

Self-supervised Representation Learning

Computer Vision
Fall 2018
Columbia University

Project

- Presentation schedule posted on Piazza.
- Review it ASAP and let us know of any problems by Wednesday
- For those presenting on December 5: OK to have some experiments in progress
- Final reports due December 10 midnight — no extensions!

GPU Credits

- If you have not requested GPU credits, do so immediately.
- We are starting to give them away...

Homeworks

- HW3 is back: median is 100% !
- HW4 grades soon
- HW5 due today

Final Grades

- We will likely curve down, but we will guarantee:
 - 90% is at least A
 - 80% is at least B
 - 70% is at least C

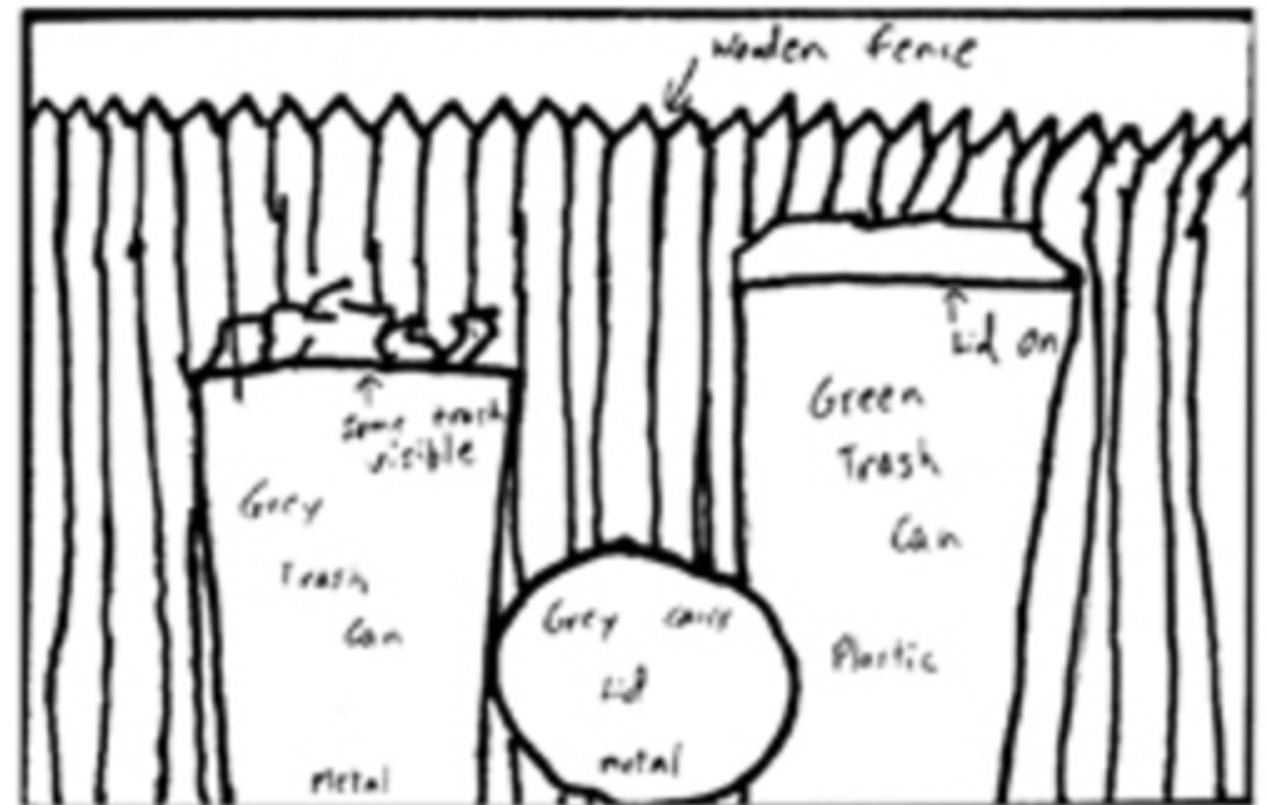
Next Semester

- E6998 Advanced Computer Vision, offered Spring 2019
- Focuses on research frontier of computer vision and applied machine learning
- Make sure to fill out survey to get off wait list

Observed image



Drawn from memory

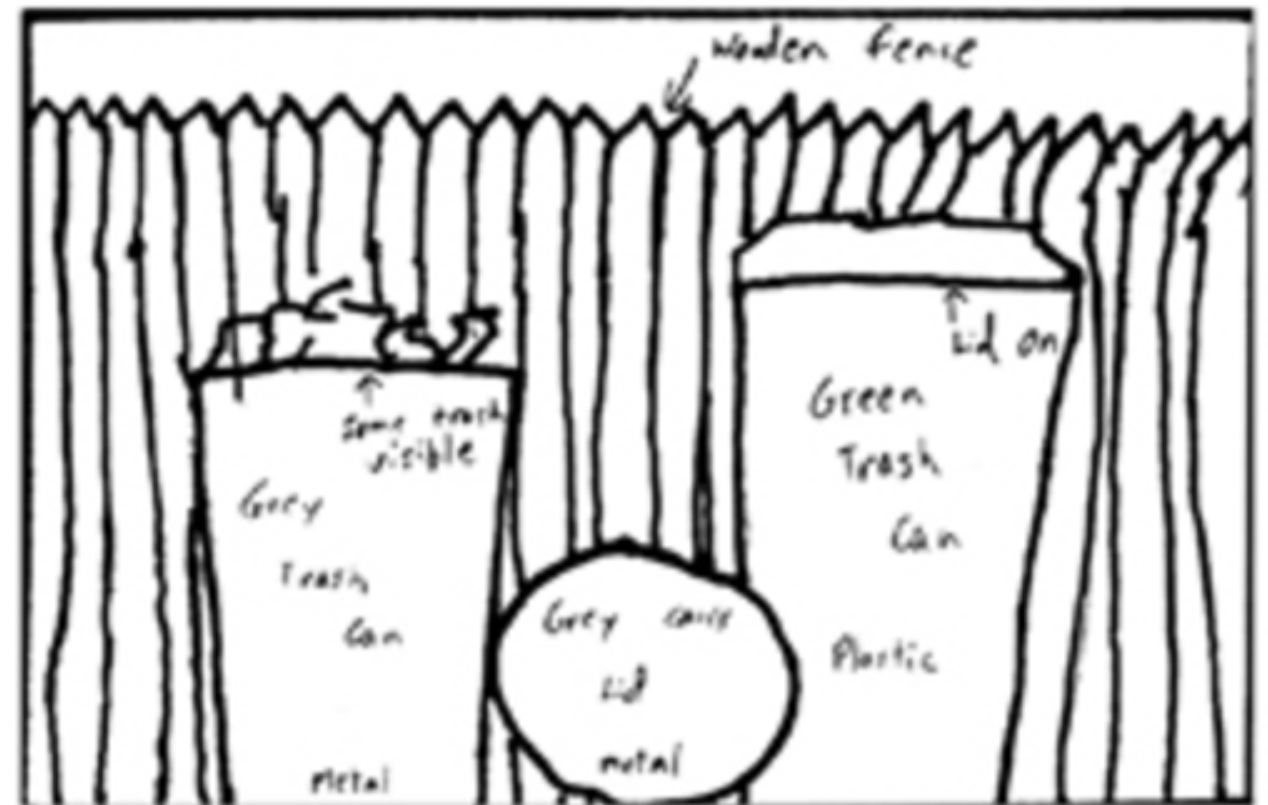


[Bartlett, 1932]
[Intraub & Richardson, 1989]

Observed image



Drawn from memory



[Bartlett, 1932]
[Intraub & Richardson, 1989]



"I stand at the window and see a house, trees, sky. Theoretically I might say there were 327 brightnesses and nuances of colour. Do I have "327"? No. I have sky, house, and trees."

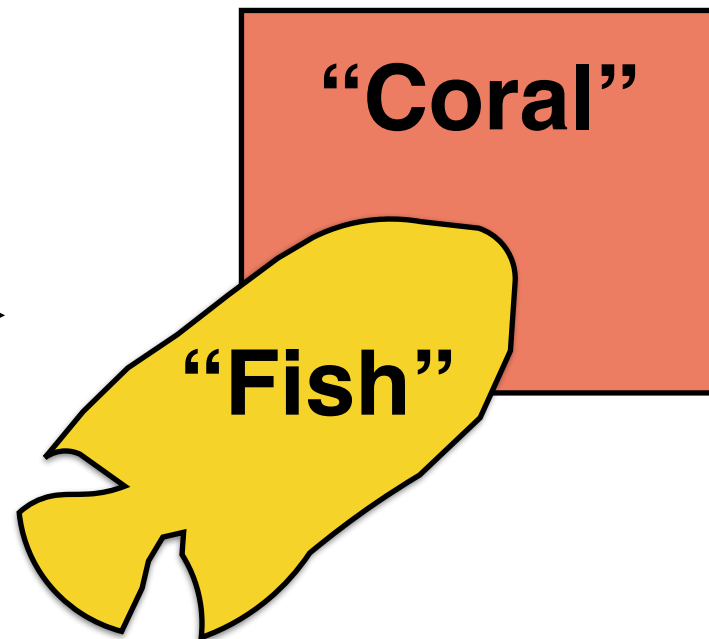
— Max Wertheimer, 1923

Representation learning

X



Image

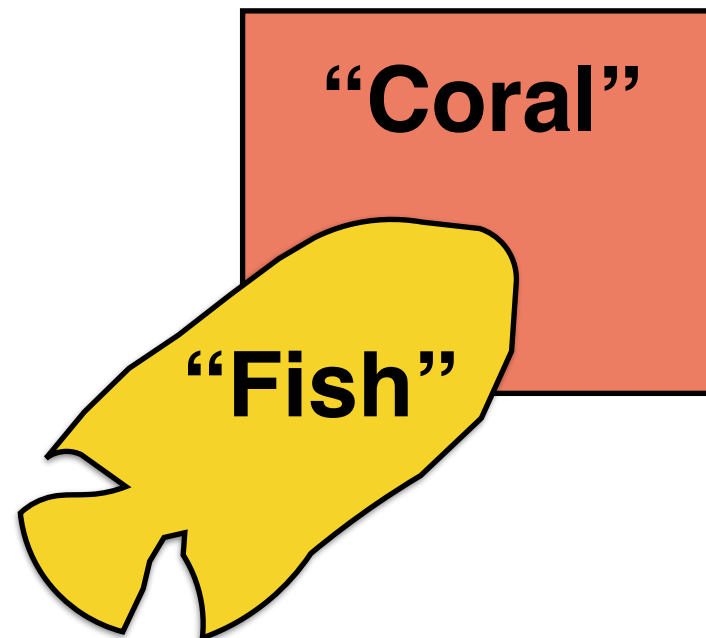


Compact mental
representation

Representation learning

Good representations are:

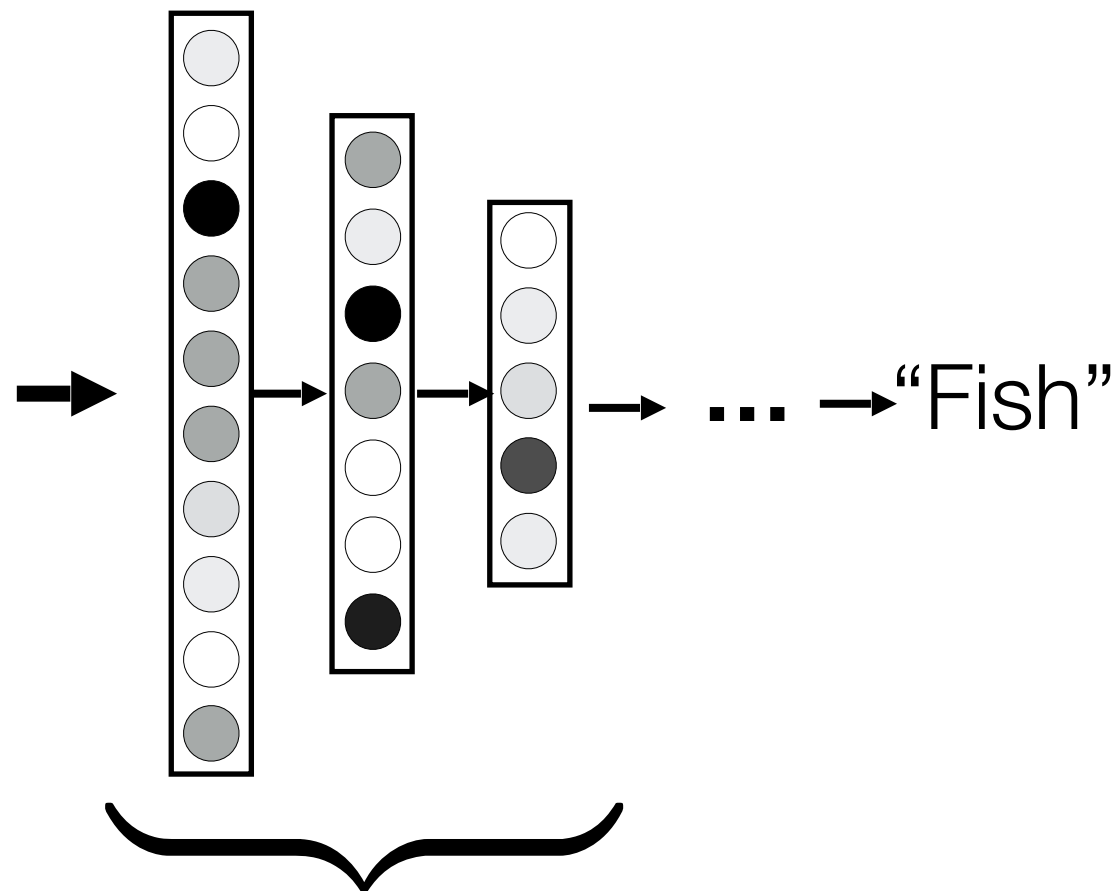
1. Compact
2. Explanatory
3. Disentangled
4. Interpretable



X



Image

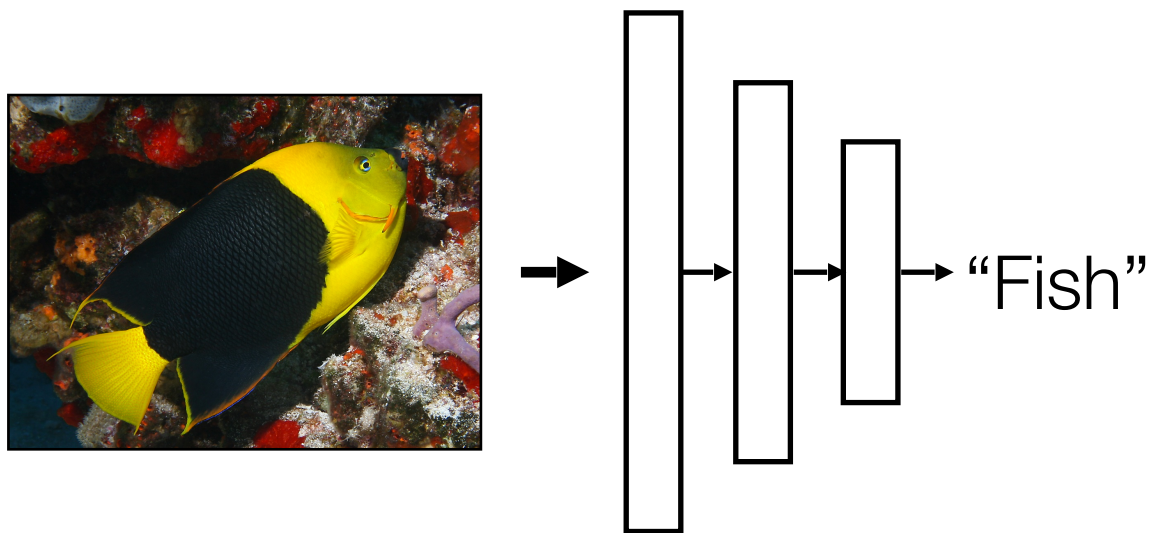


Representations!

A CNN is a multiscale, hierarchical representation of data

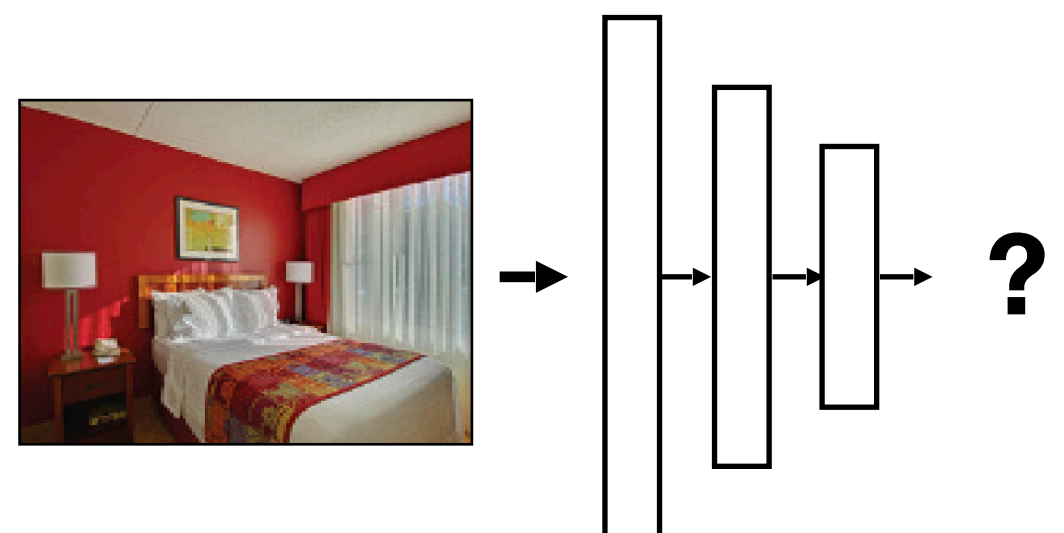
Training

Object recognition



Testing

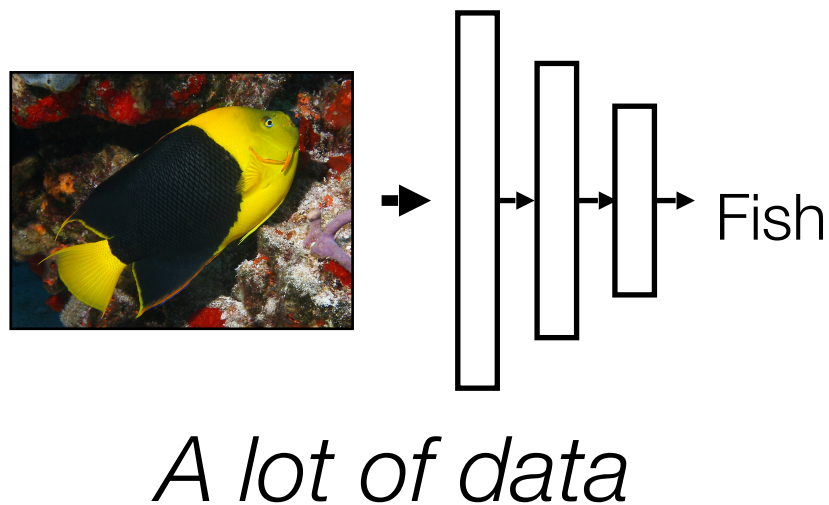
Place recognition



Often, what we will be “tested” on is to learn to do a new thing.

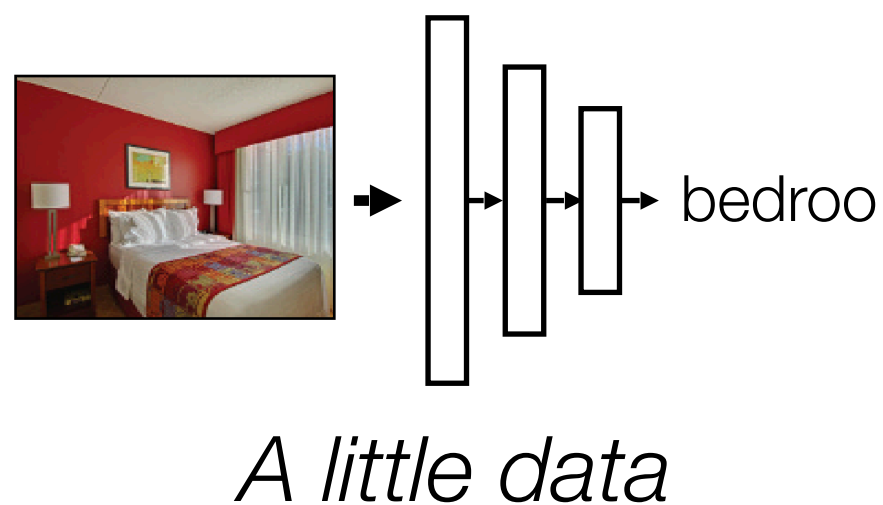
Pretraining

Object recognition



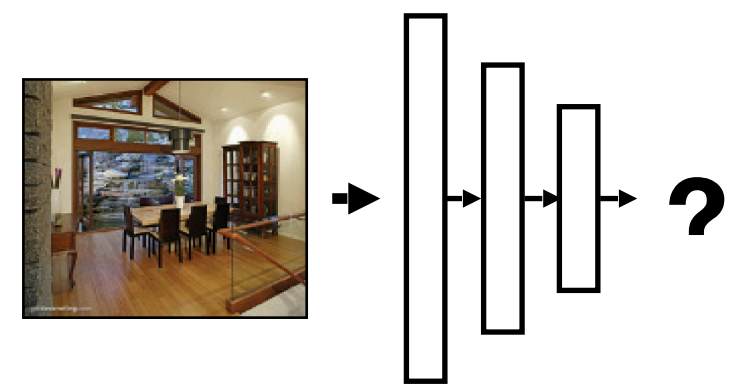
Finetuning

Place recognition



Testing

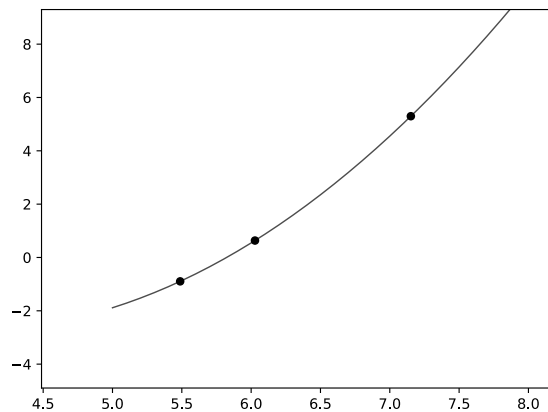
Place recognition



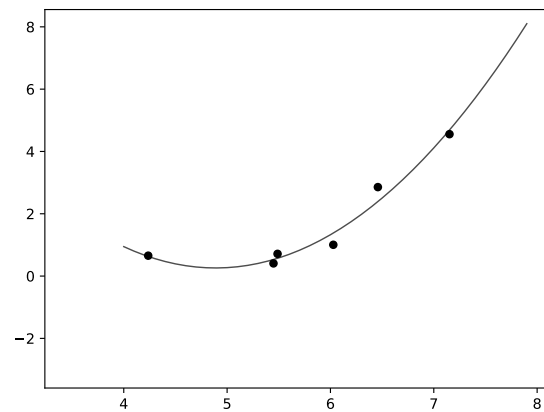
Finetuning starts with the representation learned on a previous task, and adapts it to perform well on a new task.

If we keep on finetuning for every new datapoint or task that comes our way, we get **online learning**. Humans seem to do this, we never stop learning.

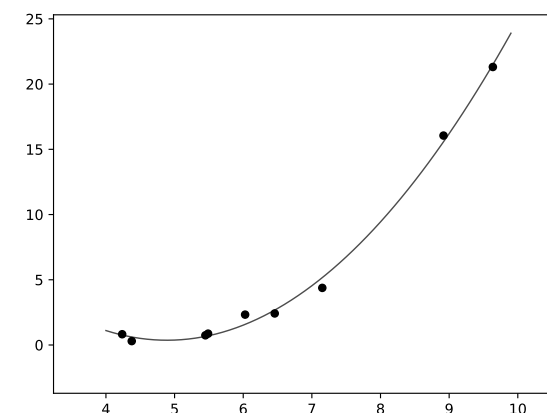
Training



More training



More training



...

