

# How much do deep neural networks understand about the images they recognize?

Jeff Clune

Assistant Professor, Computer Science, University of Wyoming  
Director, Evolving AI Lab



EVOLVING  
ARTIFICIAL  
INTELLIGENCE  
LABORATORY



UNIVERSITY  
OF WYOMING





EVOLVING  
ARTIFICIAL  
INTELLIGENCE  
LABORATORY

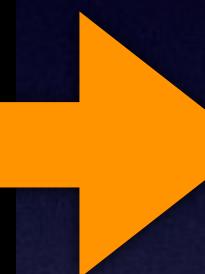
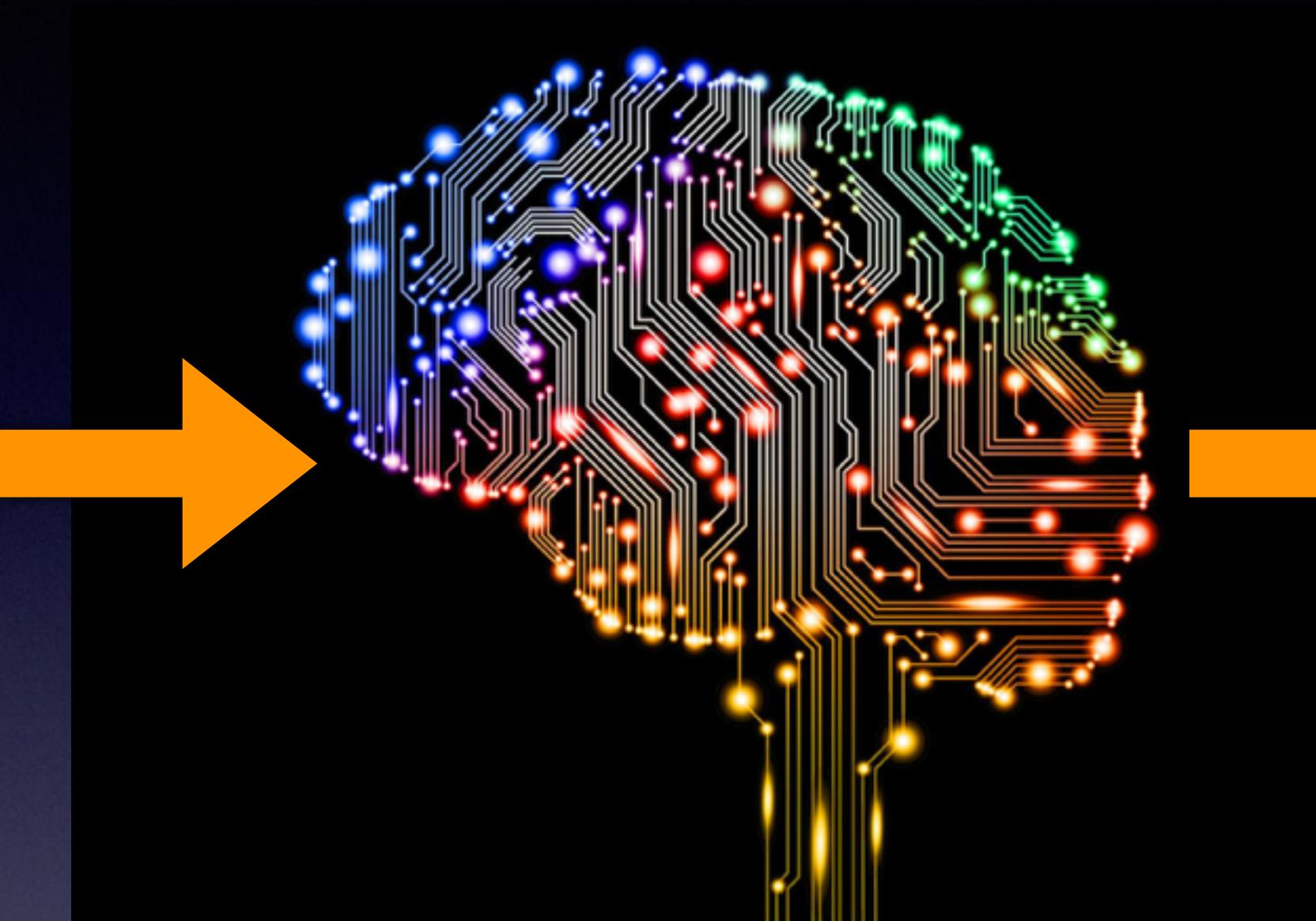
UNIVERSITY  
OF WYOMING

# Robotics & AI

- Robots that recover from damage
- Deep Learning
- Evolving neural network robot controllers
- Creative robots
- Computational biology
- etc.

