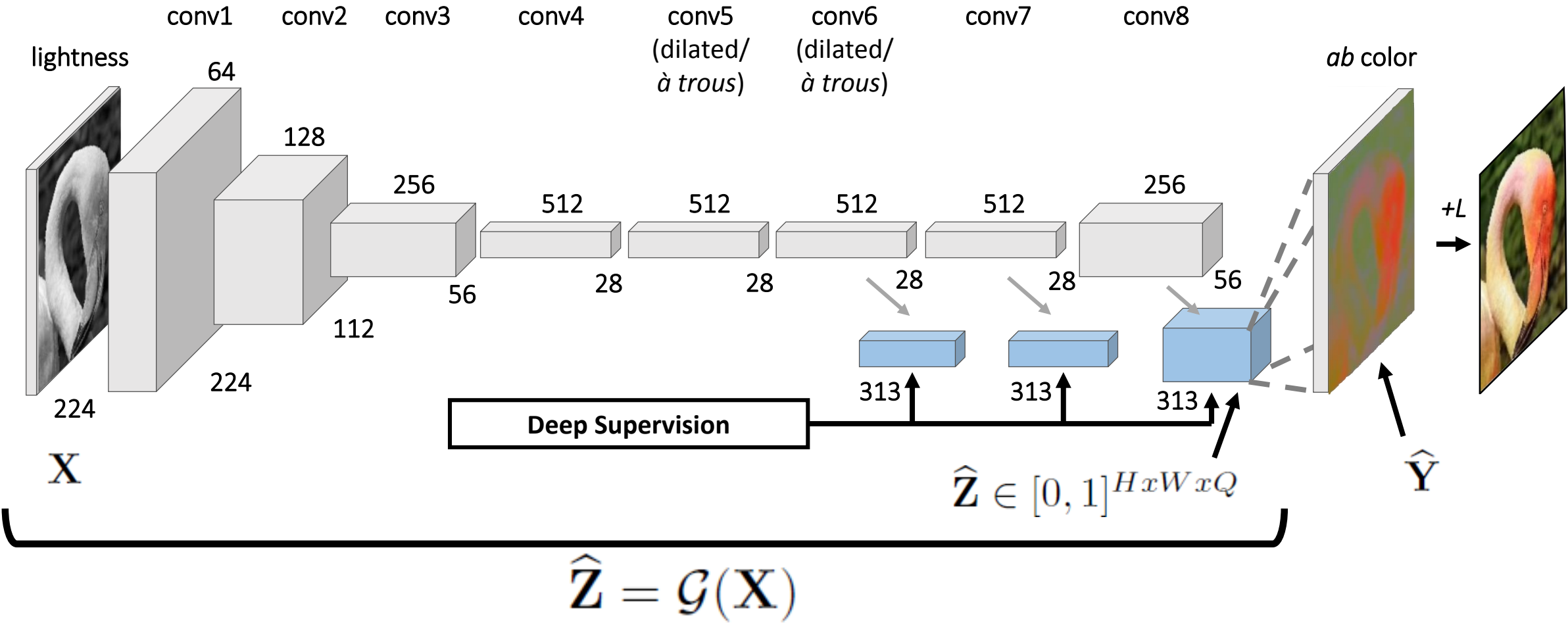
Mode  
T→0



Lowering softmax temperature T



# Network Architecture



# Network Architecture