Supervised object recognition

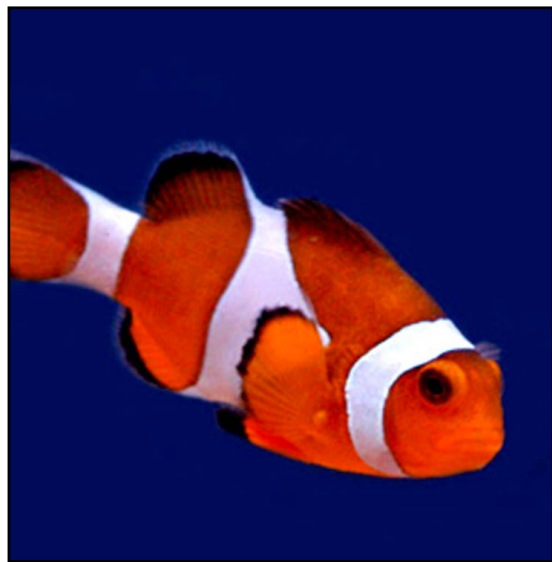


image X



Learner



"Fish"

label Y

Supervised object recognition



image X



"Fish"

label Y

Supervised object recognition



image X



"Fish"

label Y

Supervised object recognition



⋮

image X



Learner



"Fish"

label Y



Kitten Carousel

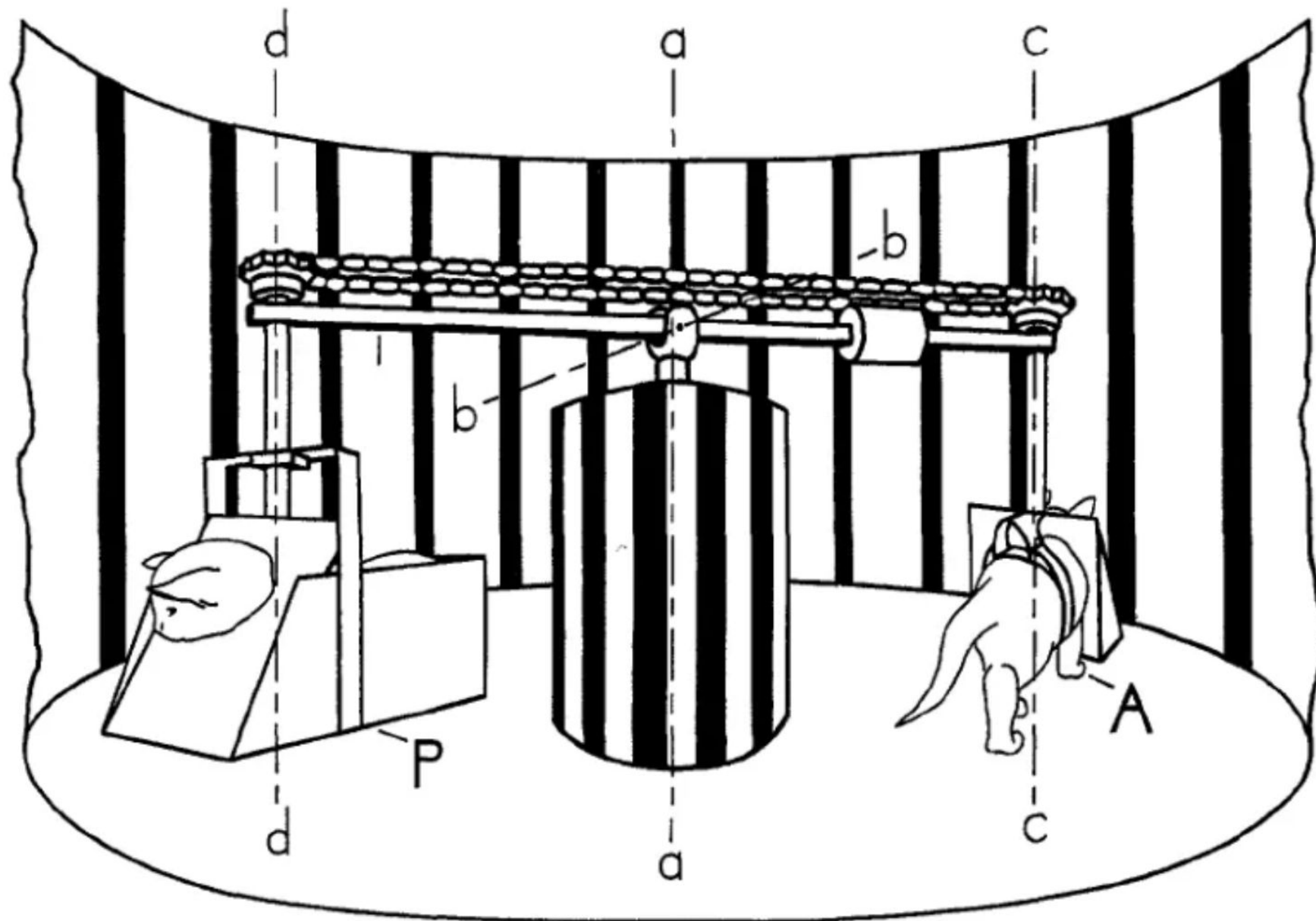


FIG. 1. Apparatus for equating motion and consequent visual feedback for an actively moving (A) and a passively moved (P) *S*.

Held and Hein, 1963

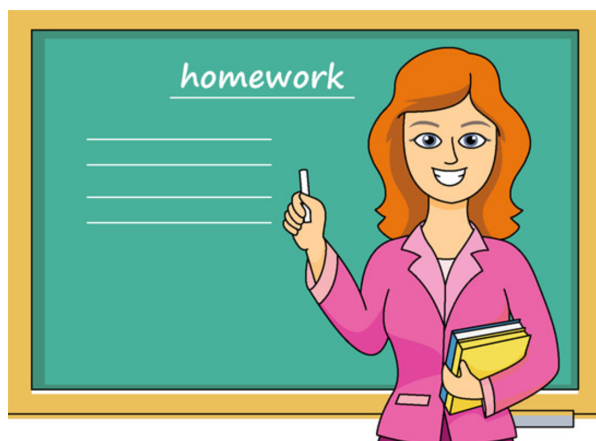
Supervised computer vision

Hand-curated training data

+ Informative

- Expensive

- Limited to teacher's
knowledge



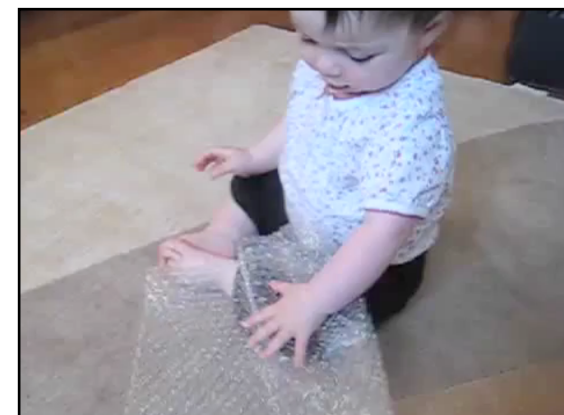
Vision in nature

Raw unlabeled training data

+ Cheap

- Noisy

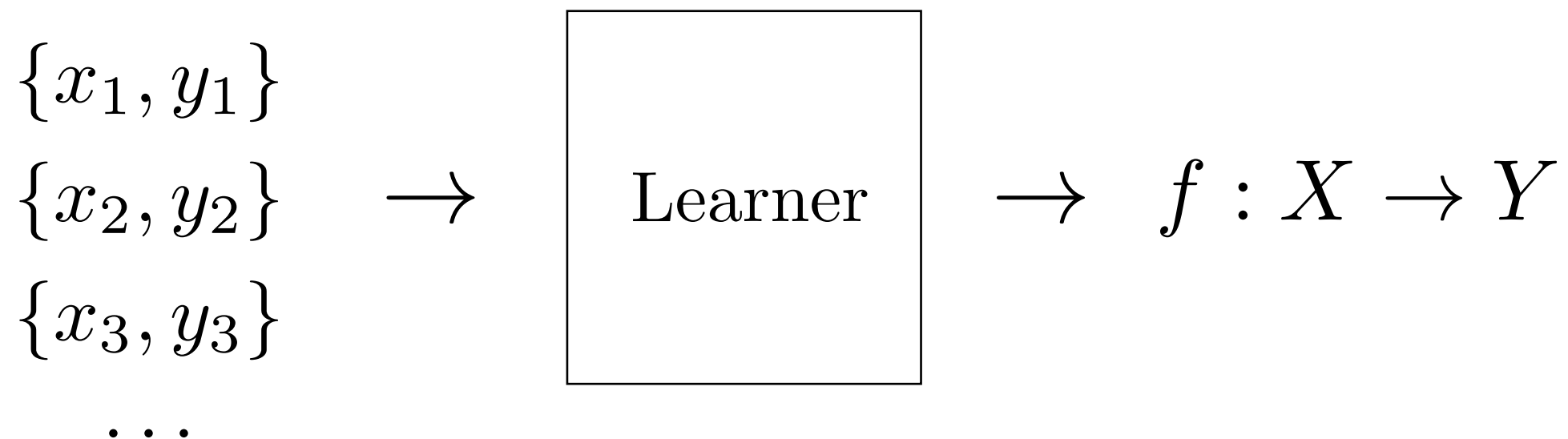
- Harder to interpret



Learning from examples

(aka **supervised learning**)

Training data



$$f^* = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^N \mathcal{L}(f(x_i), y_i)$$

Learning without examples

(includes **unsupervised learning** and **reinforcement learning**)

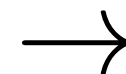
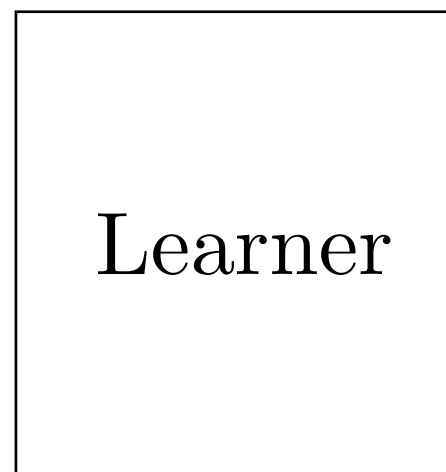
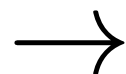
Data

$\{x_1\}$

$\{x_2\}$

$\{x_3\}$

...



?

Unsupervised Representation Learning

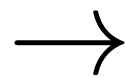
Data

$\{x_1\}$

$\{x_2\}$

$\{x_3\}$

...



Learner

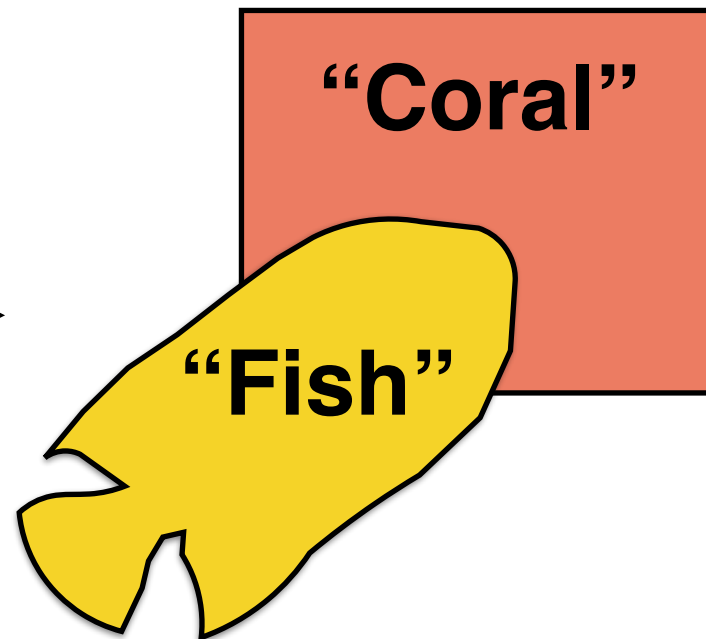
→ Representations

Unsupervised Representation Learning

X

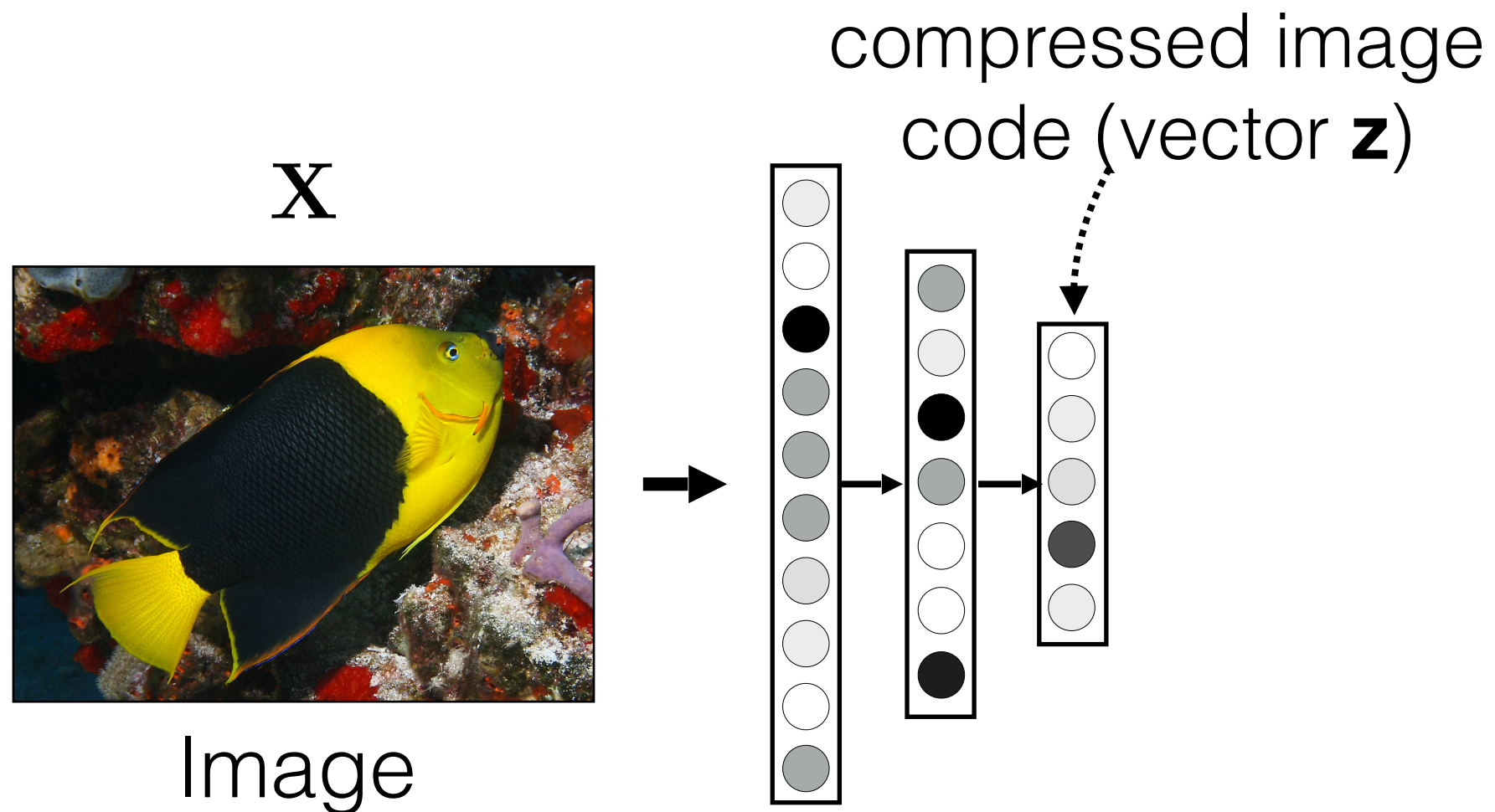


Image

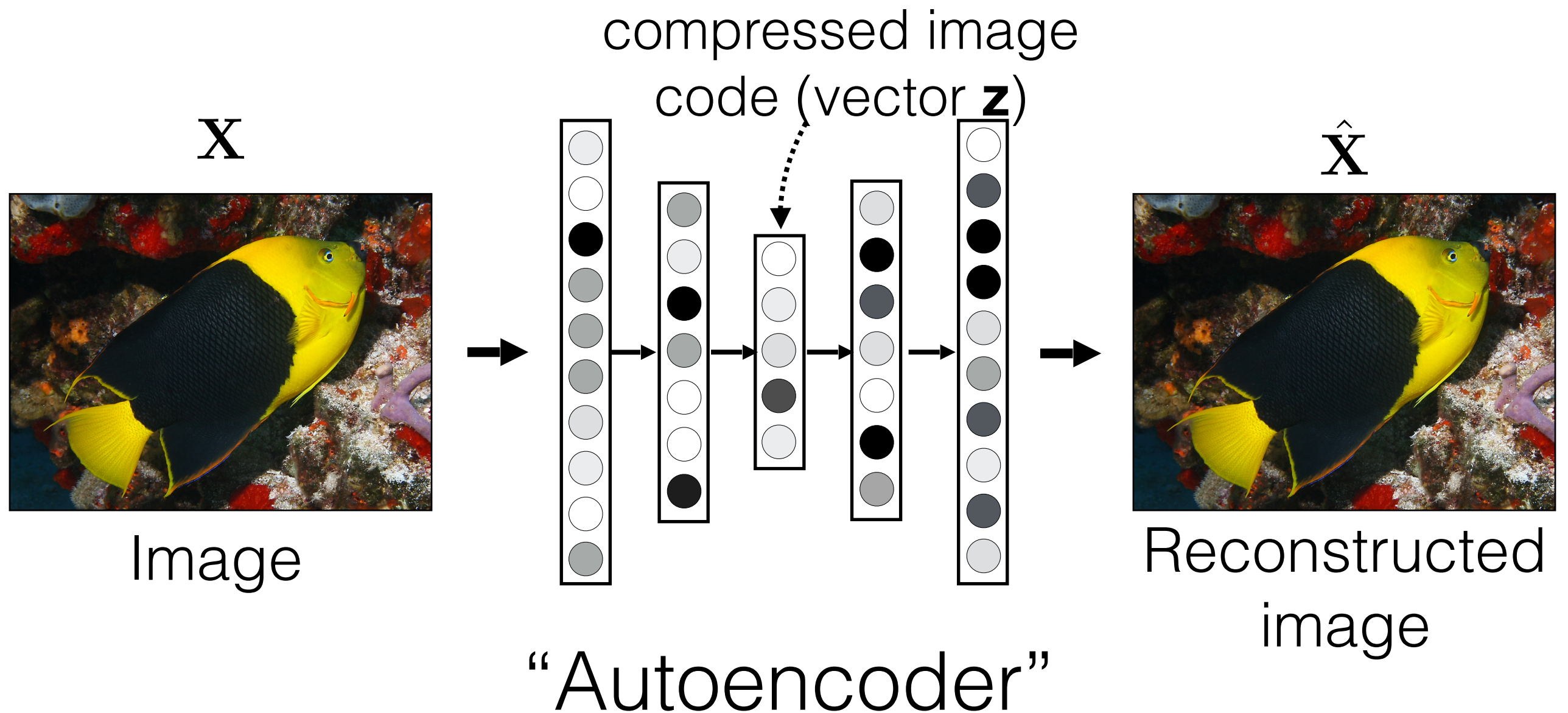


Compact mental
representation

Unsupervised Representation Learning

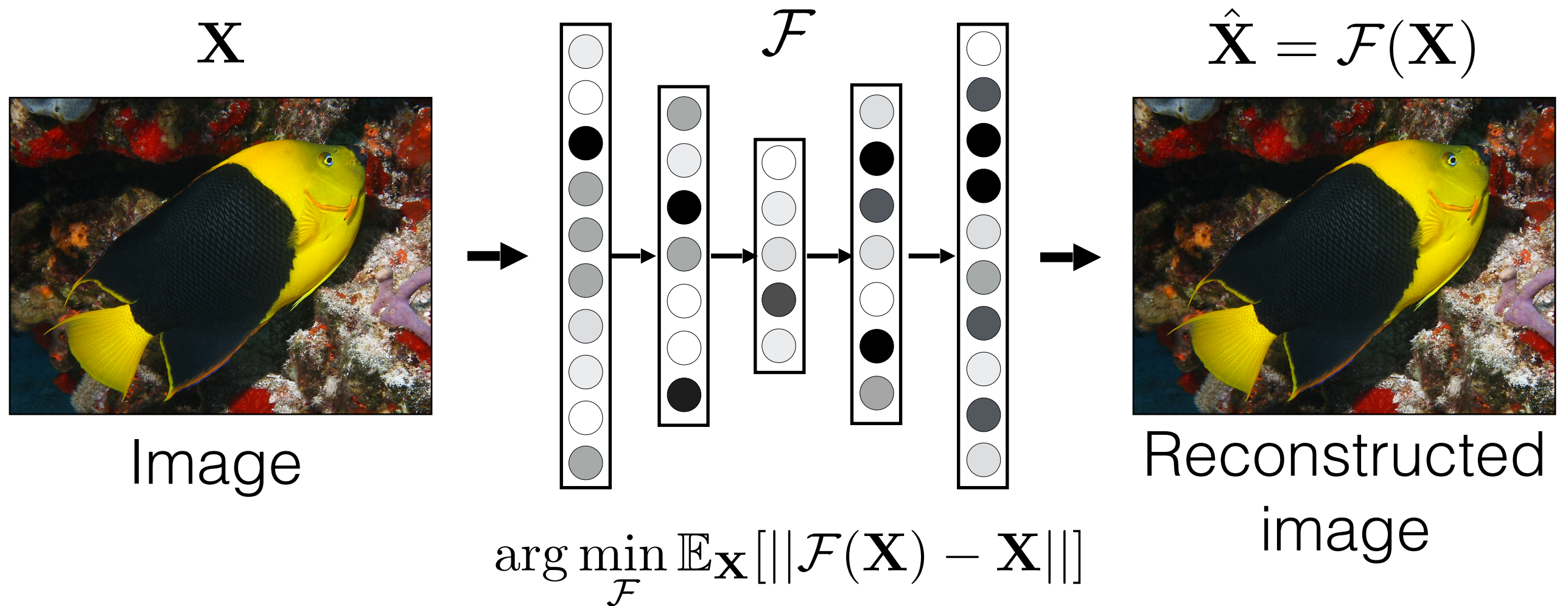


Unsupervised Representation Learning

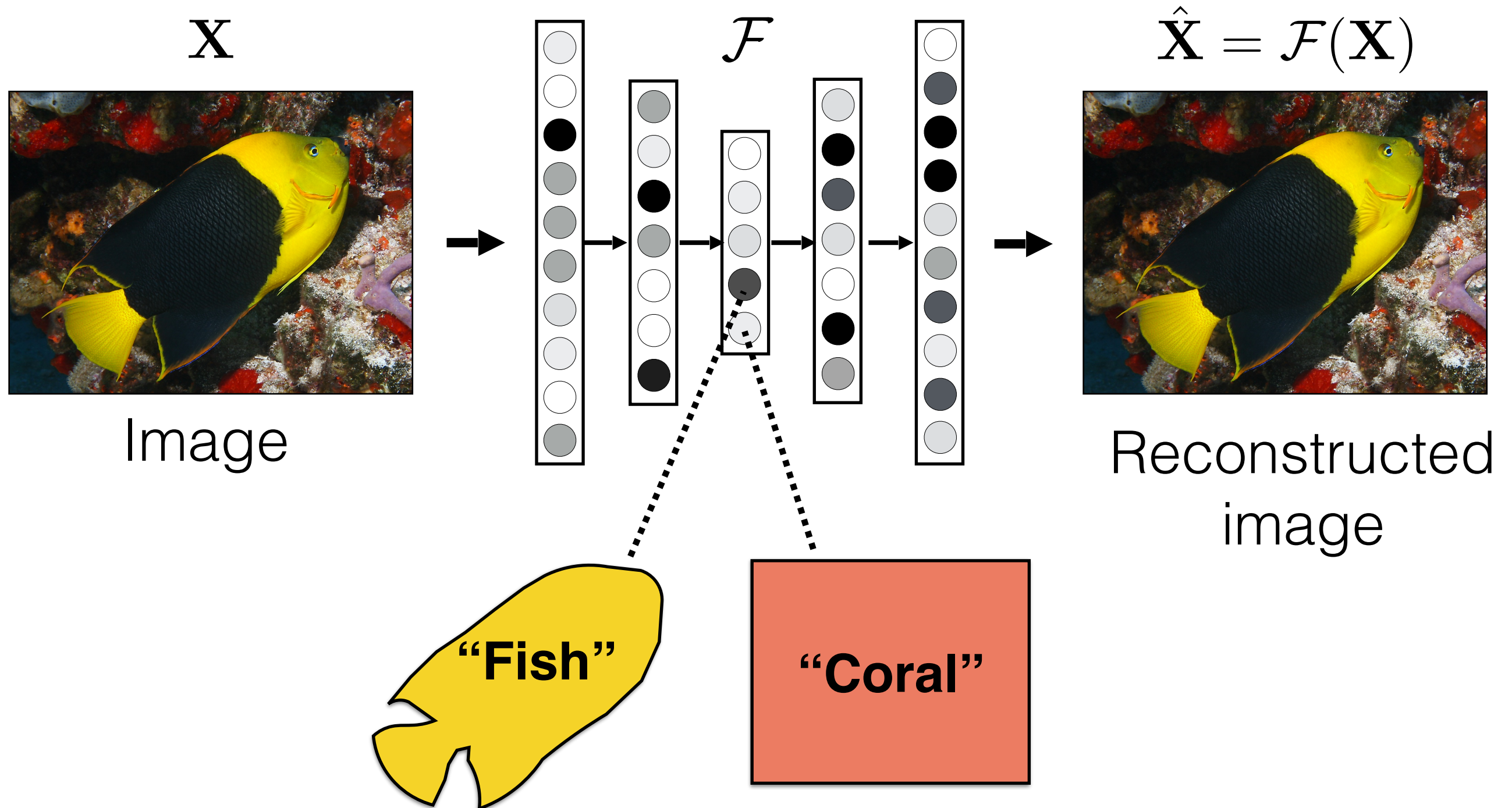


[e.g., Hinton & Salakhutdinov, Science 2006]