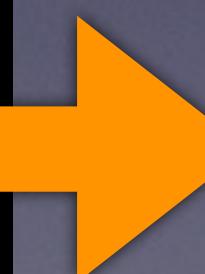
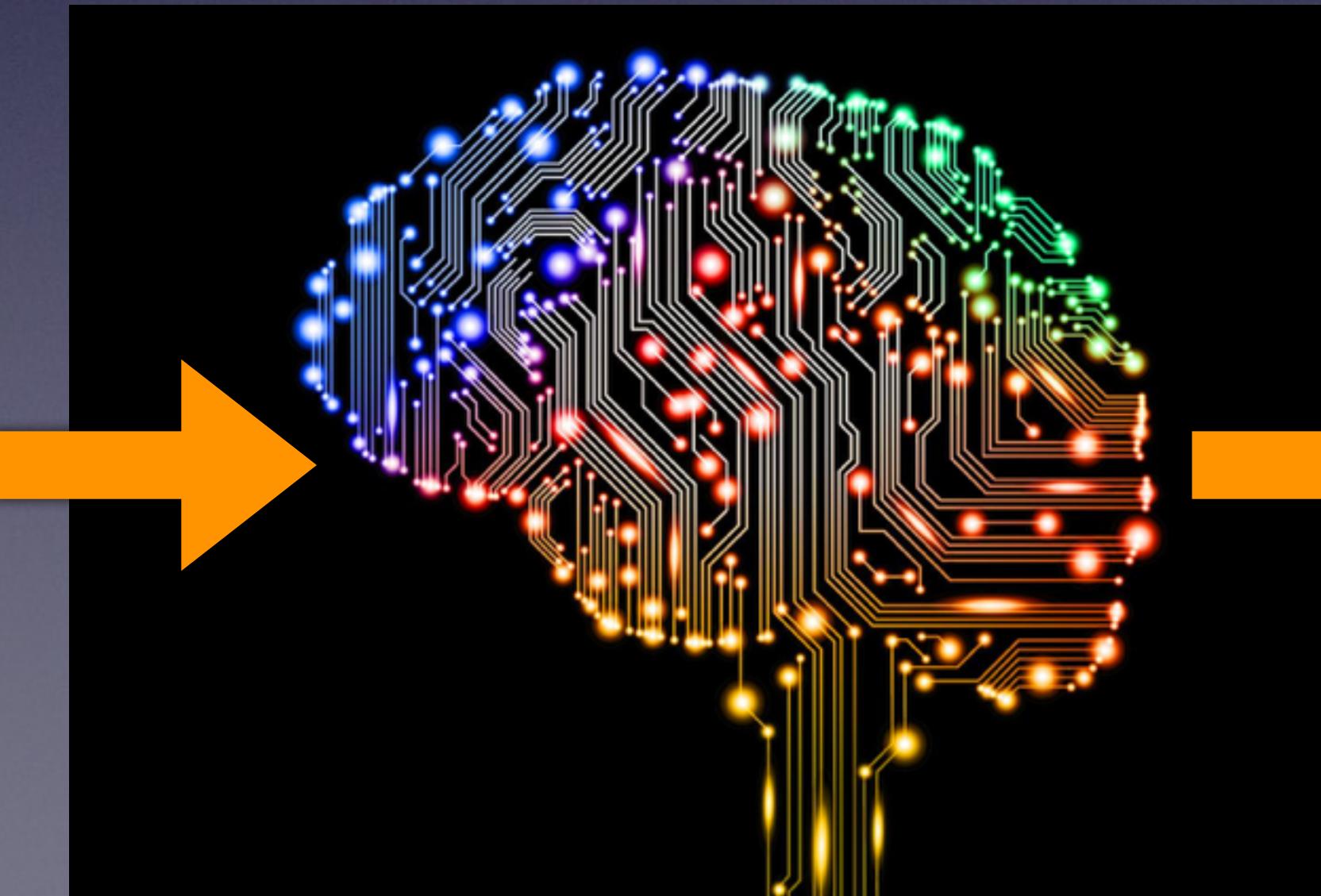
The cover of the journal **nature** (THE INTERNATIONAL WEEKLY JOURNAL OF SCIENCE) features a large image of a black and silver robotic arm with multiple cameras and sensors, positioned on a grid-like surface that appears to be a heatmap or a sensor array. The title 'nature' is written in a large, lowercase, sans-serif font across the top. A small circular logo in the top right corner reads 'INSIGHT Machine intelligence'. Below the main image, the headline 'Back on its feet' is displayed, followed by the subtext 'Using an intelligent trial-and-error learning algorithm this robot adapts to injury in minutes'. A small caption 'PAGES 426 & 503' is also present. At the bottom of the cover, there are three columns of text: 'COGNITION WHY FISH NEED TO BE CLEVER Social behaviours need plenty of brainpower PAGE 412', 'ARTIFICIAL INTELLIGENCE LIVING WITH ROBOTS AI researchers' ethics prescriptions PAGE 415', and 'HUMAN EVOLUTION ANOTHER FACE IN THE CROWD A new hominin from Ethiopia's middle Pliocene PAGES 432 & 483'. The journal's logo 'NATURE.COM/NATURE' and the issue details '28 May 2015 £10 Vol. 521, No. 7555' are at the very bottom.