Autoencoder

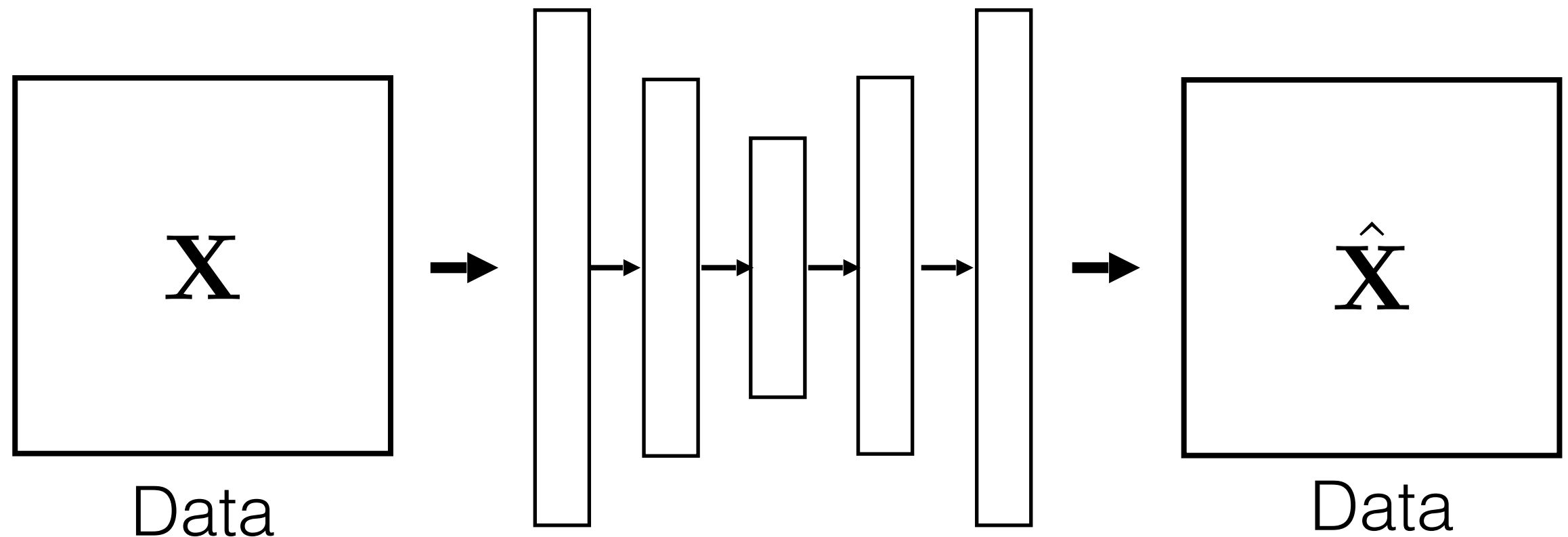


[e.g., Hinton & Salakhutdinov, Science 2006]



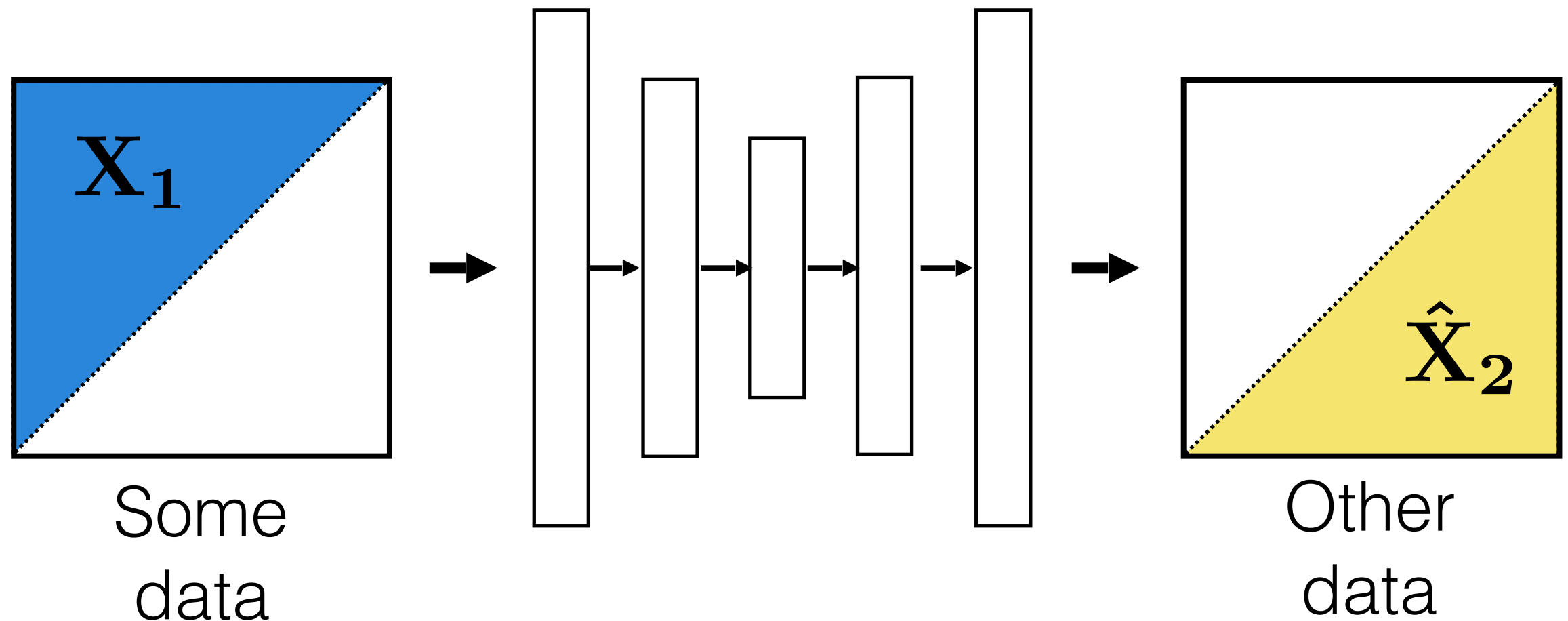
[e.g., Hinton & Salakhutdinov, Science 2006]

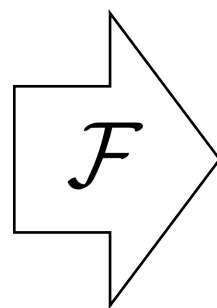
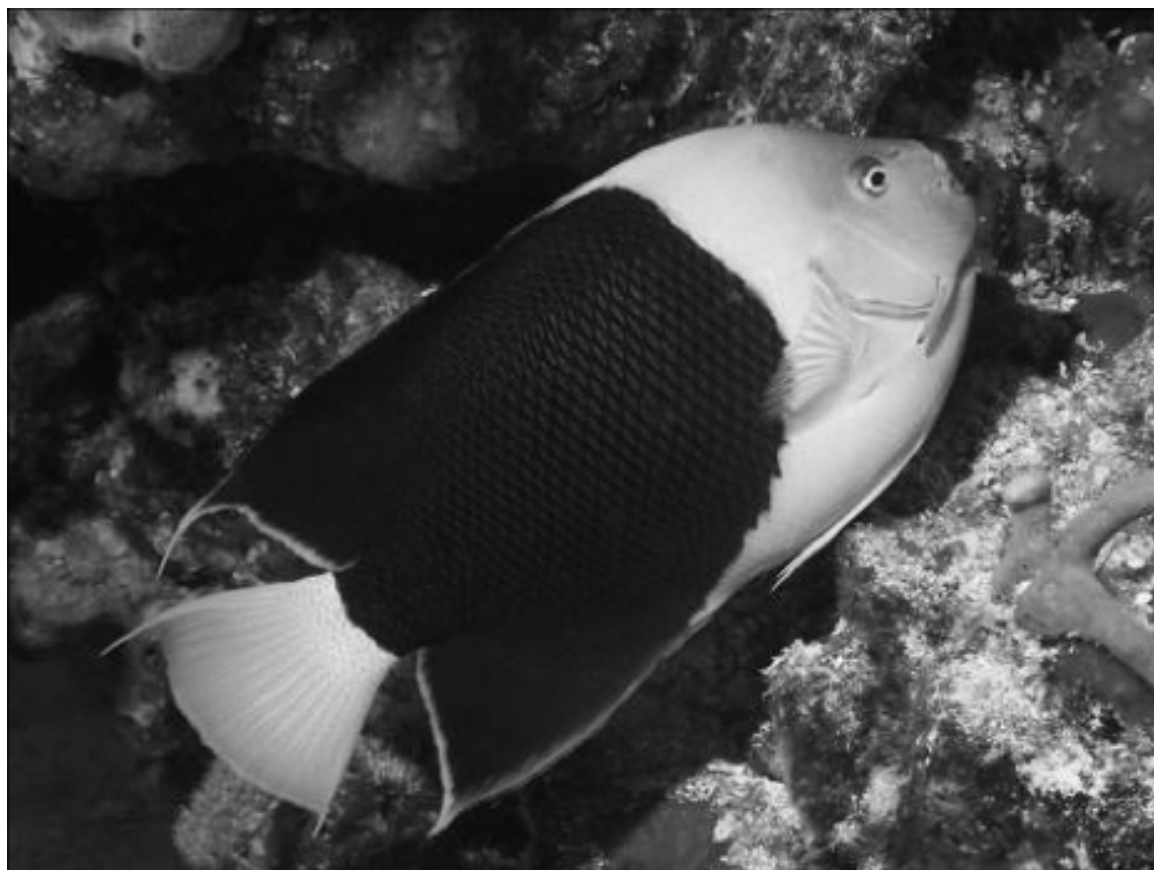
Data compression



[e.g., Hinton & Salakhutdinov, Science 2006]

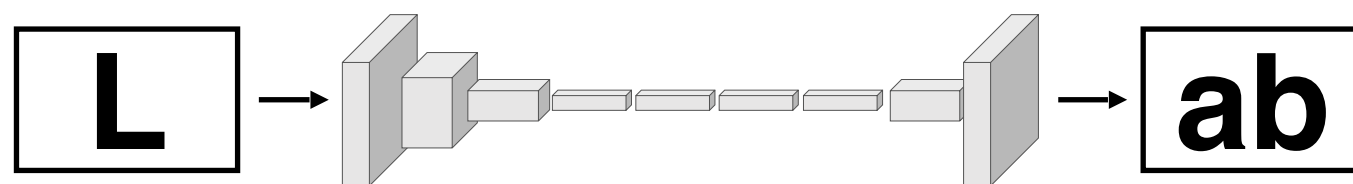
Data prediction

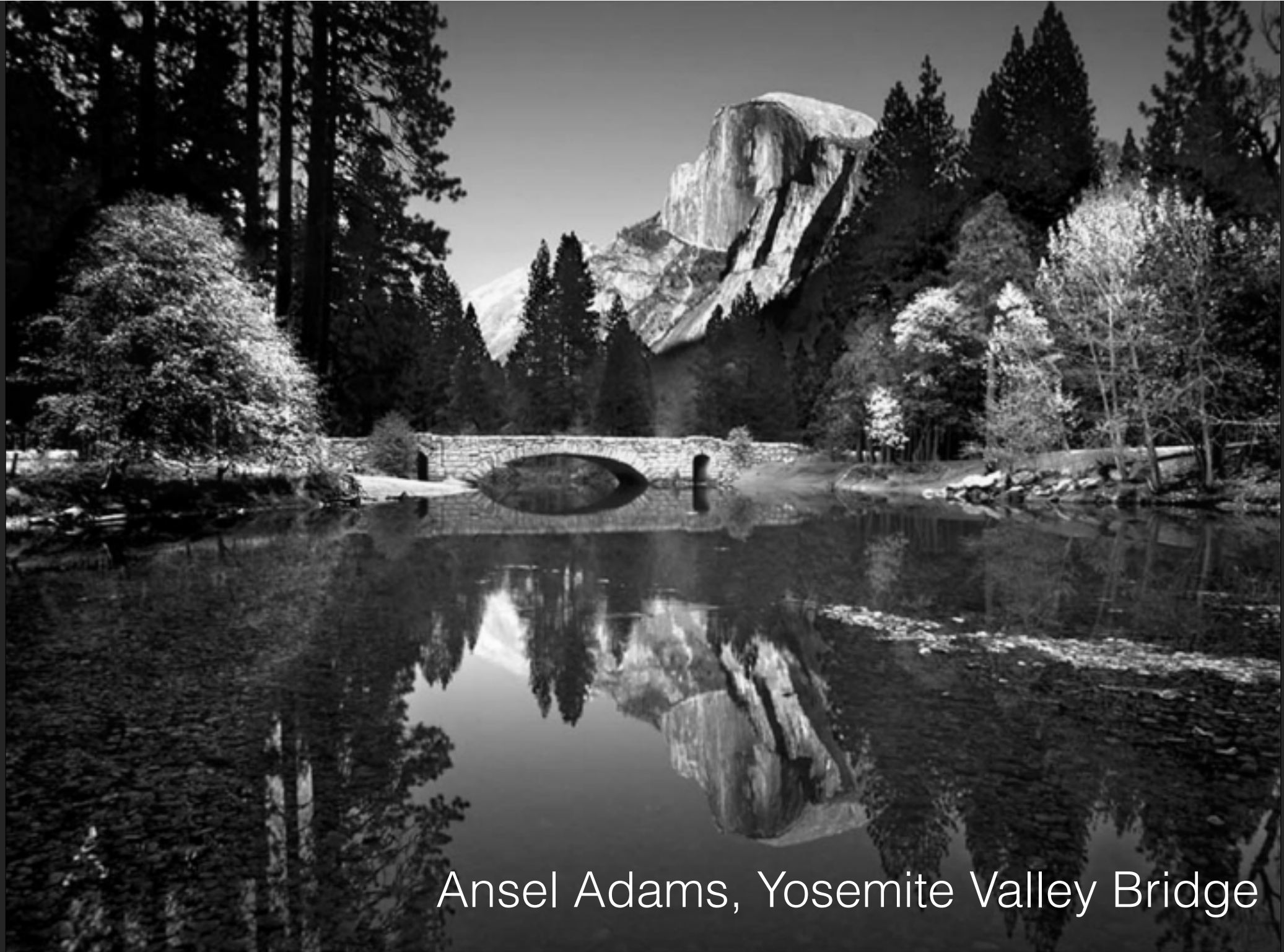




Grayscale image: L channel
 $\mathbf{X} \in \mathbb{R}^{H \times W \times 1}$

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





Ansel Adams, Yosemite Valley Bridge

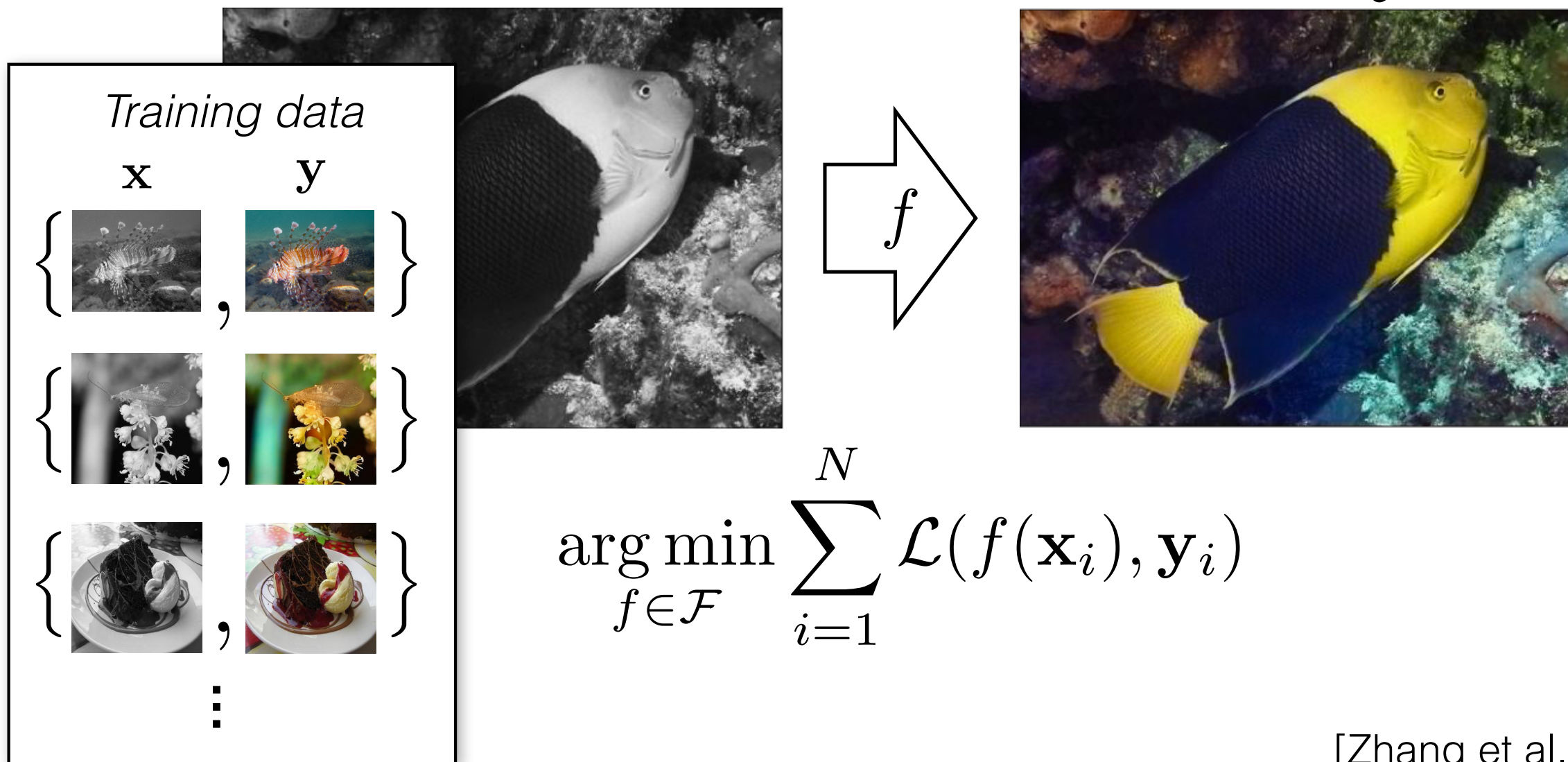


Result of [Zhang et al., ECCV 2016]

Image colorization

Input \mathbf{x}

Output \mathbf{y}



[Zhang et al., ECCV 2016]

Choosing loss and representation

Input



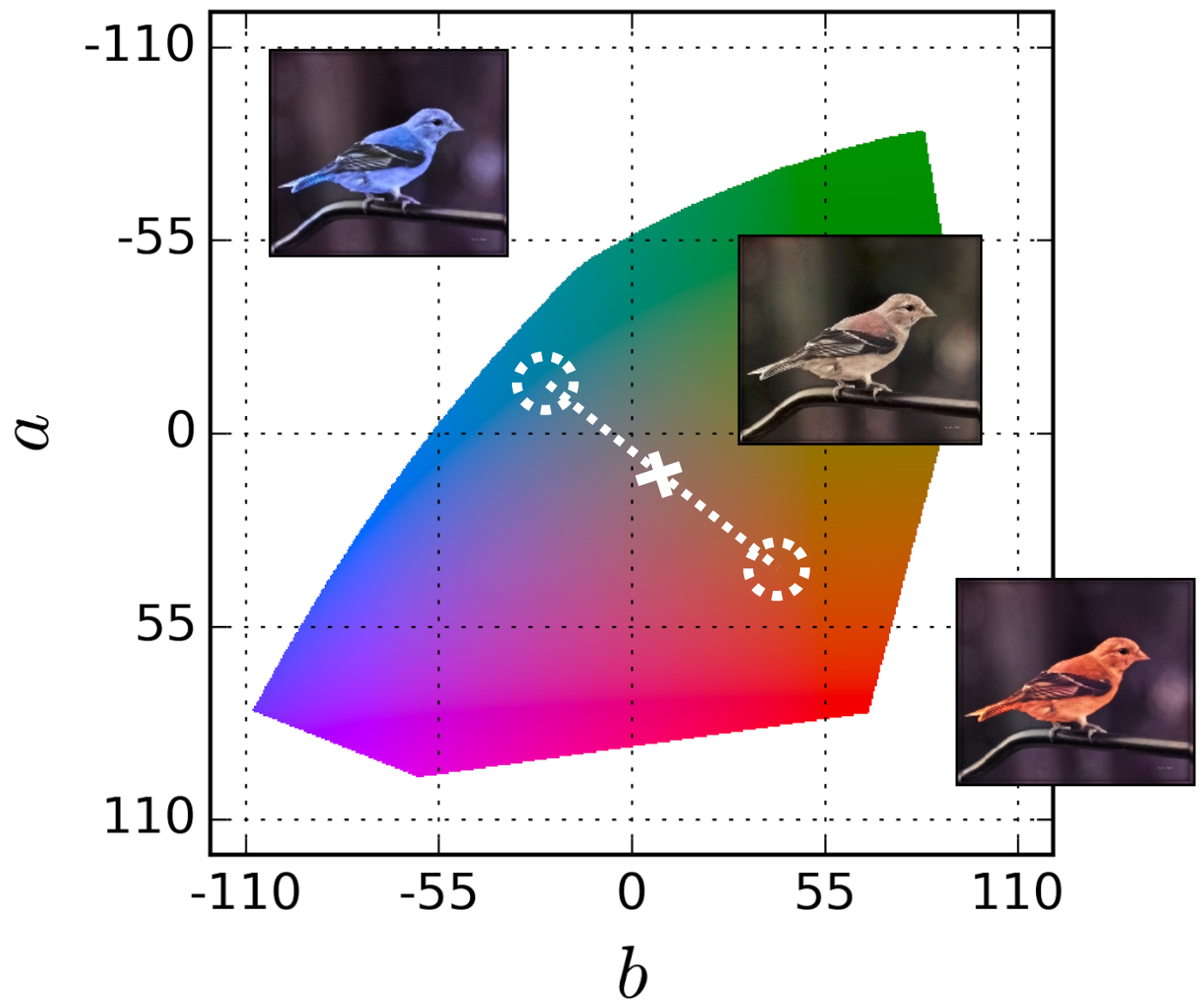
Output



Ground truth

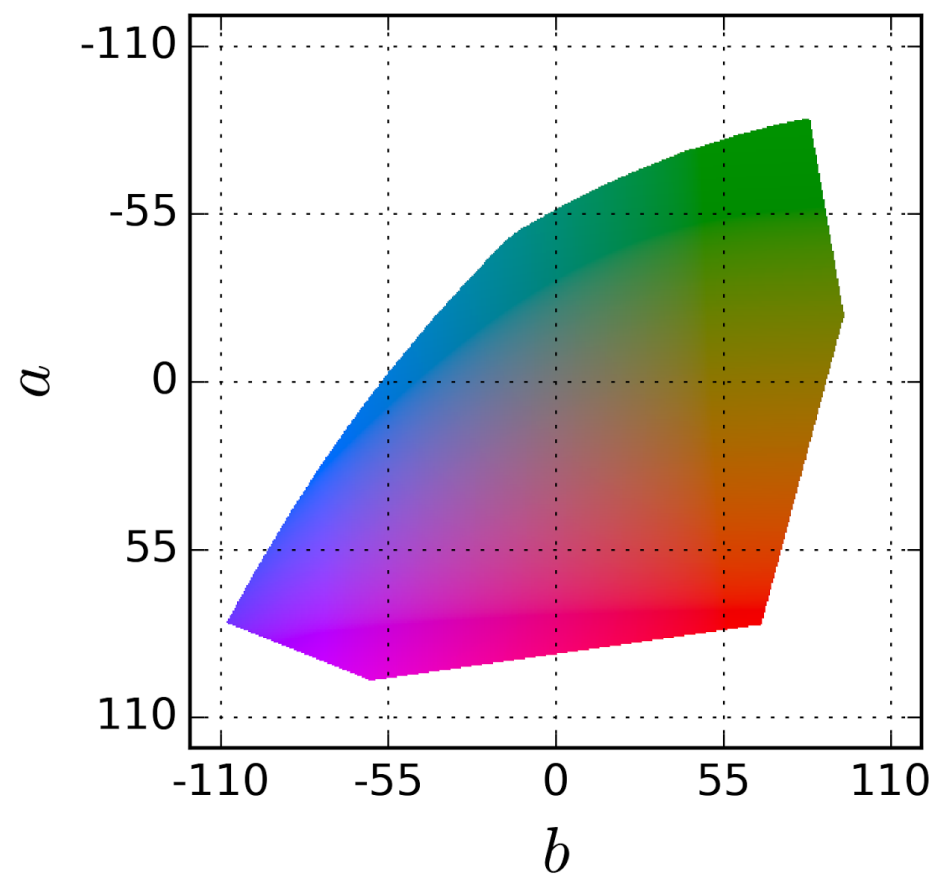


$$\mathcal{L}(f(\mathbf{x}), \mathbf{y}) = \|f(\mathbf{x}) - \mathbf{y}\|_2^2$$

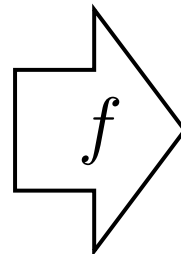


$$\mathcal{L}(f(\mathbf{x}), \mathbf{y}) = \|f(\mathbf{x}) - \mathbf{y}\|_2^2$$

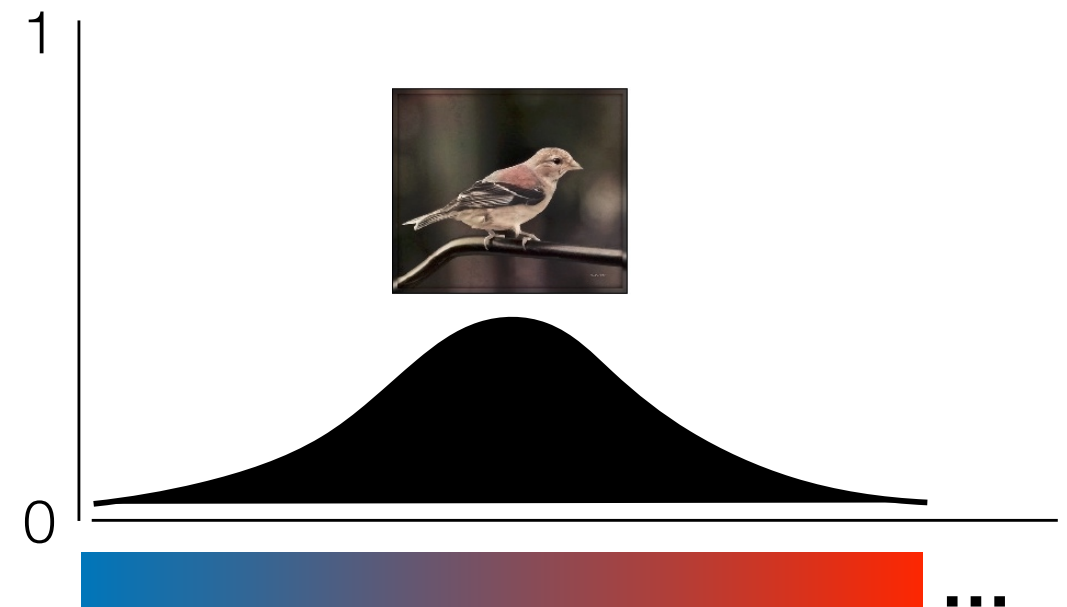
$$\mathbf{y} \in \mathbb{R}^{H \times W \times 2}$$



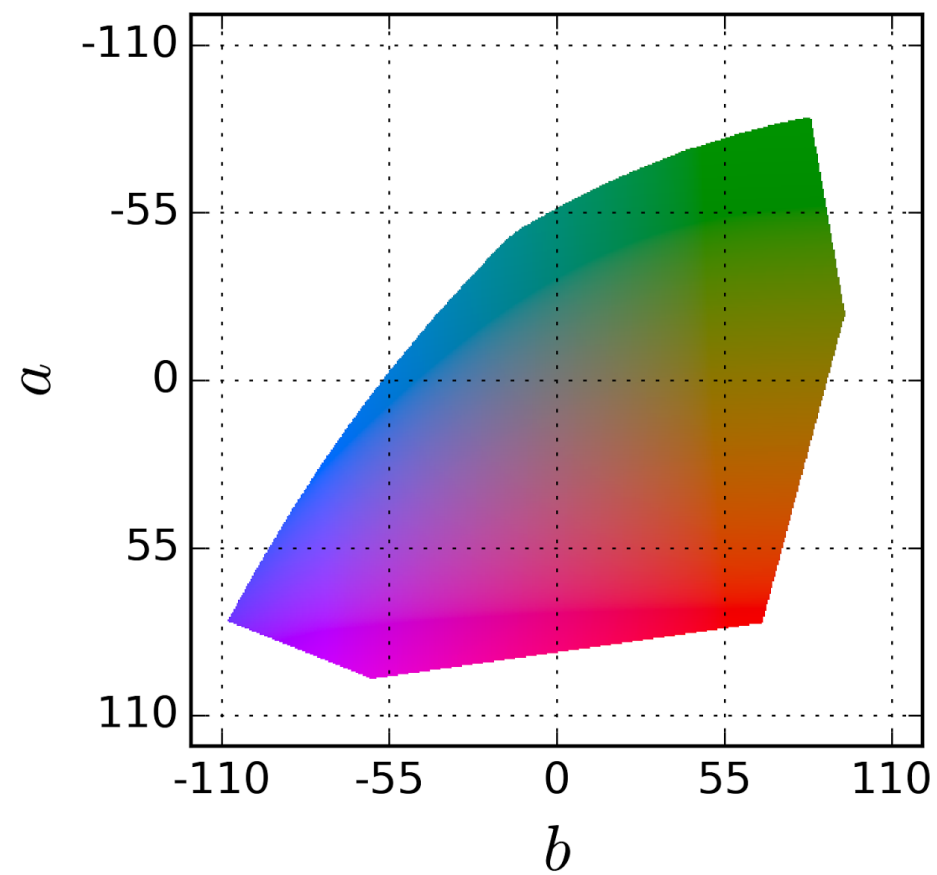
$$\mathcal{L}(f(\mathbf{x}), \mathbf{y}) = \|f(\mathbf{x}) - \mathbf{y}\|_2^2$$



Prediction for a single pixel i, j

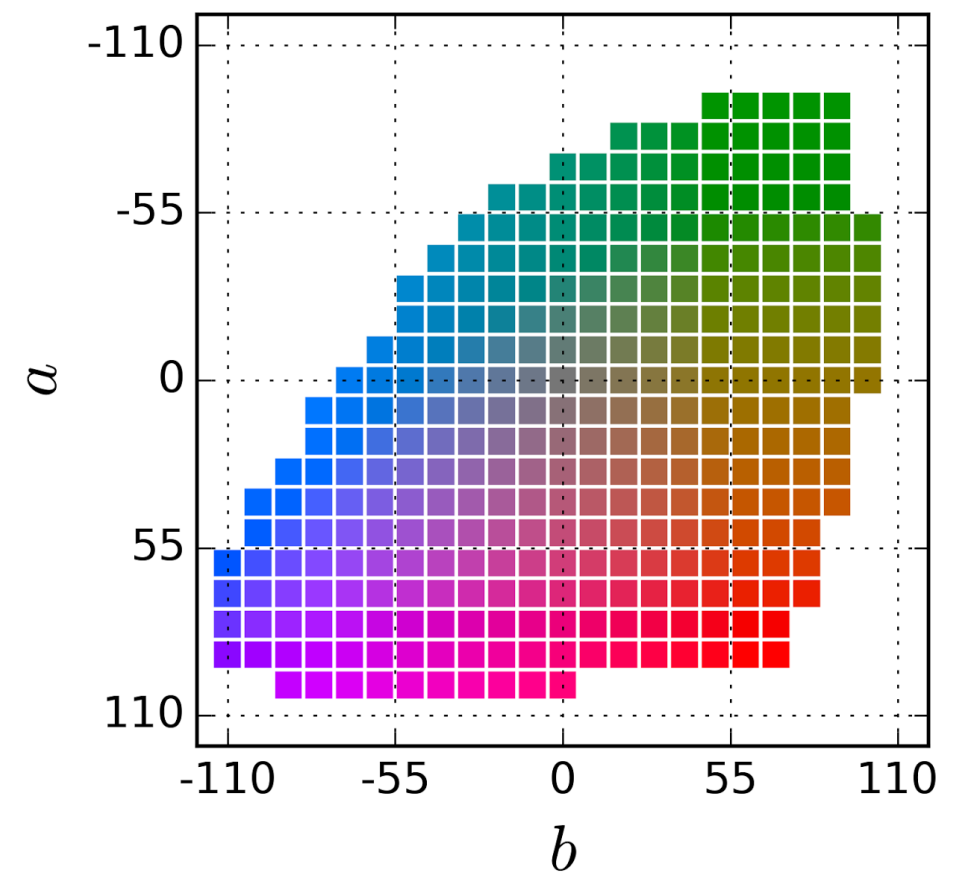


$$\mathbf{y} \in \mathbb{R}^{H \times W \times 2}$$

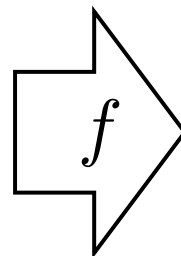
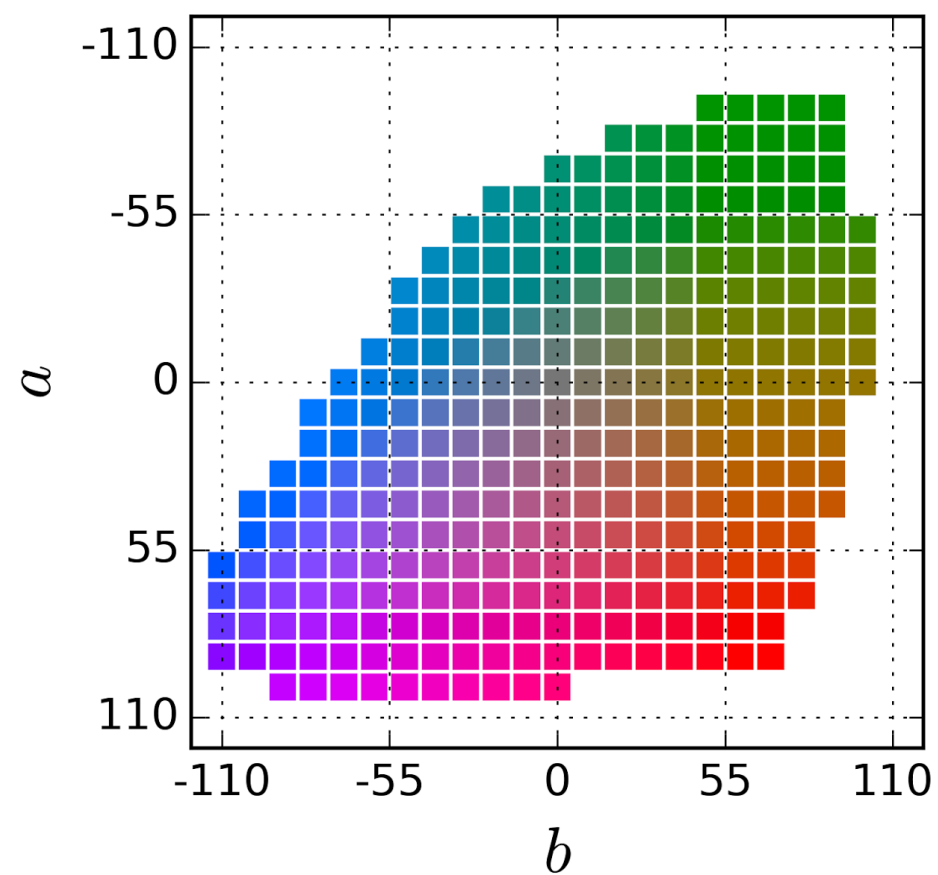


one-hot representation of K discrete classes

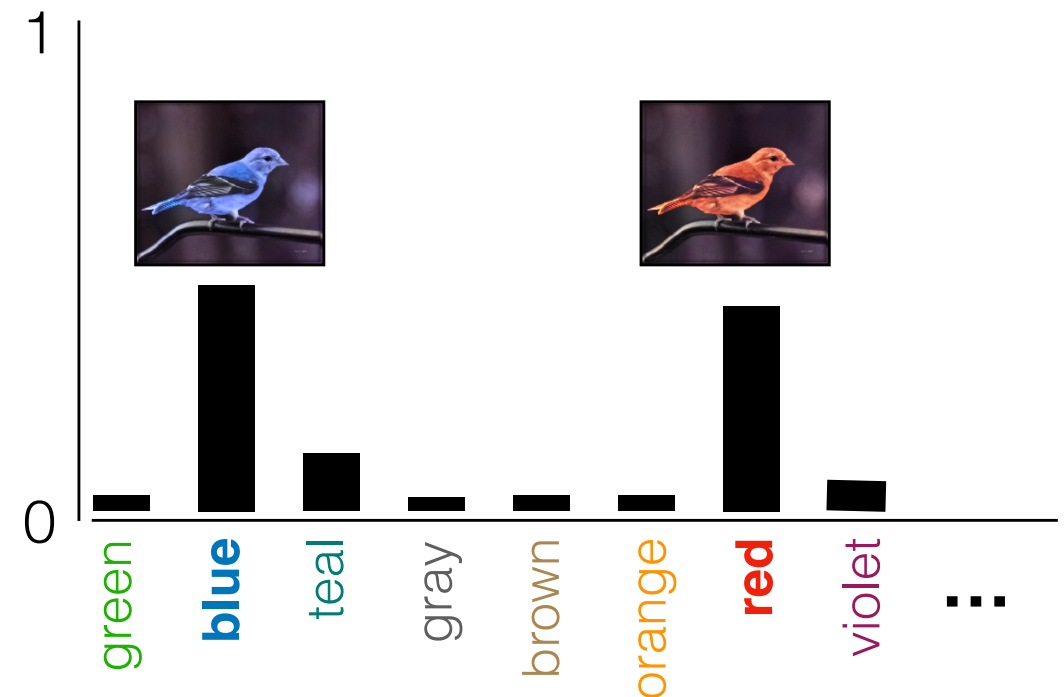
$$\mathbf{y} \in \mathbb{R}^{H \times W \times K}$$



$$\mathbf{y} \in \mathbb{R}^{H \times W \times K}$$



Prediction for a single pixel i, j



$$\mathcal{L}(\mathbf{y}, f_{\theta}(\mathbf{x})) = H(\mathbf{y}, \text{softmax}(f_{\theta}(\mathbf{x})))$$

Input



Zhang et al. 2016

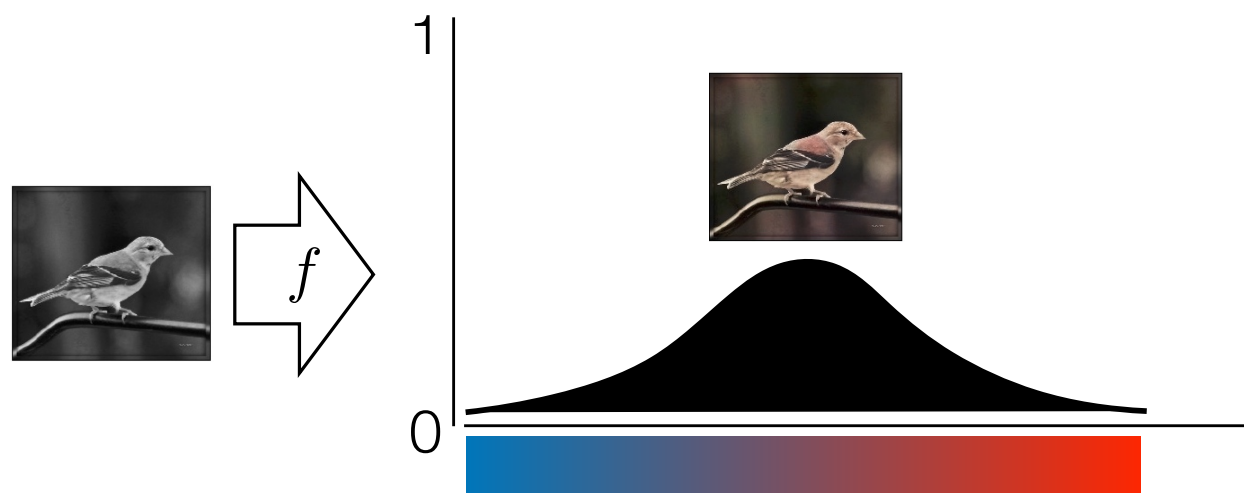


Ground truth



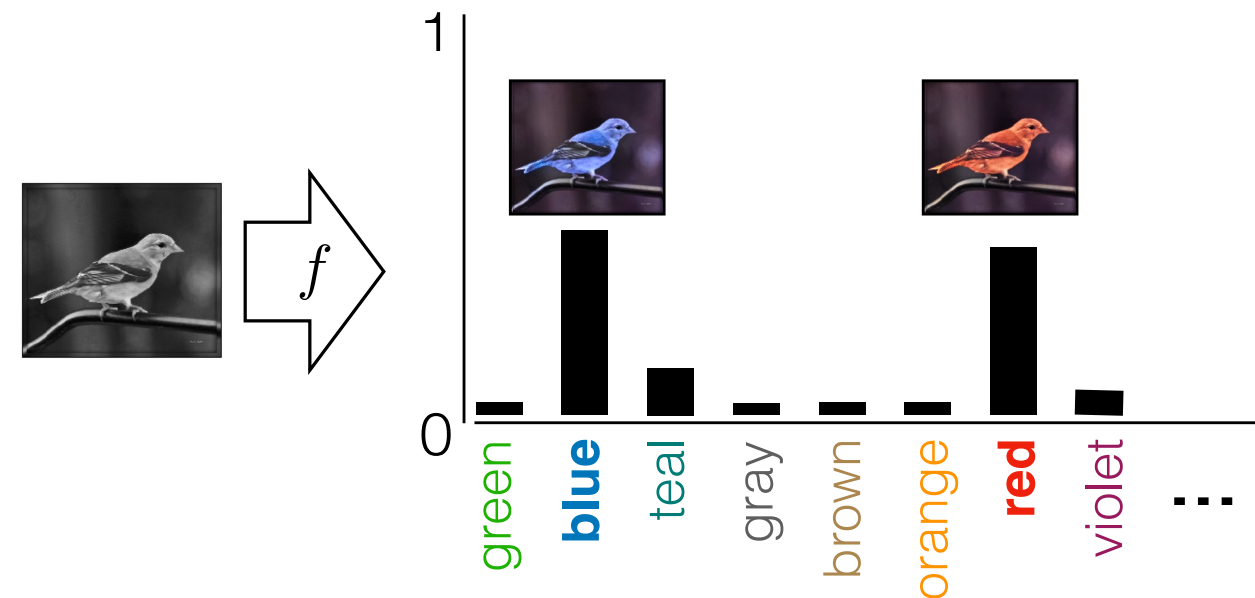
$$\mathcal{L}(\mathbf{x}, \mathbf{y}) = H(\mathbf{y}, \text{softmax}(f_{\theta}(\mathbf{x})))$$

“Regression”



- Continuous-valued prediction
- (Usually) models unimodal distribution

“Classification”



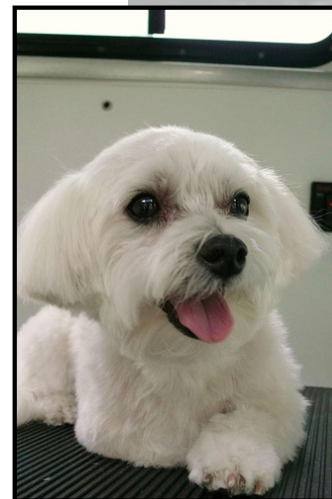
- Discrete-valued prediction
- Models multimodal distribution



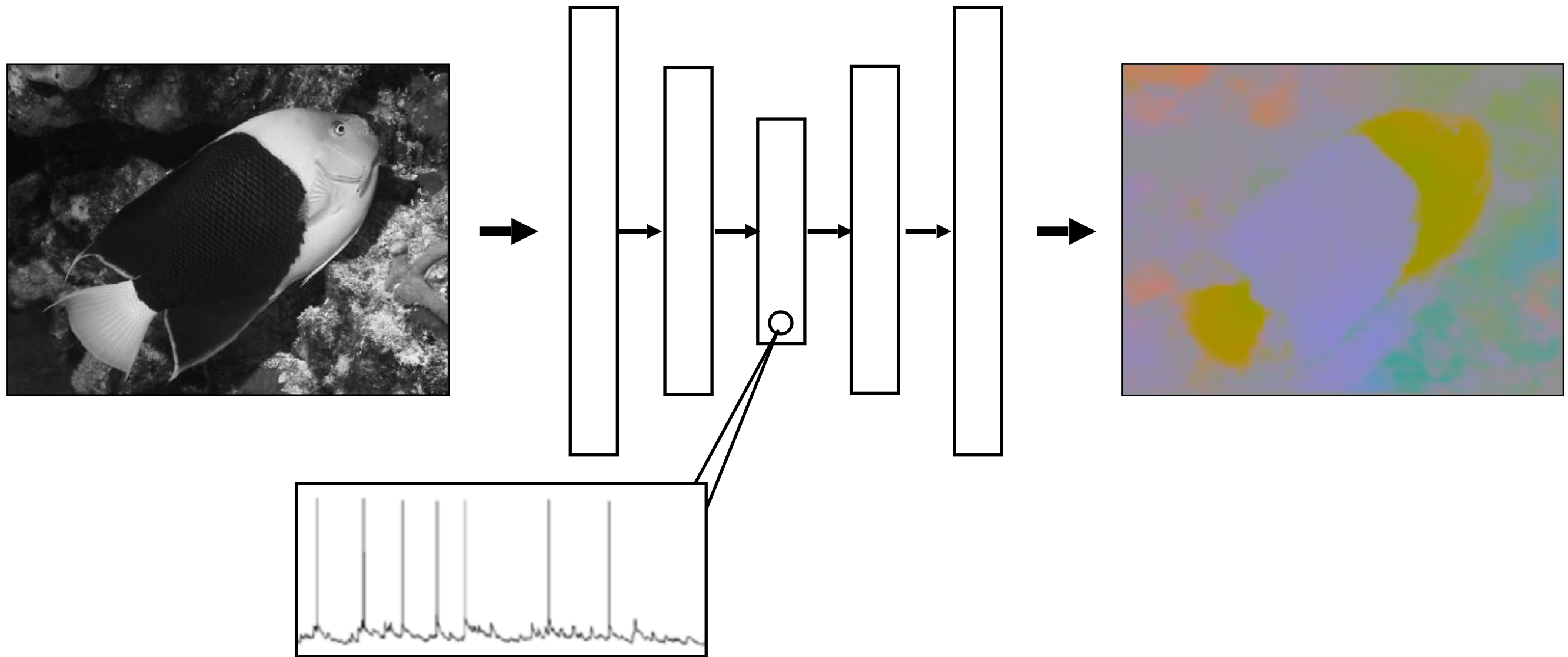
Instructive failure



Instructive failure



Deep Net “Electrophysiology”

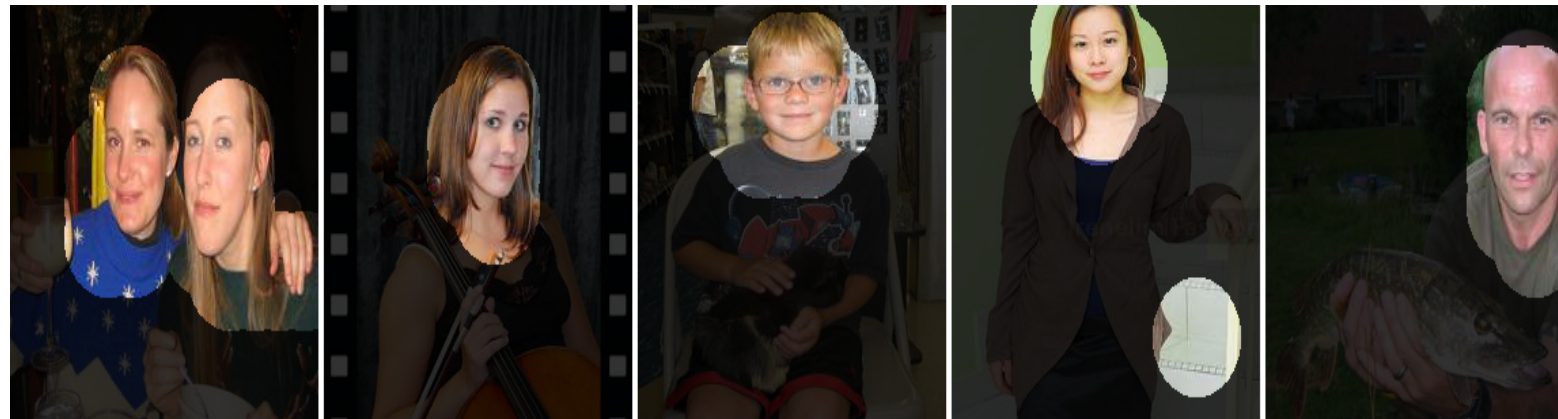


[Zeiler & Fergus, ECCV 2014]

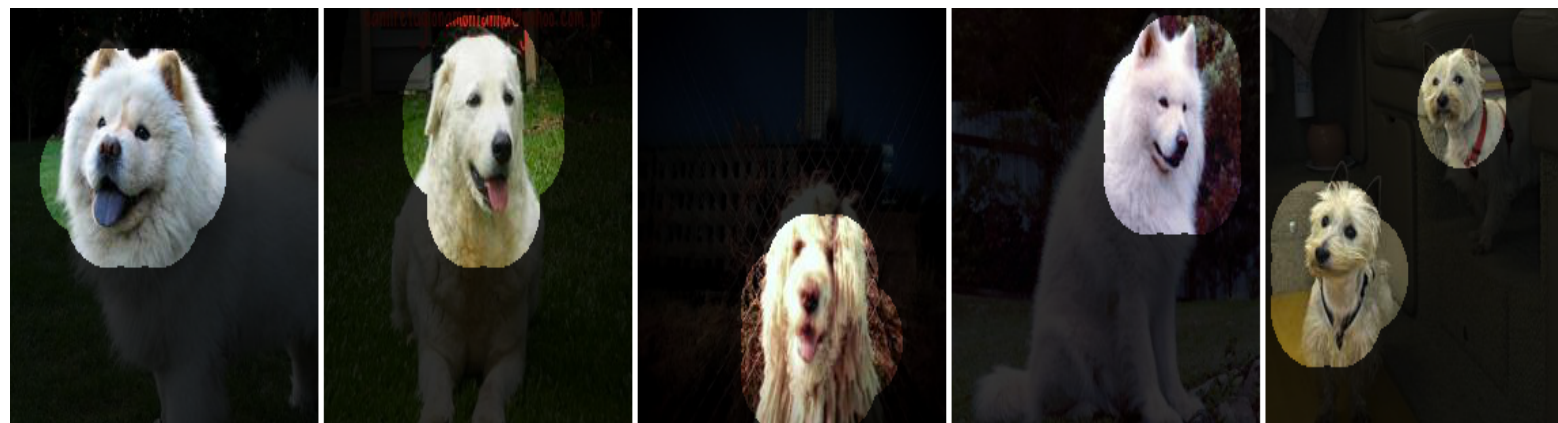
[Zhou et al., ICLR 2015]

Stimuli that drive selected neurons (conv5 layer)

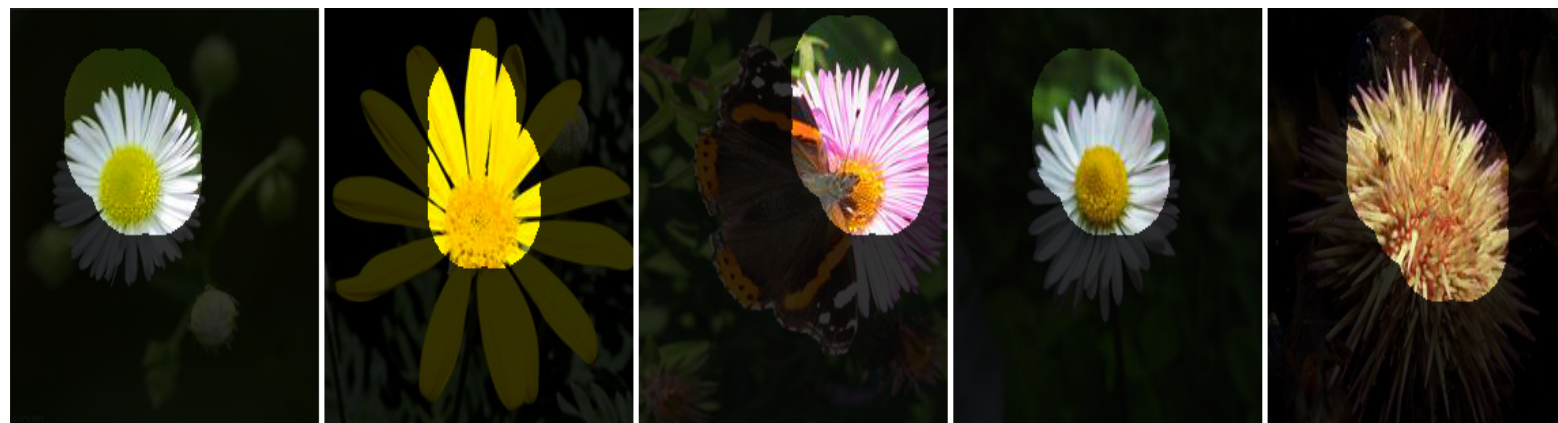
faces

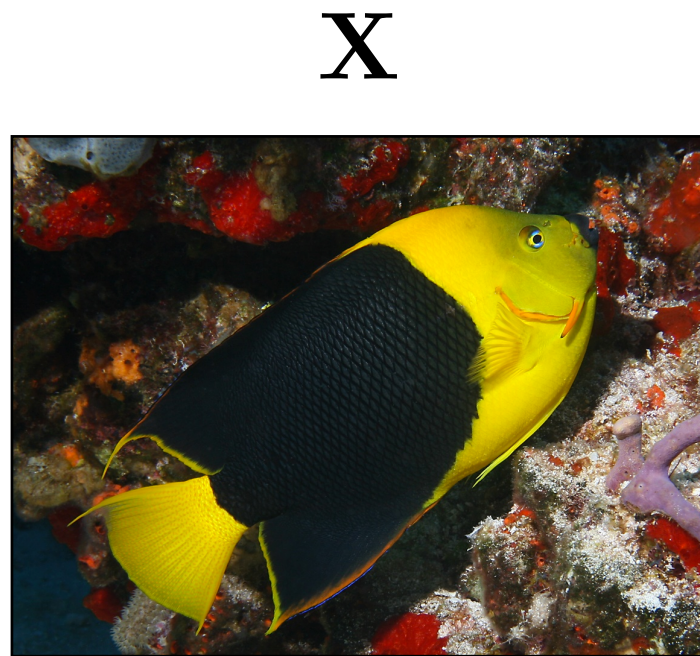


dog
faces

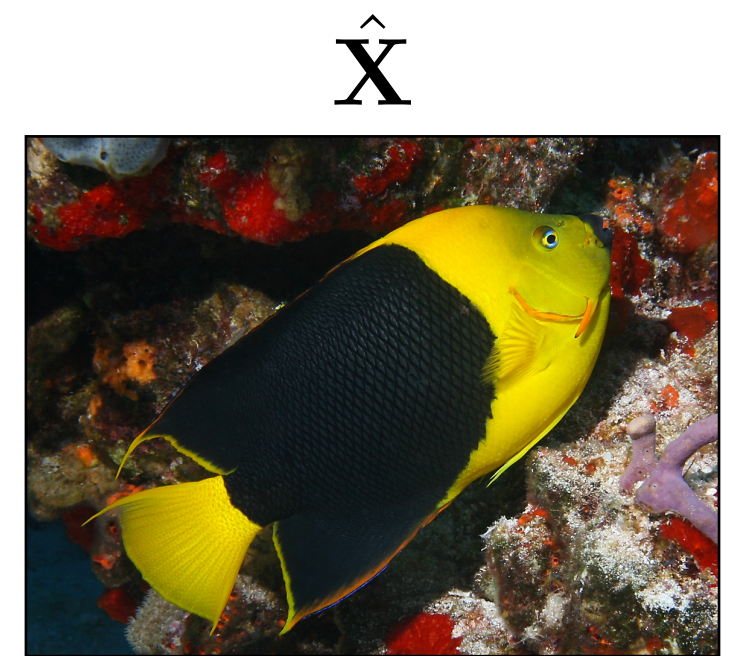
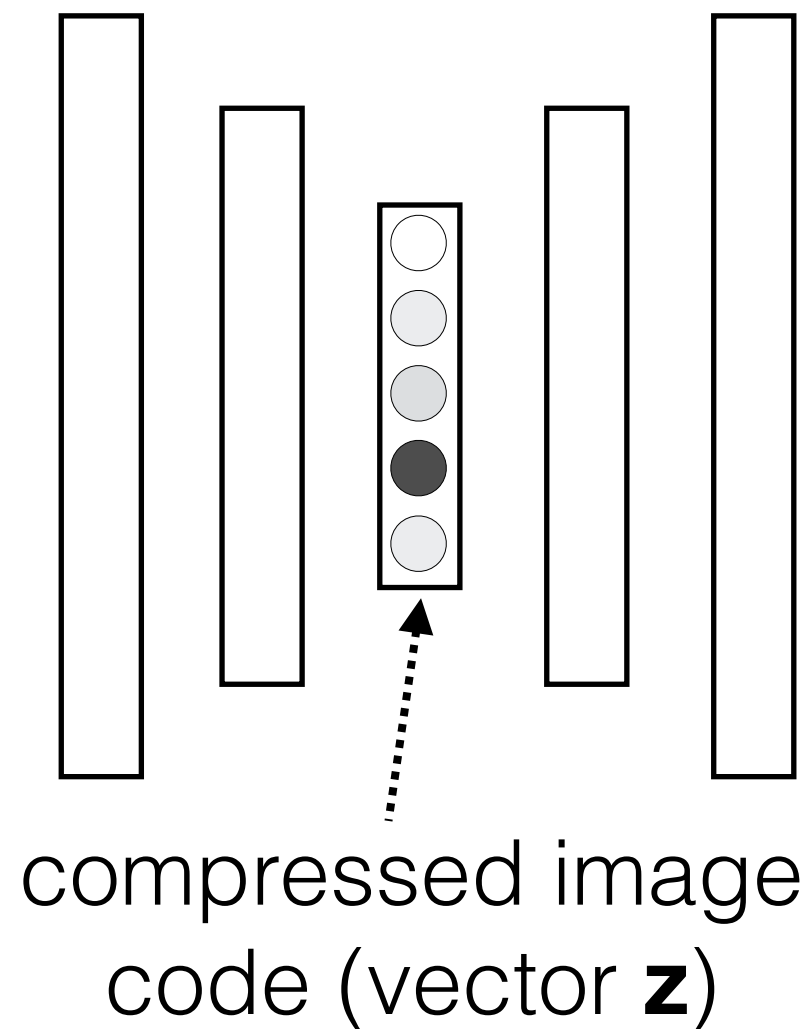


flowers





Image

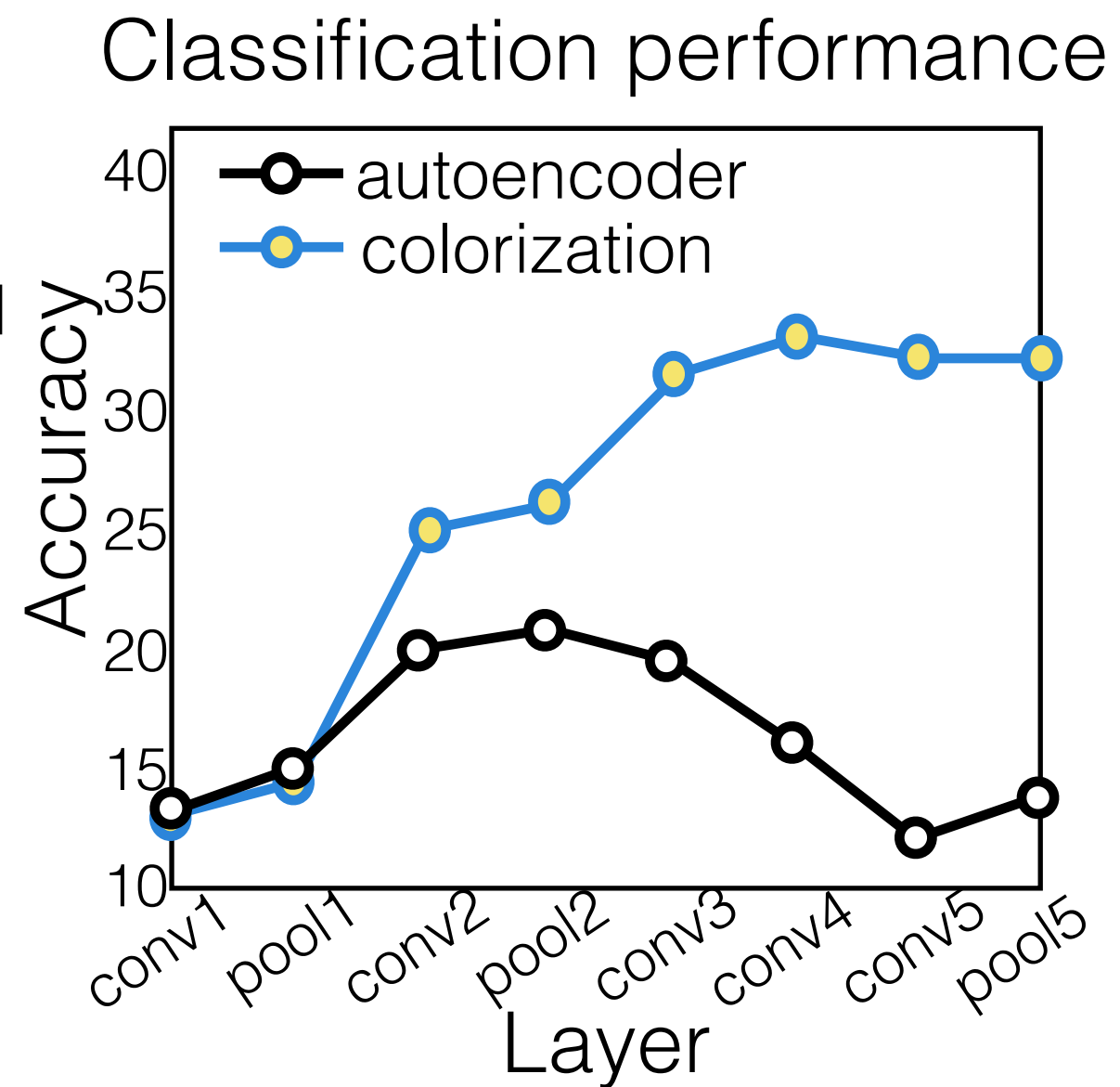
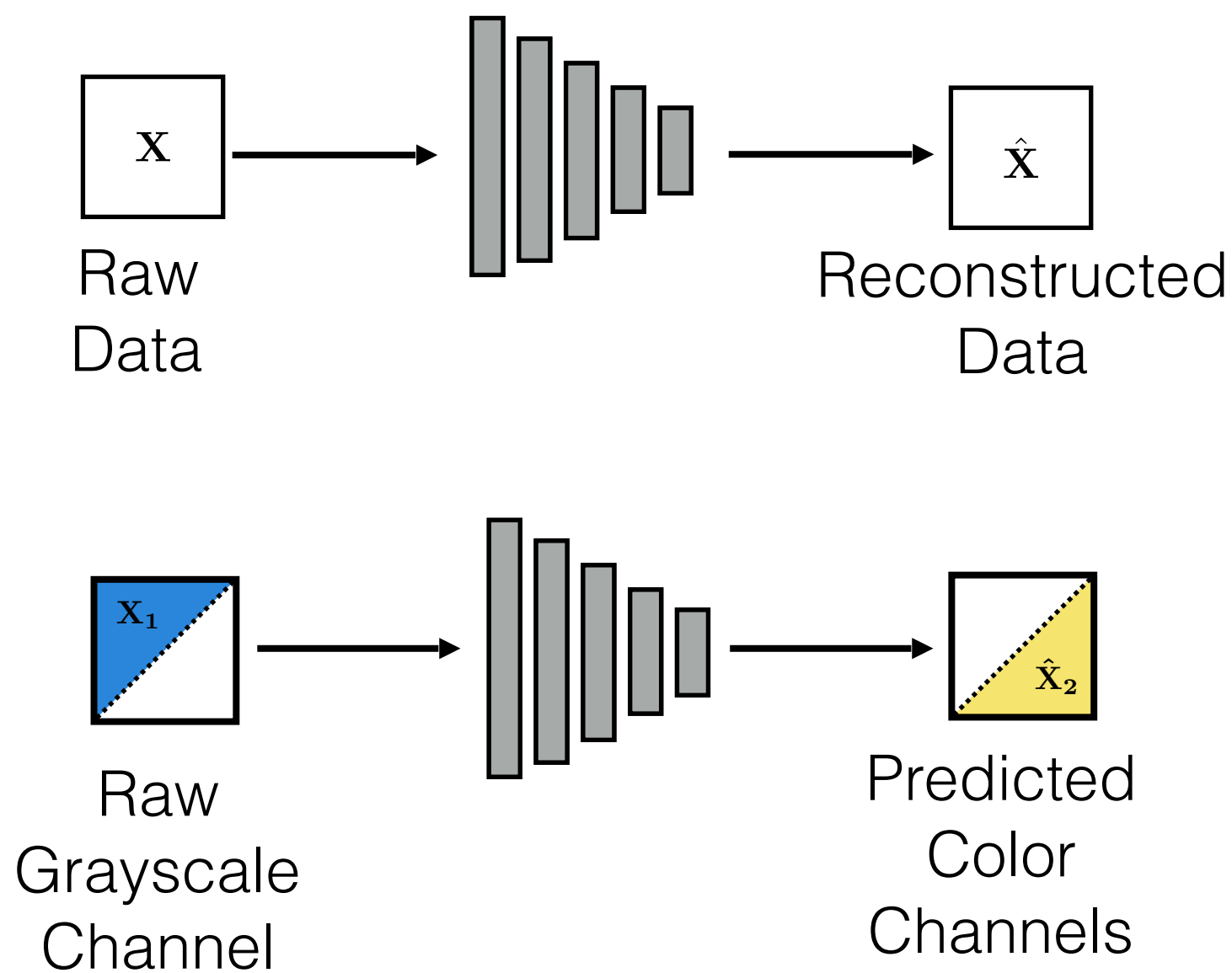


Reconstructed image

Is the code informative about object class y ?

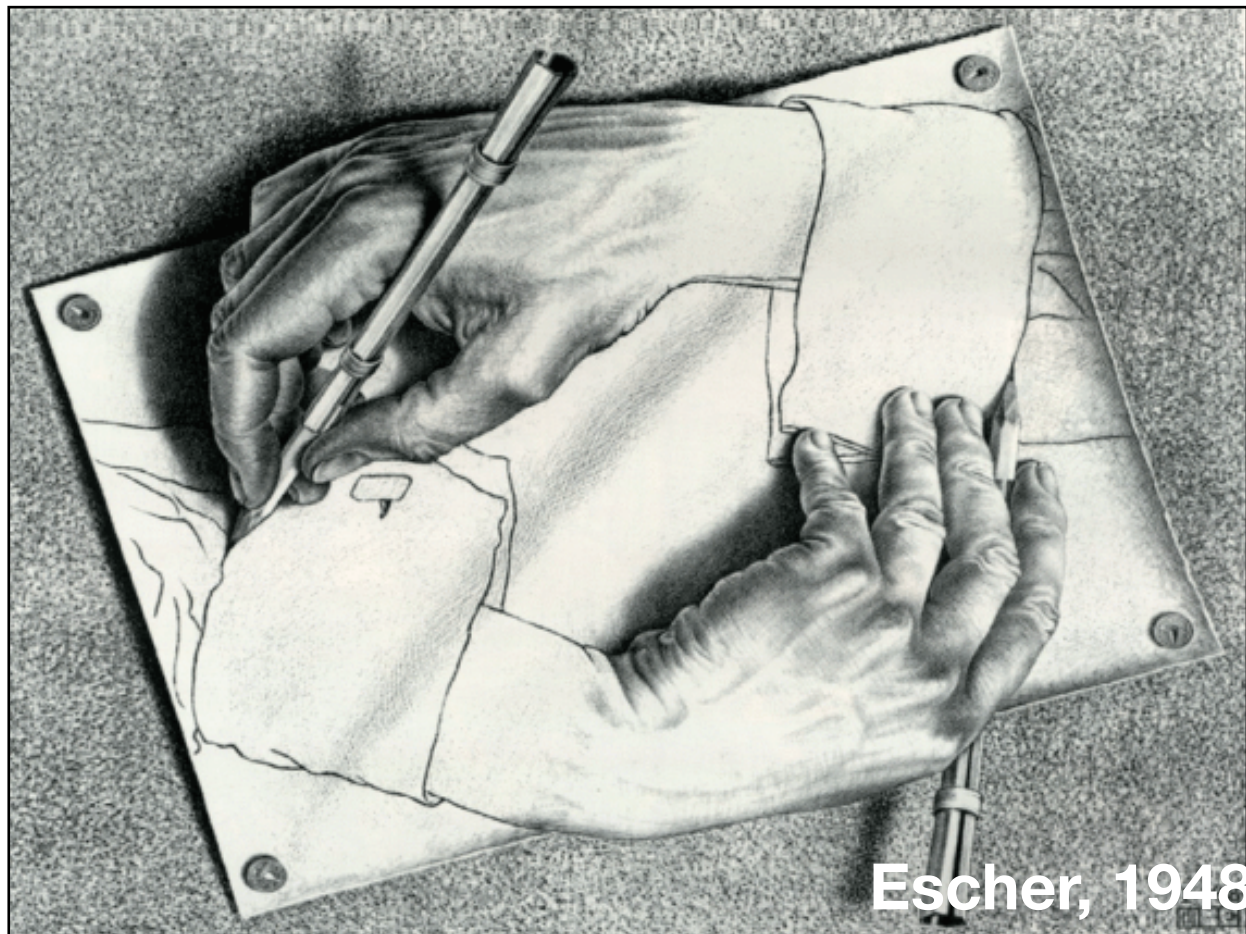
Logistic regression:

$$y = \sigma(\mathbf{W}\mathbf{z} + \mathbf{b})$$



Task from [Russakovsky et al. 2015]

Self-supervised learning



Common trick:

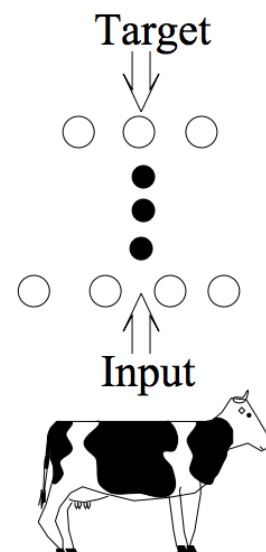
- Convert “unsupervised” problem into “supervised” empirical risk minimization
- Do so by cooking up “labels” (prediction targets) from the raw data itself

Multisensory self-supervision

Supervised

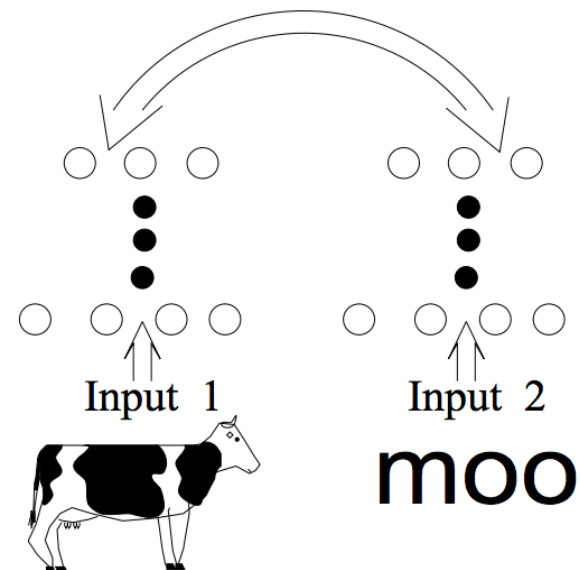
- implausible label

"COW"



Self-Supervised

- derives label from a co-occurring input to another modality



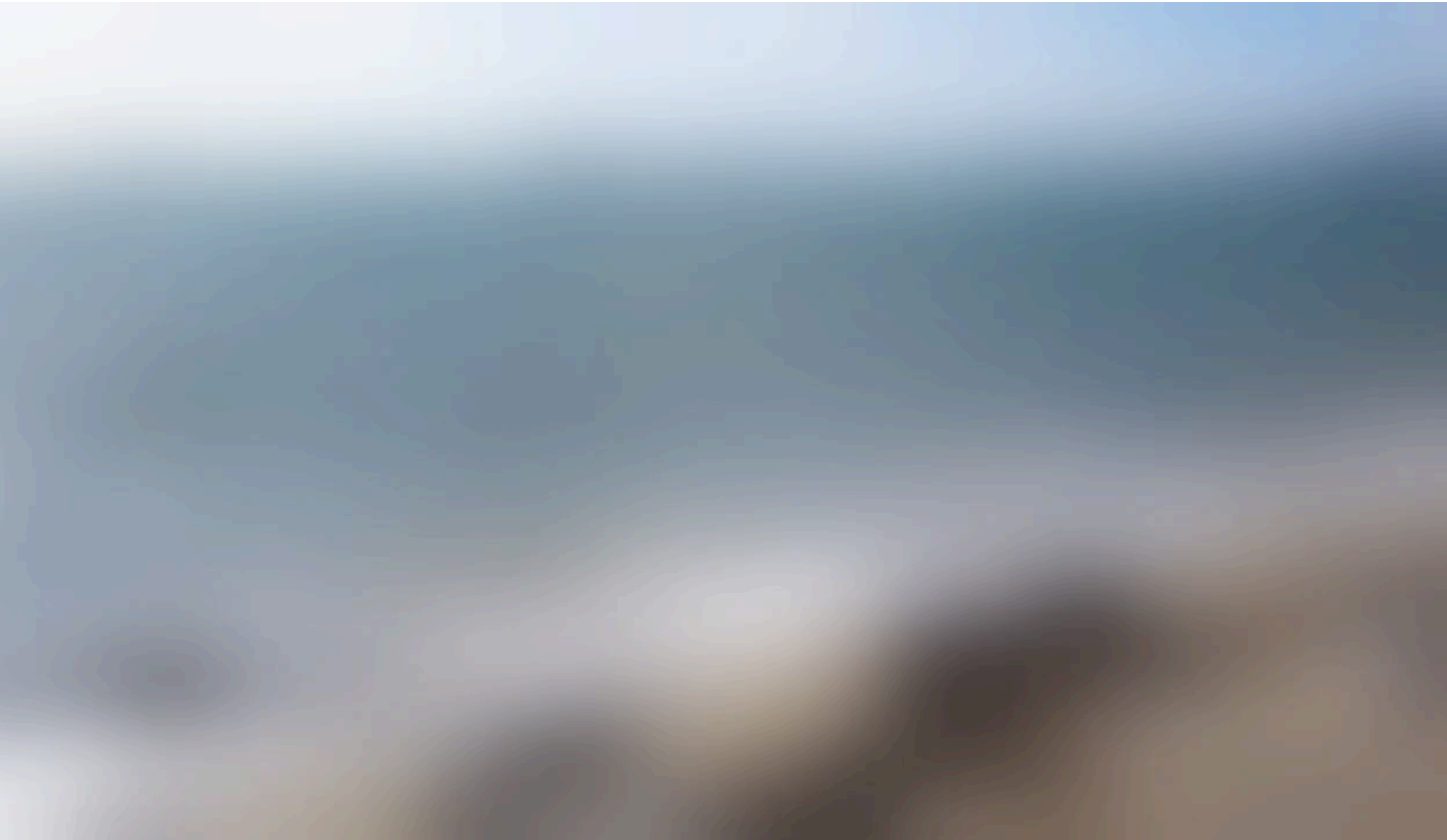
Virginia de Sa. *Learning Classification with Unlabeled Data*. NIPS 1994.

[see also "Six lessons from babies", Smith and Gasser 2005]

Ambient Sound Provides Supervision for Visual Learning

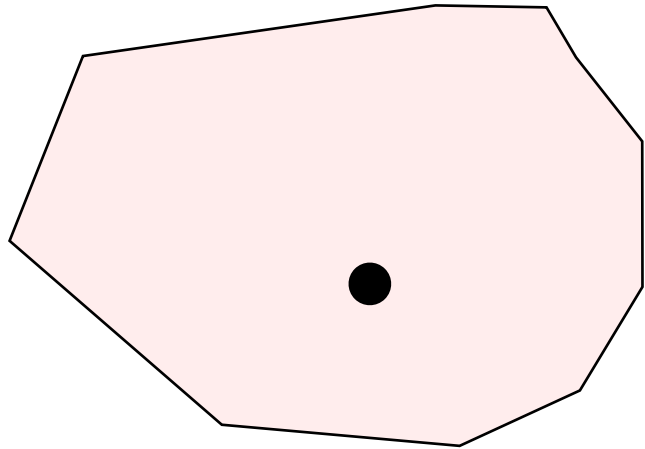
Andrew Owens Jiajun Wu Josh McDermott
William Freeman Antonio Torralba

MIT, Google Research

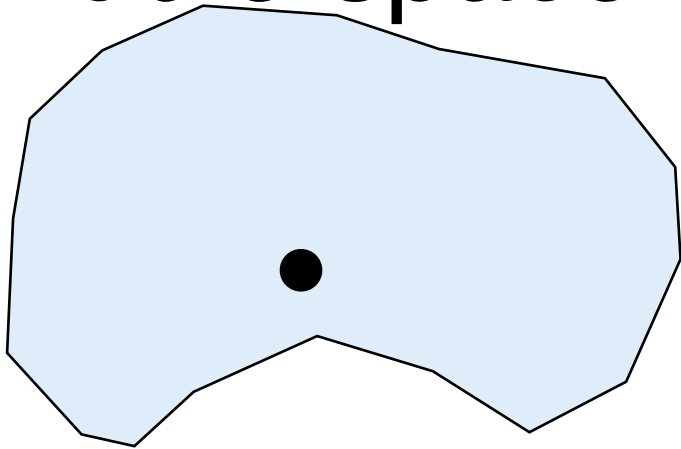


Audio is invariant to many visual transformations

Image space

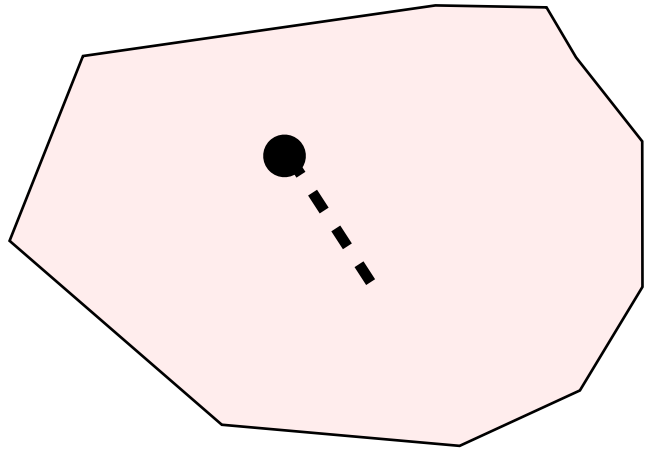


Audio space

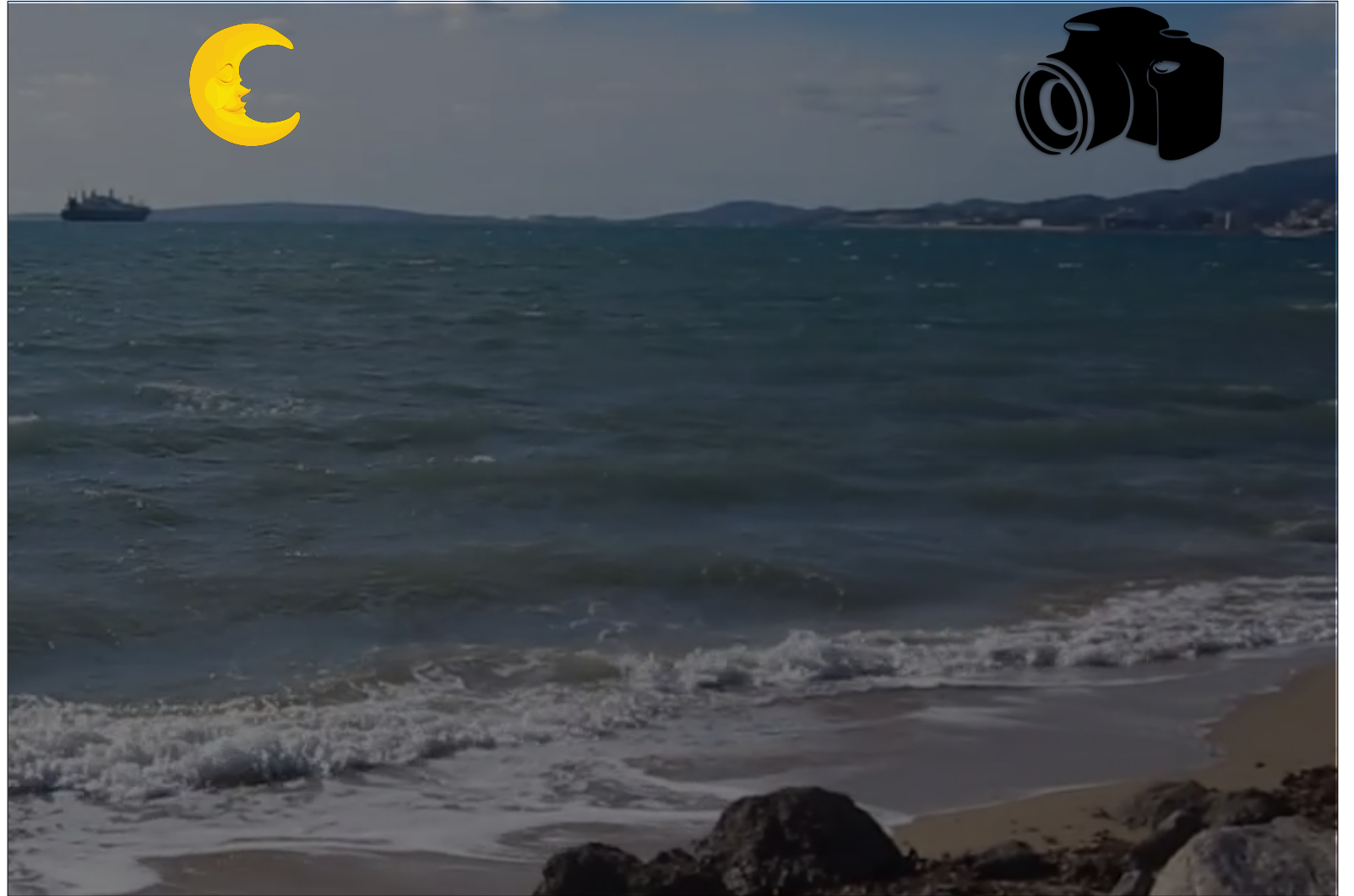
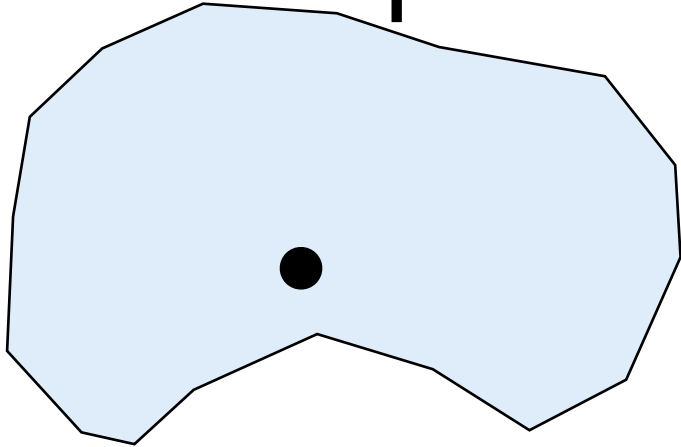


Audio is invariant to many visual transformations

Image space

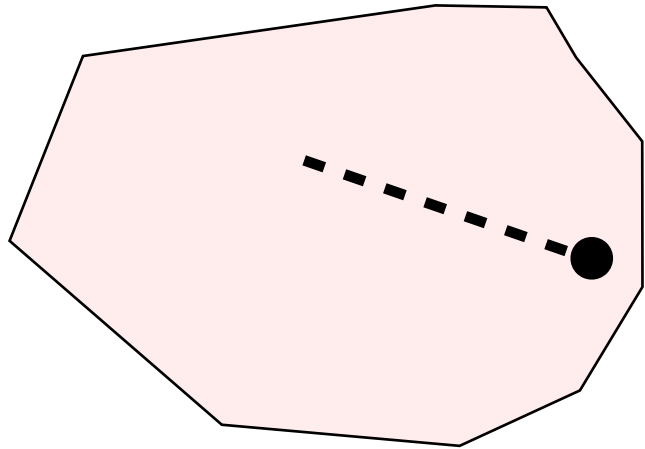


Audio space

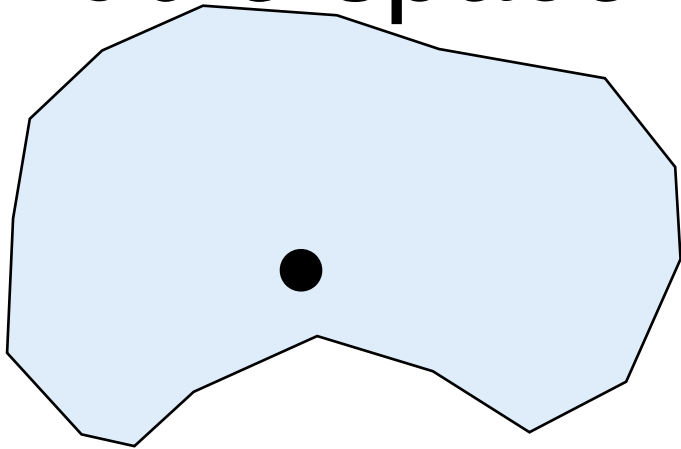


Audio is invariant to many visual transformations

Image space

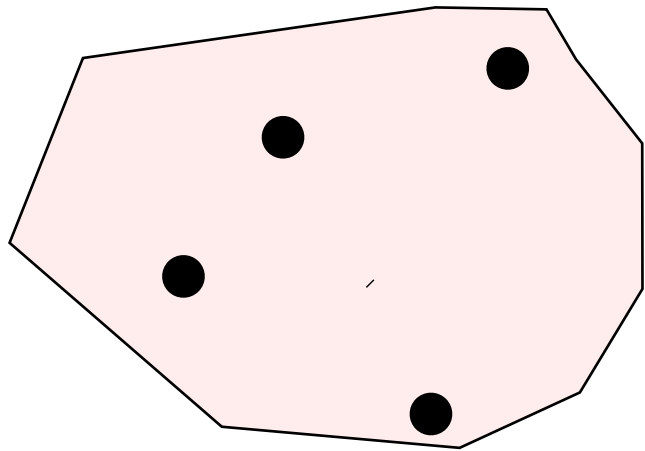


Audio space

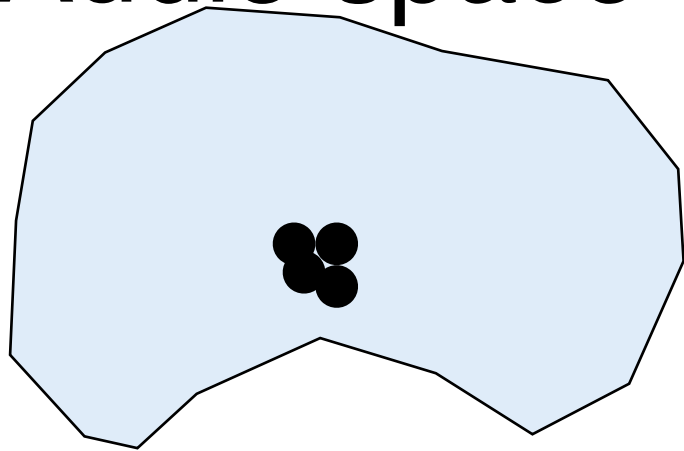


Audio is invariant to many visual transformations

Image space

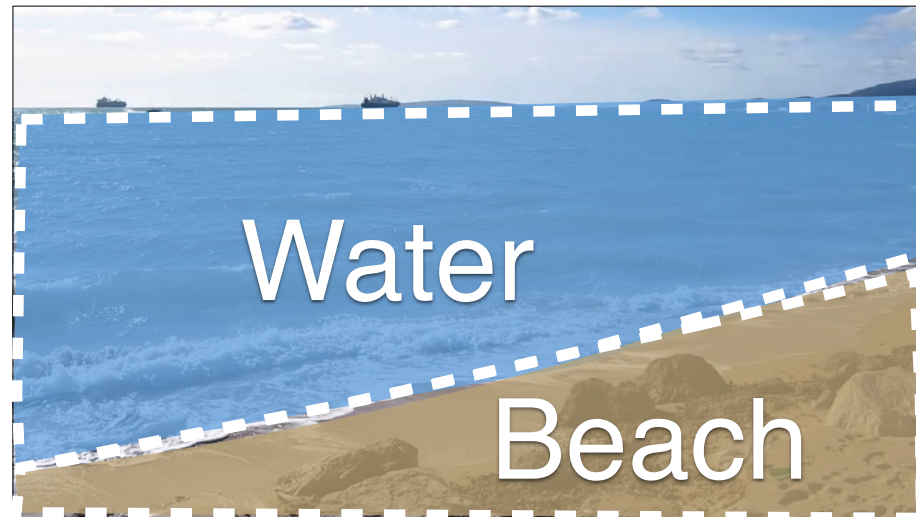
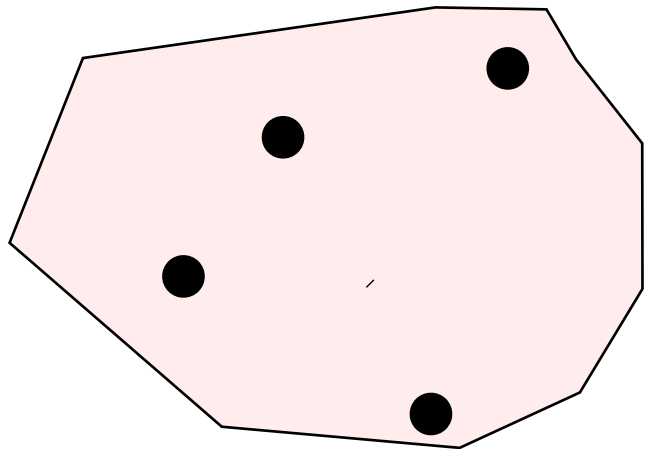


Audio space

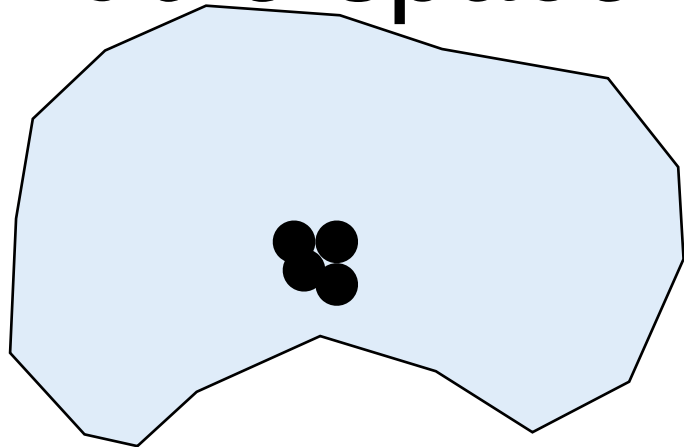


Audio is invariant to many visual transformations

Image space



Audio space

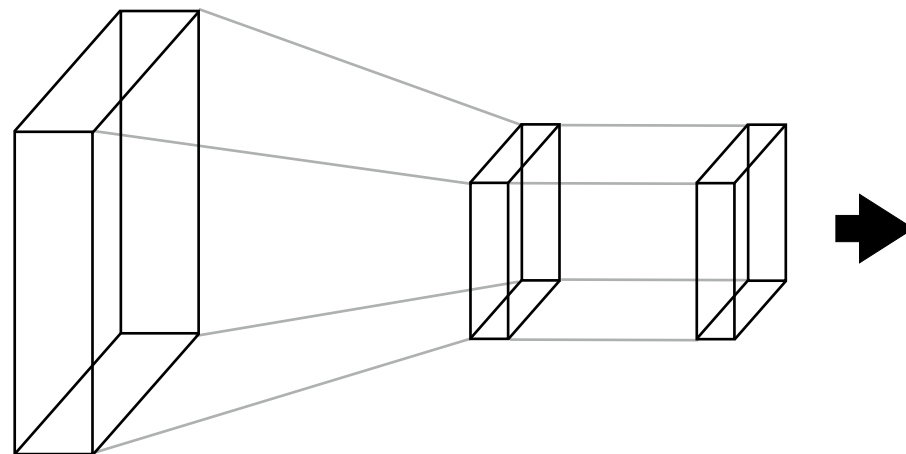


Predicting sound

- Flickr video dataset.
- 180K videos, 10 random frames from each.
- Trained from scratch



Video frame



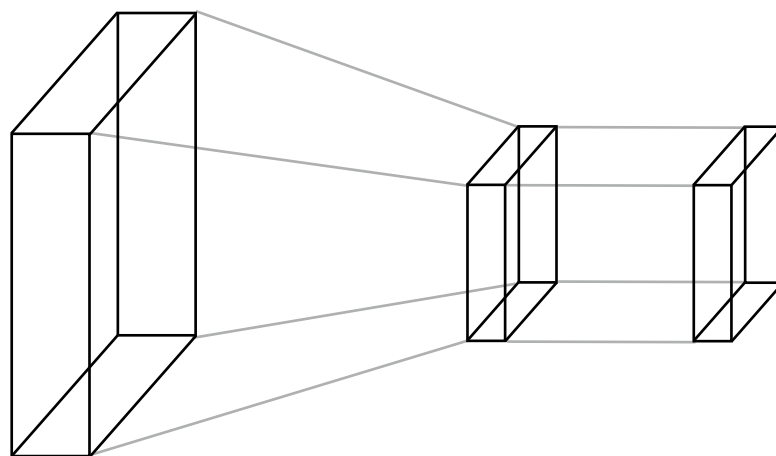
ConvNet

Sound texture
[McDermott &
Simoncelli 2011]

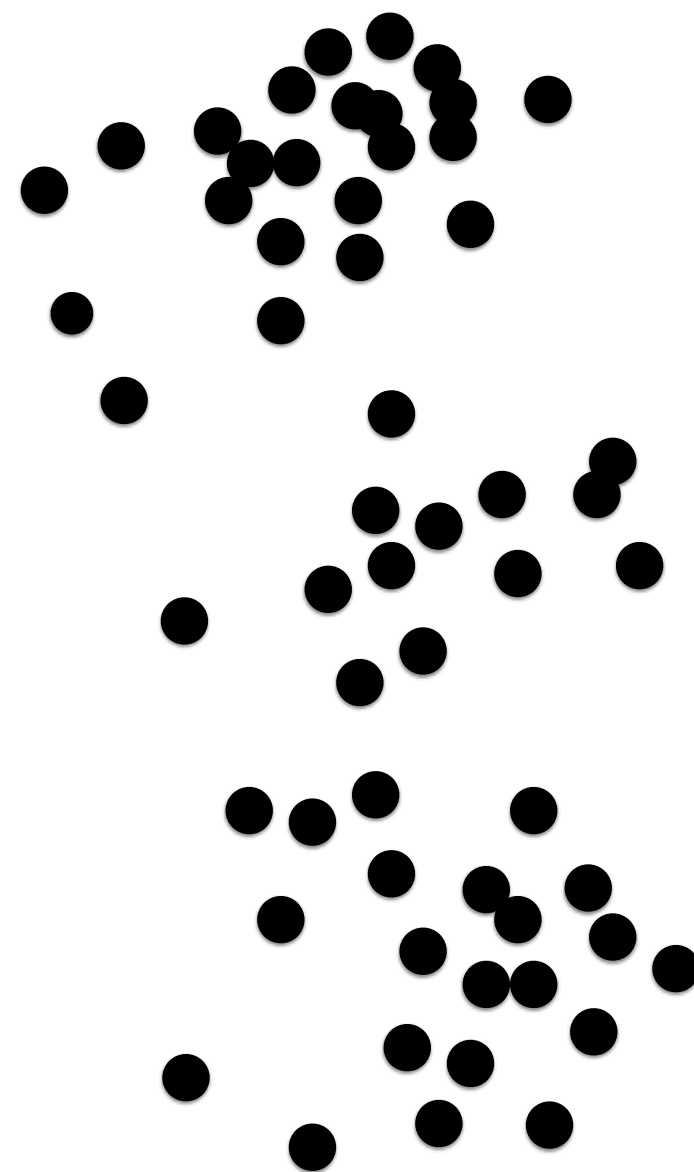
Predicting sound



Video frame



ConvNet

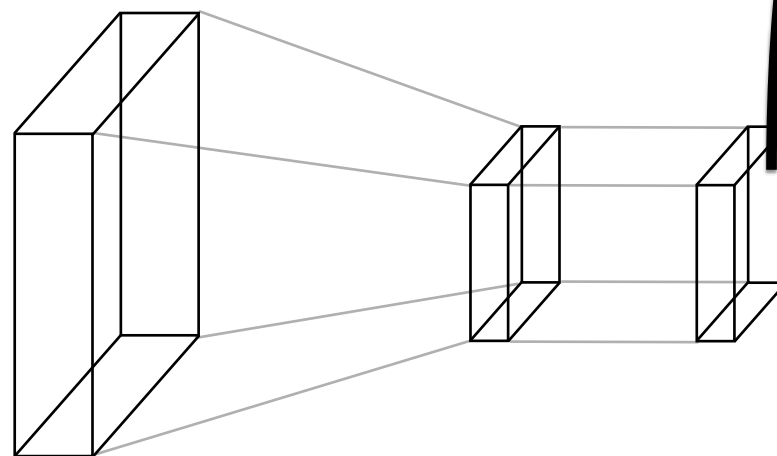


Sound feature

Predicting sound

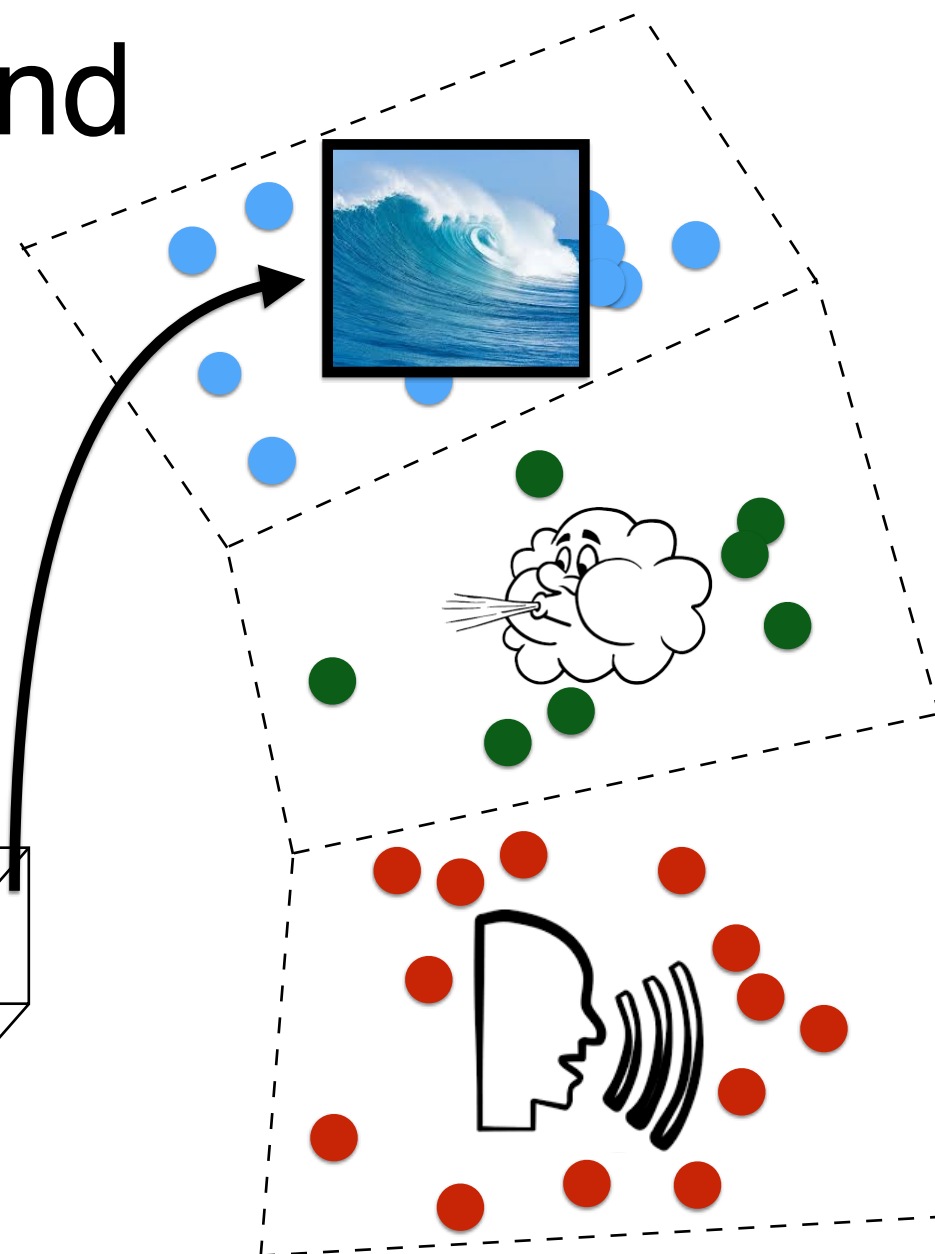


Video frame

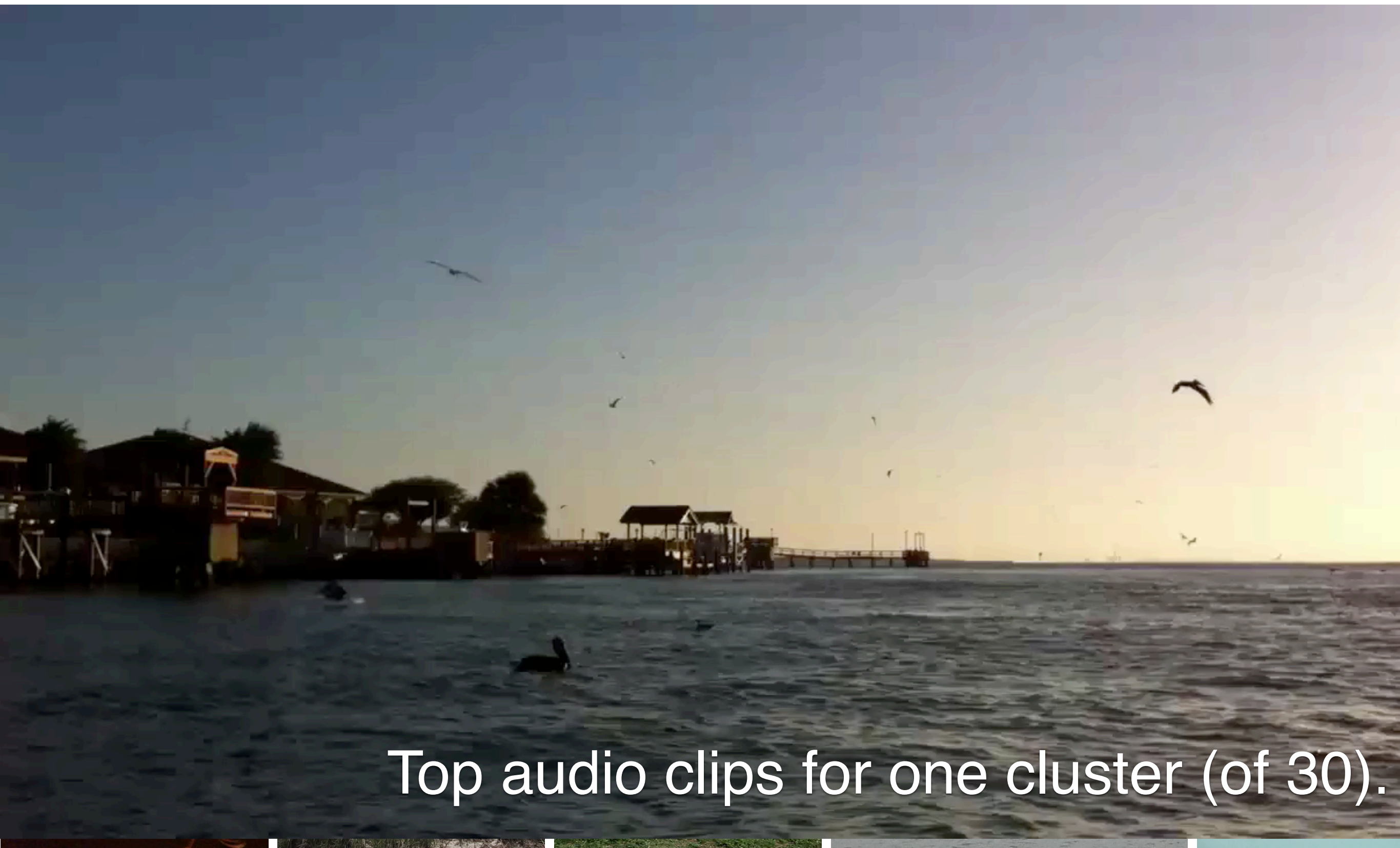


ConvNet

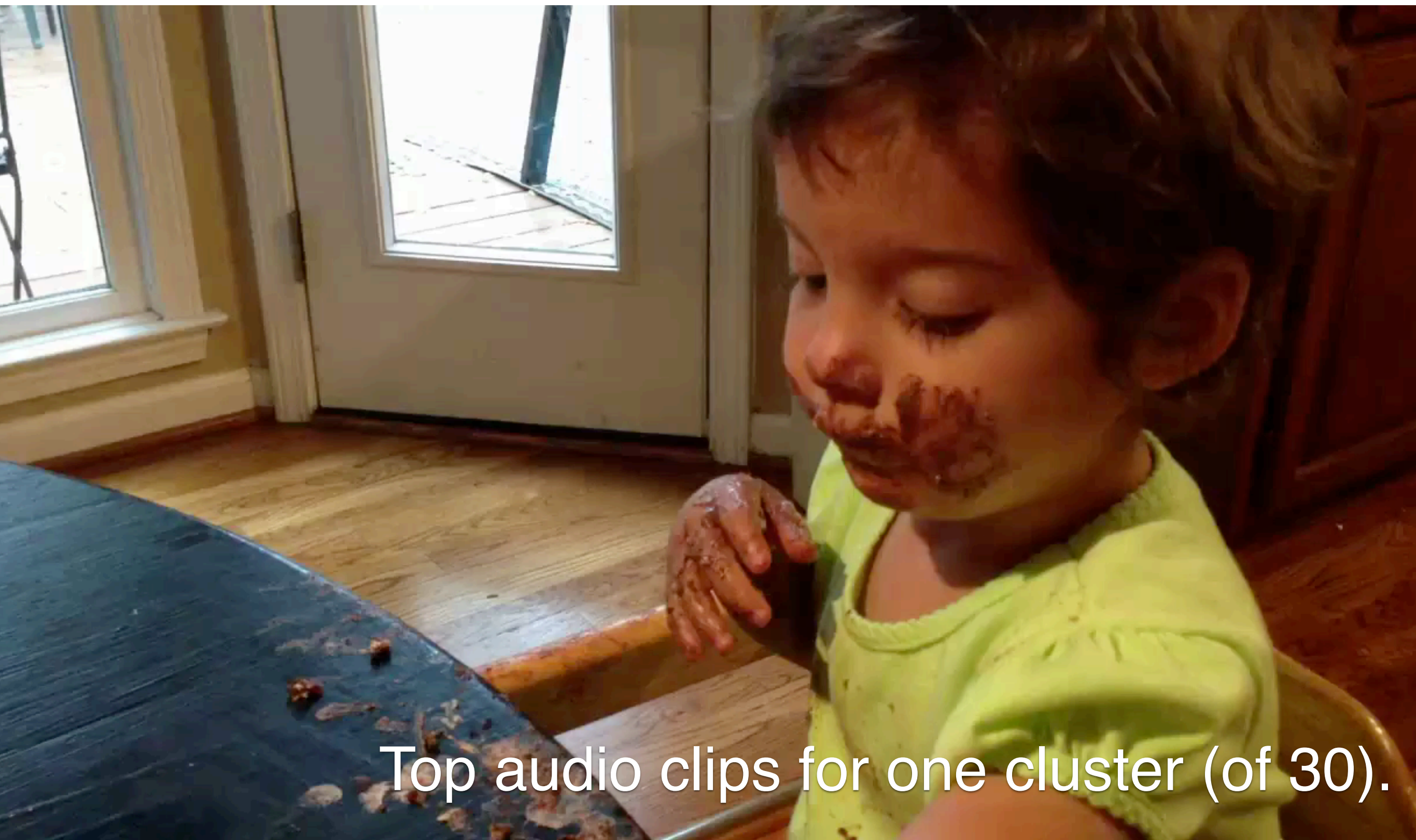
Audio



K-means or

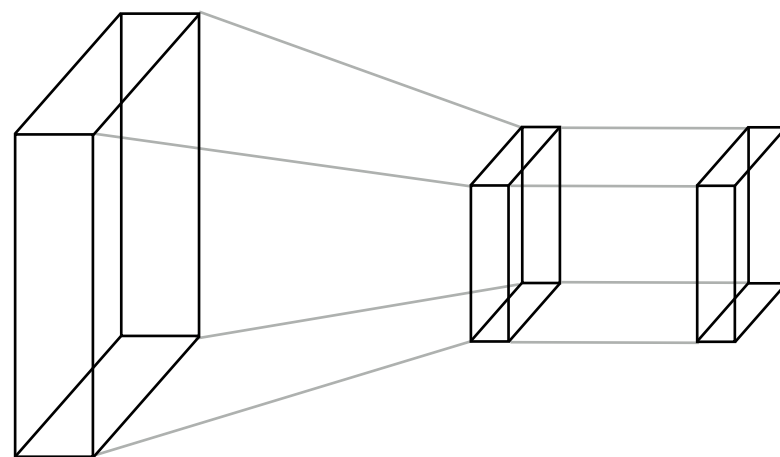


Top audio clips for one cluster (of 30).



Top audio clips for one cluster (of 30).

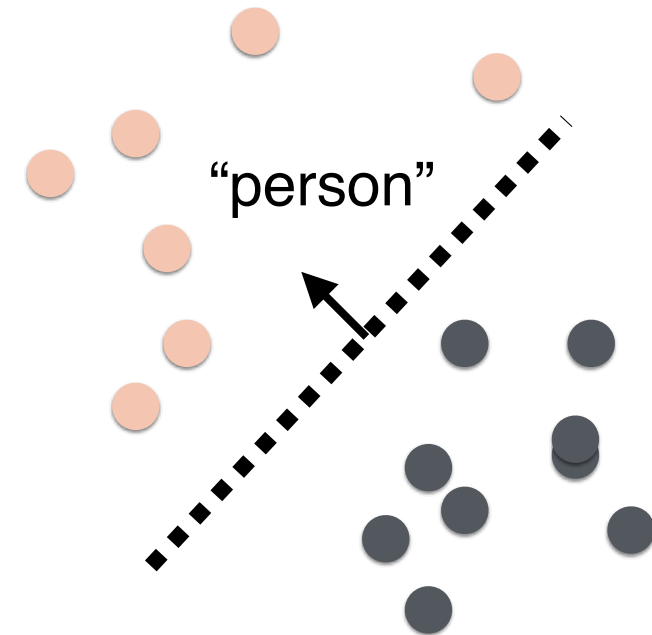
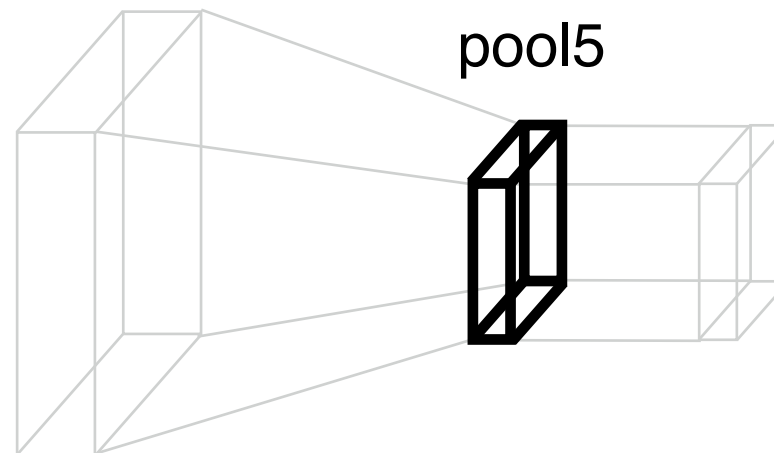
What did the network learn?



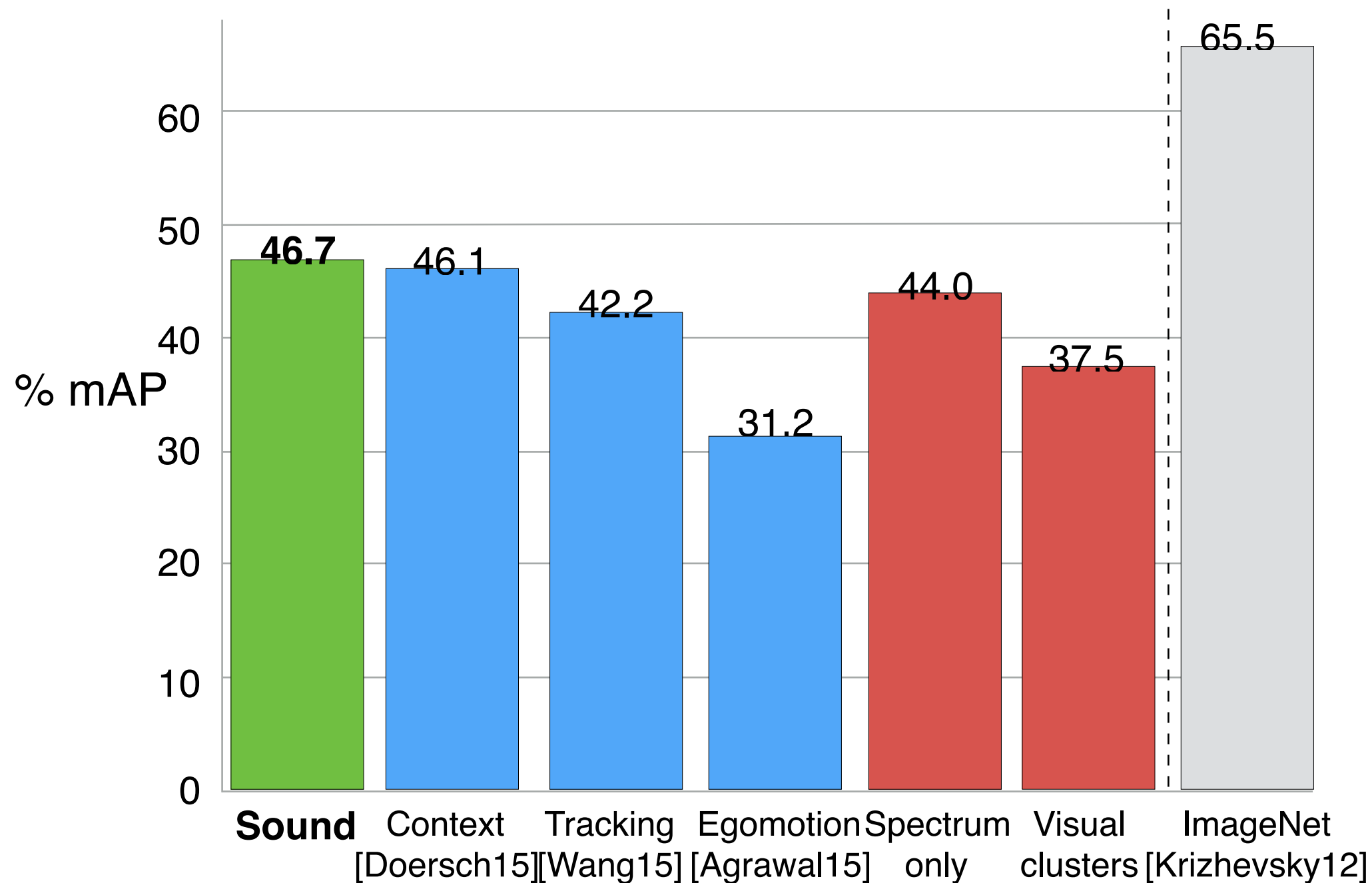
Audio label

What did the network learn?

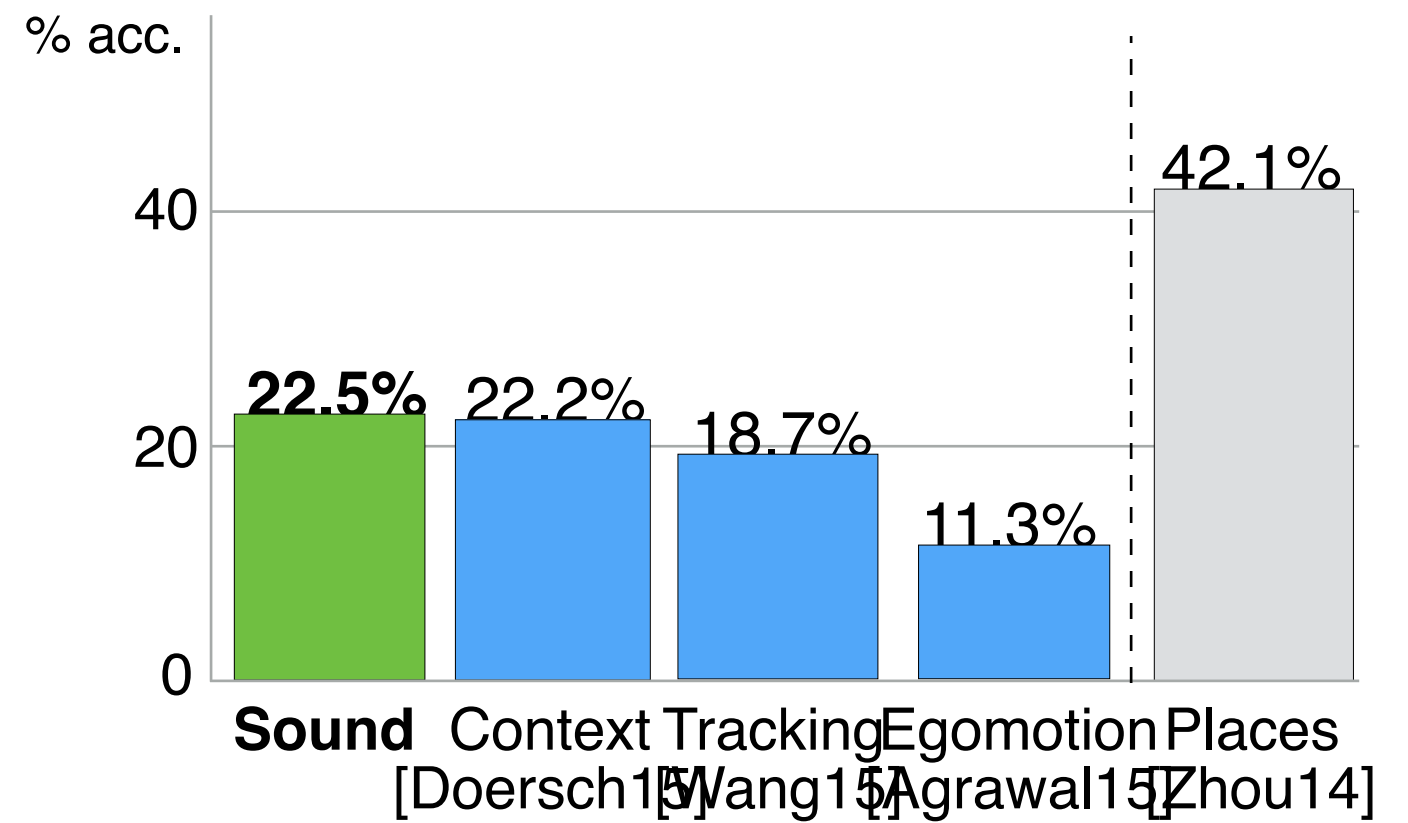
PASCAL VOC 2007



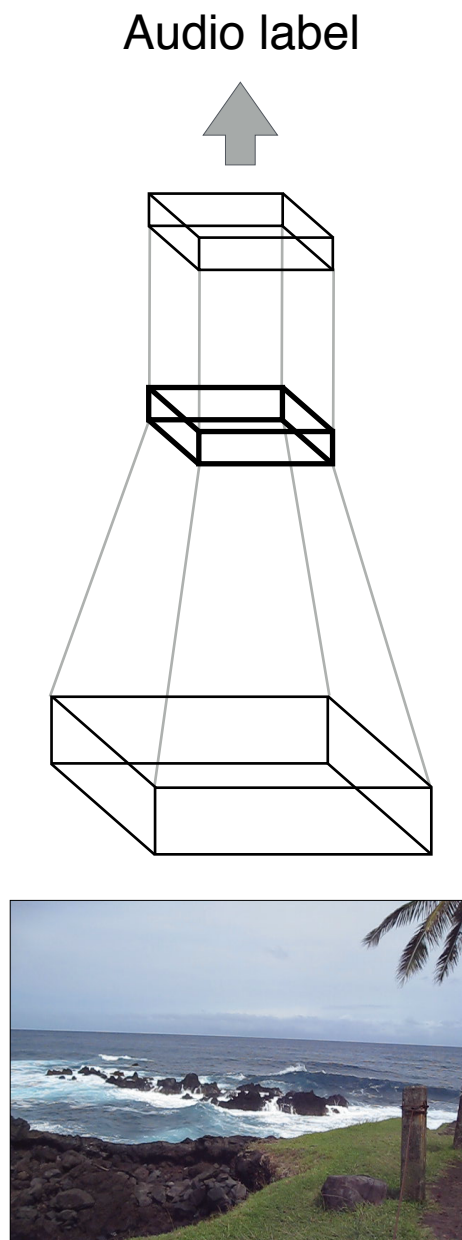
PASCAL VOC Classification



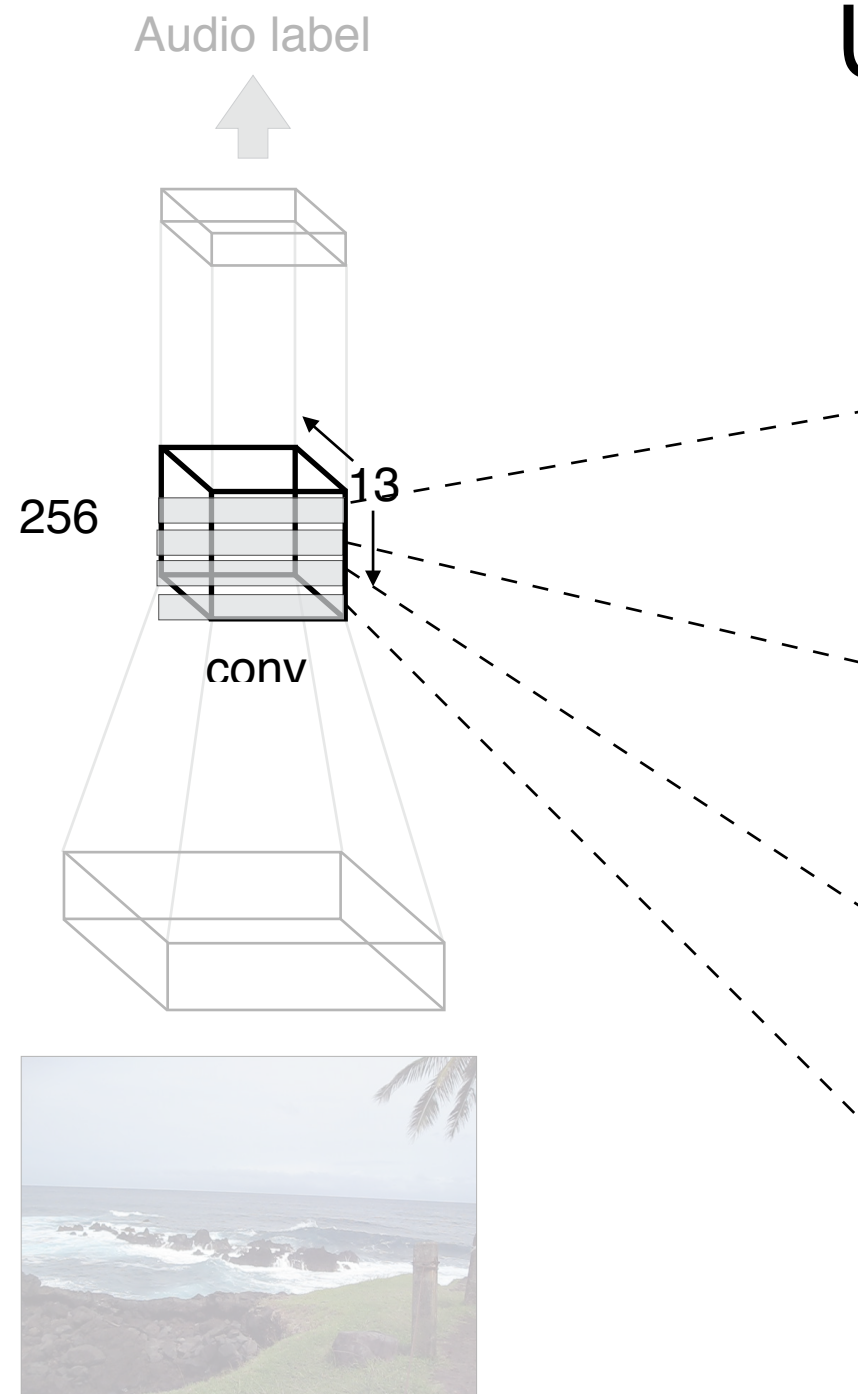
SUN397 Scene Recognition



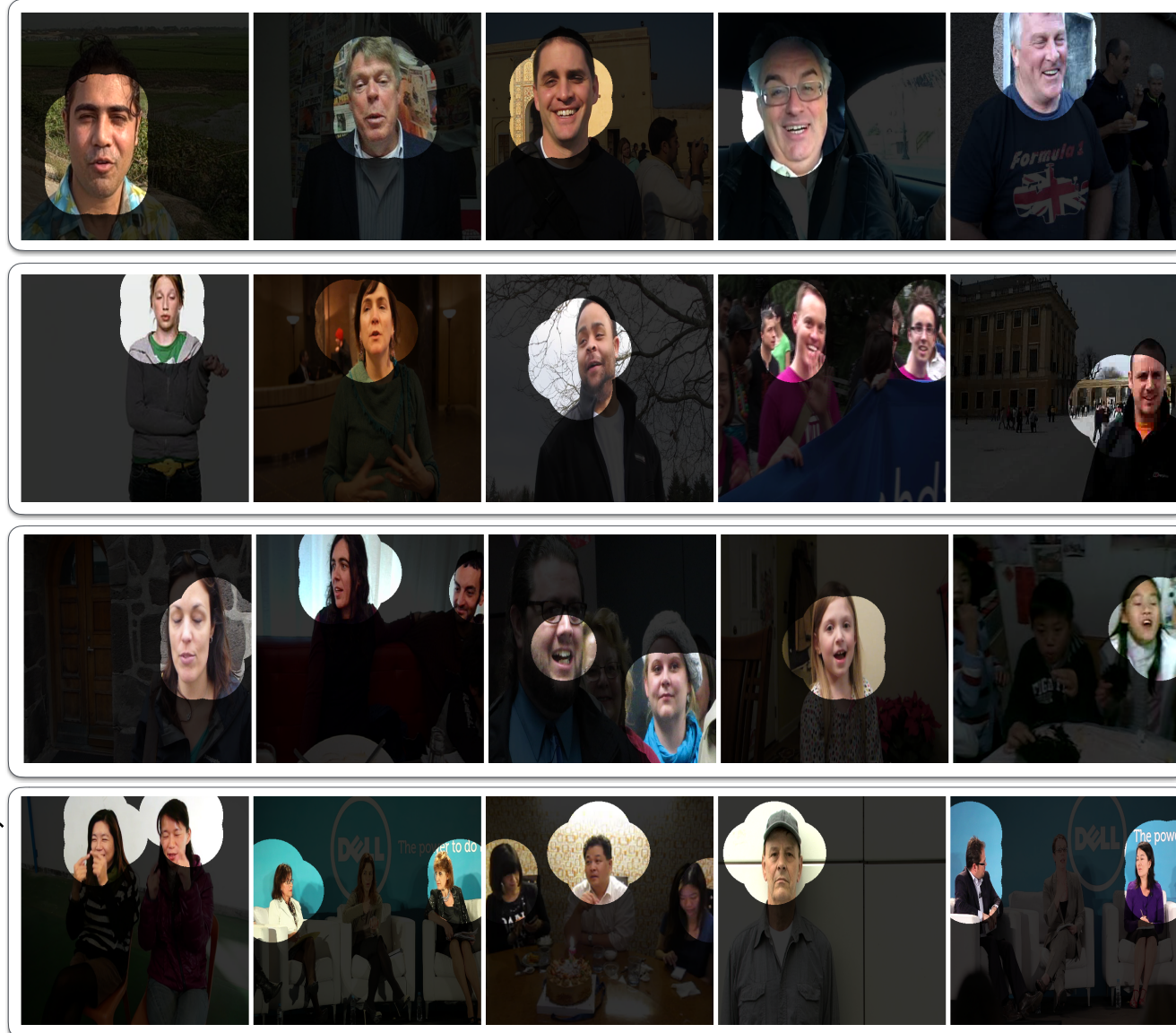
What did the network learn?



Unit visualizations



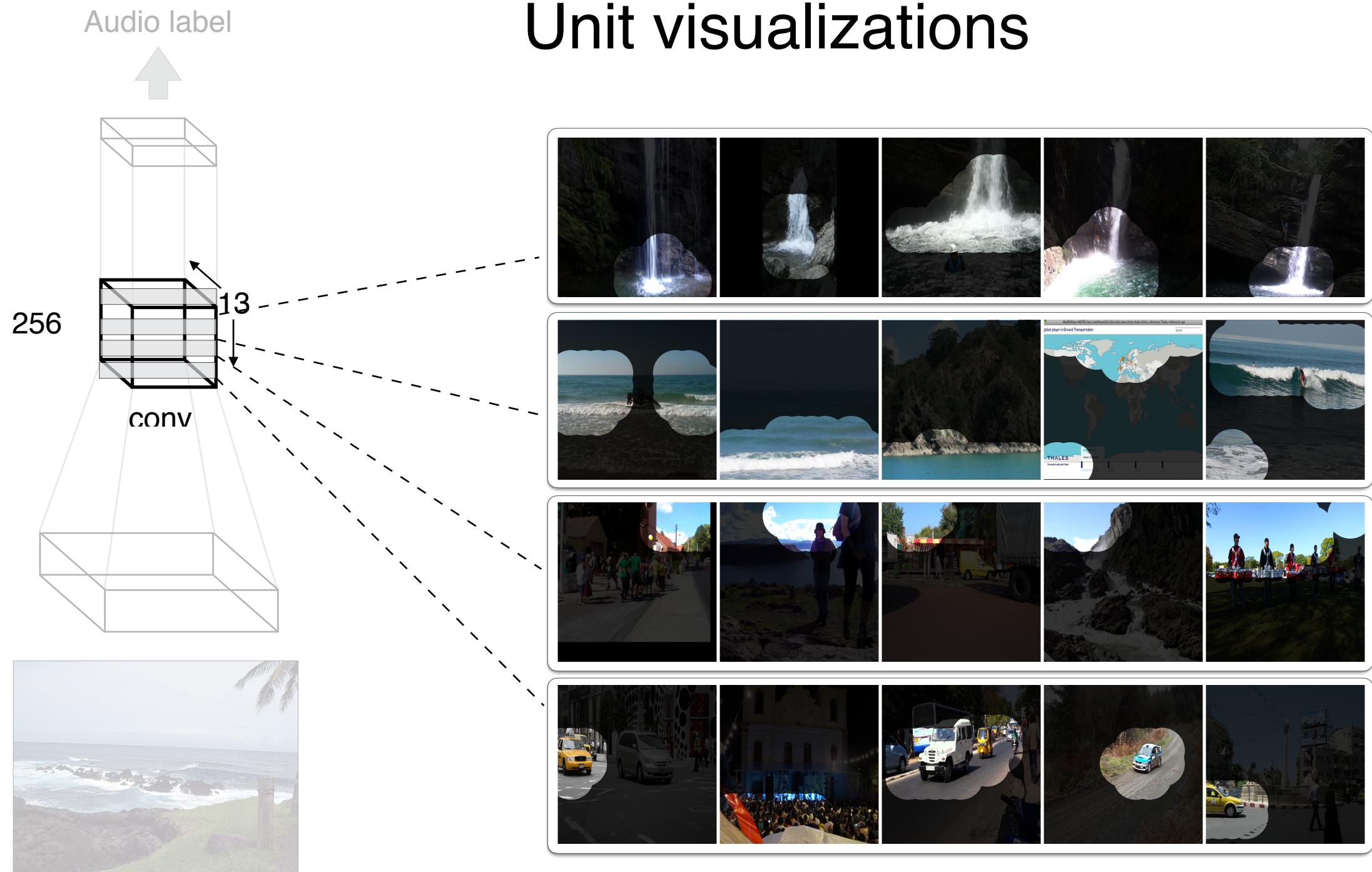
Top responses (unit #90)



Unit visualizations



Unit visualizations



Unsupervised visual representation learning by context prediction

[Carl Doersch, Abhinav Gupta, Alexei A. Efros, ICCV 2015]

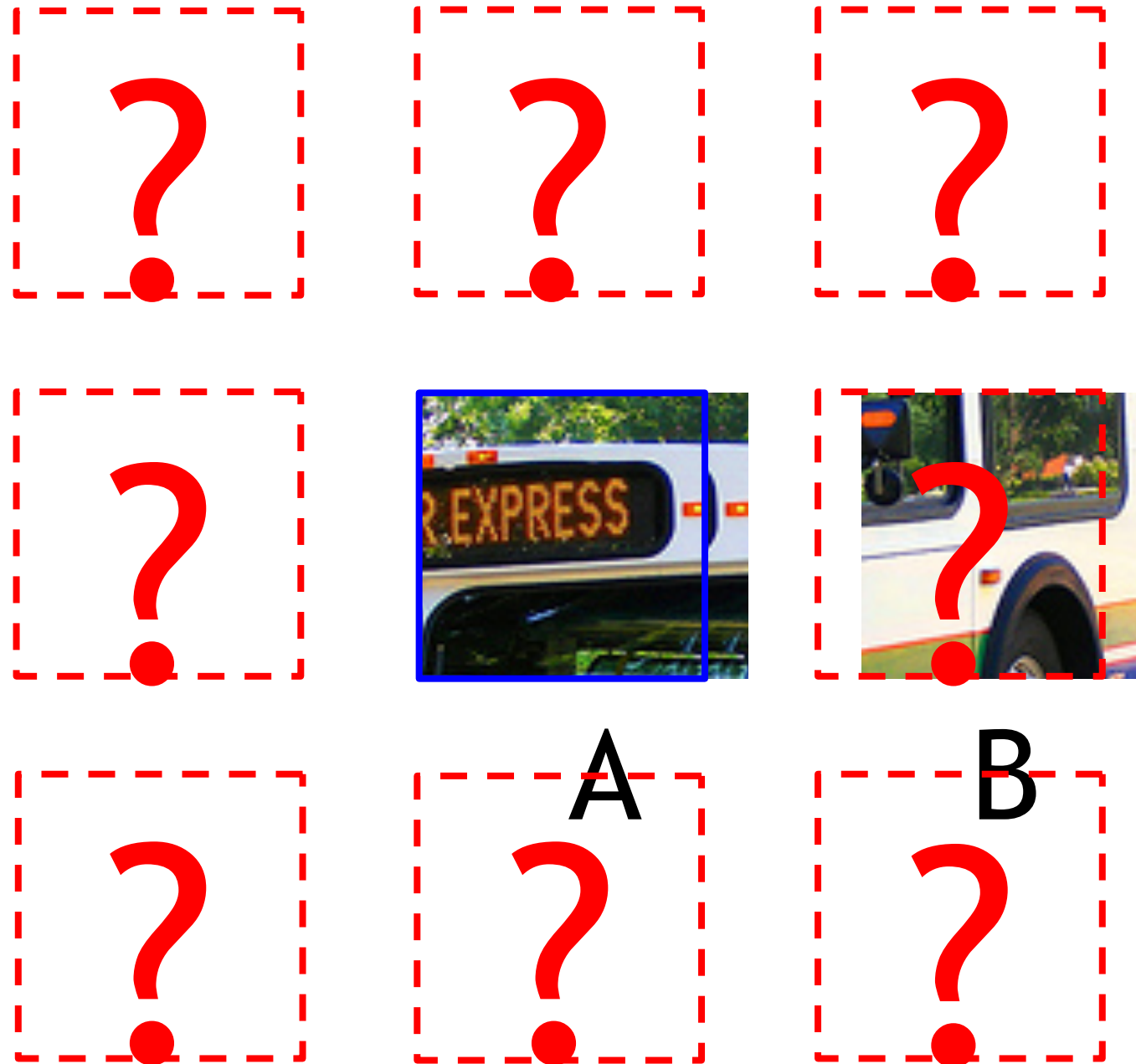
Context as Supervision

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

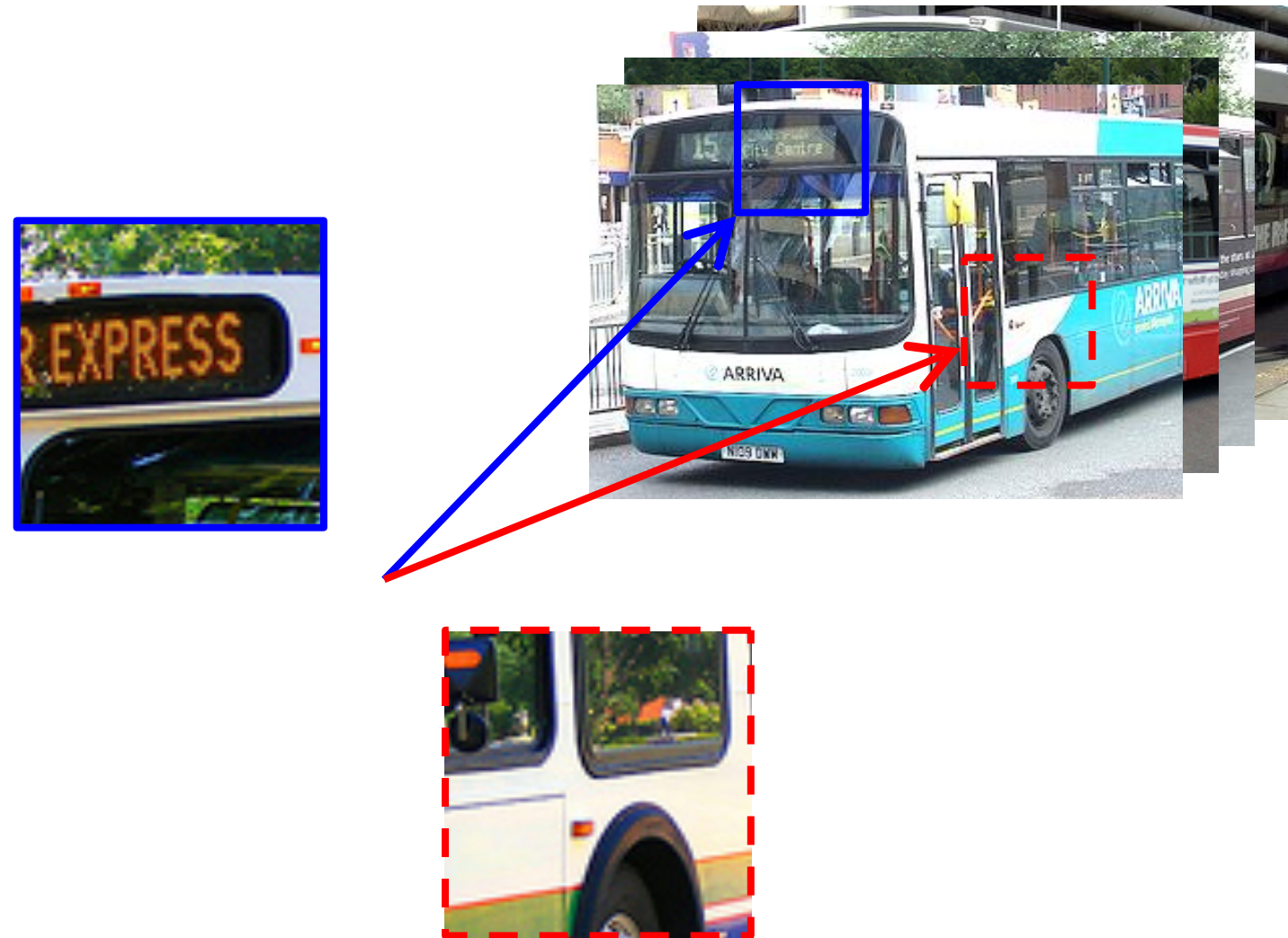
house, where the professor lived without his wife and child; or so he said jokingly sometimes: "Here's where I live. My house." His daughter often added, without resentment, for the visitor's information, "It started out to be for me, but it's really his." And she might reach in to bring forth an inch-high 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

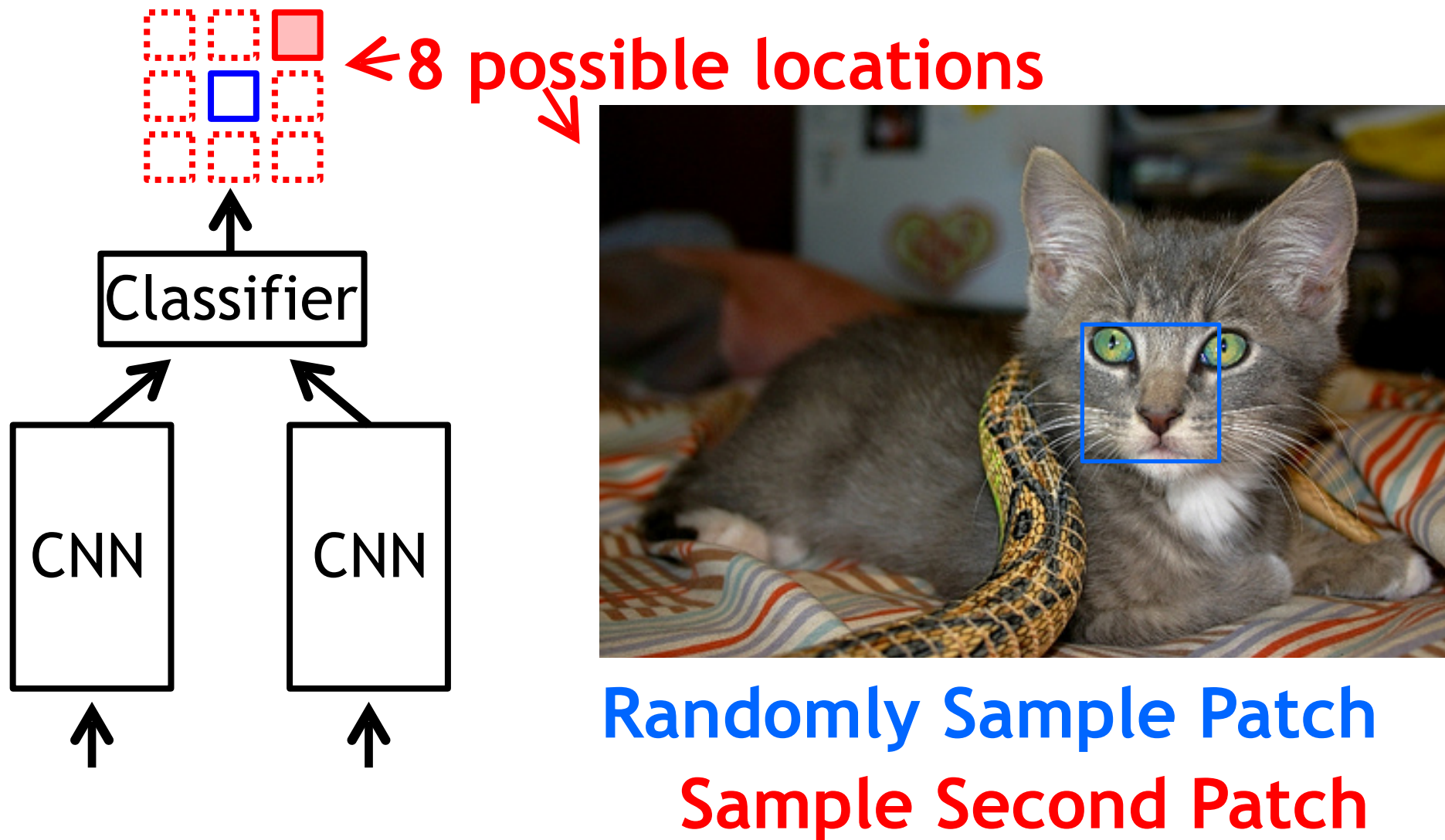
Context Prediction for Images

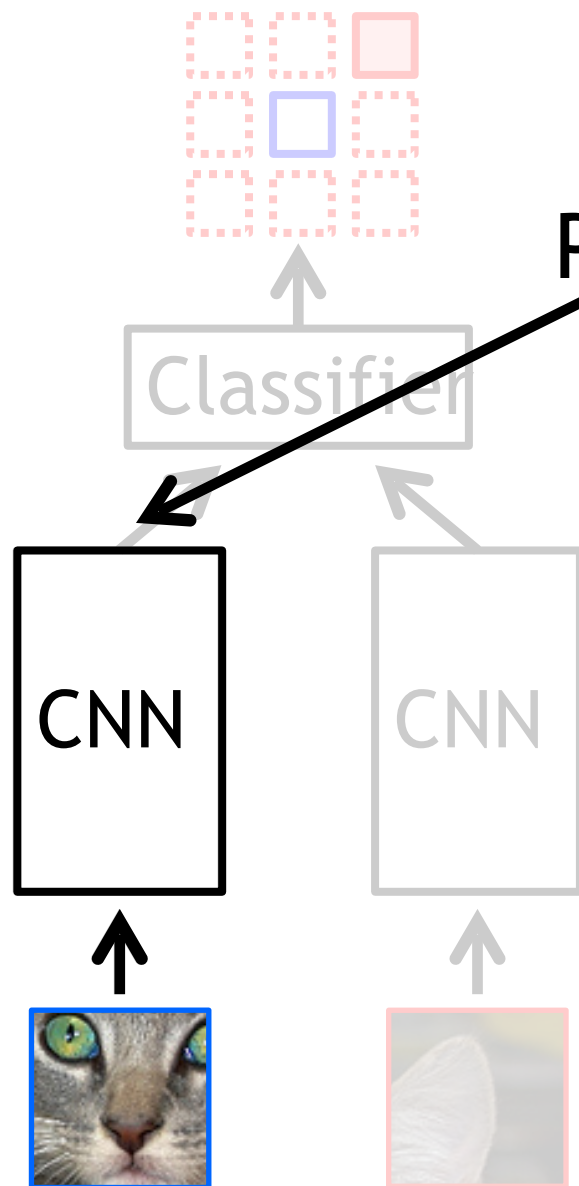


Semantics from a non-semantic task



Relative Position Task





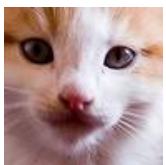
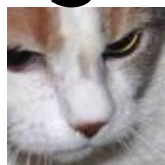
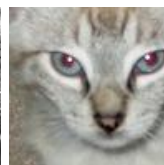
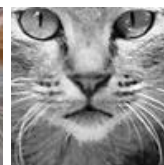
Patch Embedding (representation)

Input



!

Nearest Neighbors



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

Learning by Rotating



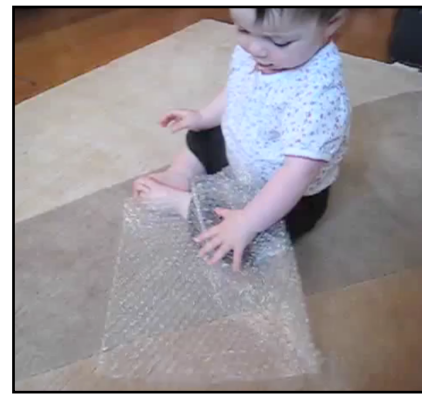
Unsupervised Representation Learning by Predicting Image Rotations
Spyros Gidaris, Praveer Singh, Nikos Komodakis

How are we doing?

	Classification	Detection	Segmentation
ImageNet	78.2%	56.8%	48.0%
Context	55.3%	46.6%	-
Jigsaw Puzzle	67.6%	53.2%	37.6%
Inpainting	56.5%	44.5%	30.0%
Colorization	61.5%	46.9%	35.6%
Tracking	58.7%	47.4%	-
Counting	67.7%	51.4%	36.6%
Rotation	72.9%	54.4%	39.1%

PASCAL VOC 2007

Prediction hypothesis



1. To survive, biological agents are constantly trying to anticipate, to predict sensations
2. This trains up representations useful for prediction — surfaces, objects, events!



Henri Cartier-Bresson



Yann LeCun's cake:

1. Cake is unsupervised representation learning
2. Frosting is supervised transfer learning
3. Cherry on top is reinforcement learning (model-based RL)