# Deep Neural Networks/Deep Learning



**“Lion”**

**“Cogito,  
ergo sum”**



**“I think,  
therefore I am”**

# DNNs as good as humans at image recognition

**predictions for natural images**



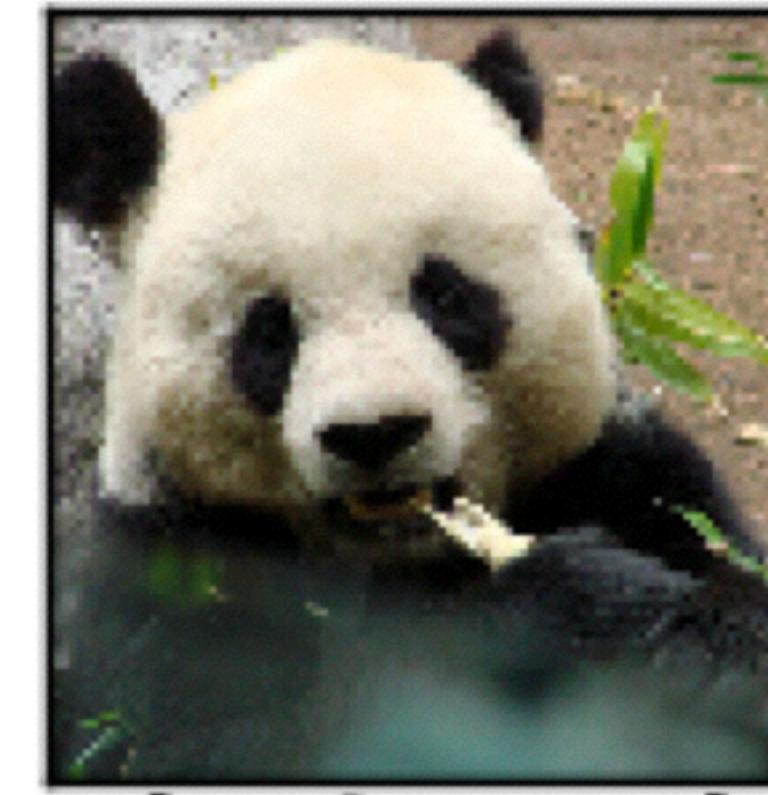
**ladybug**



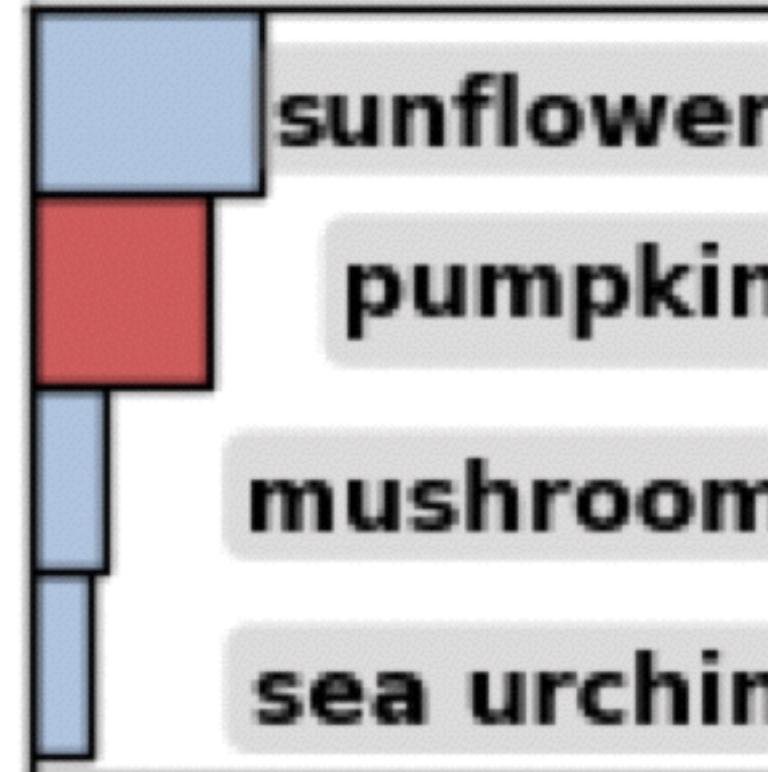
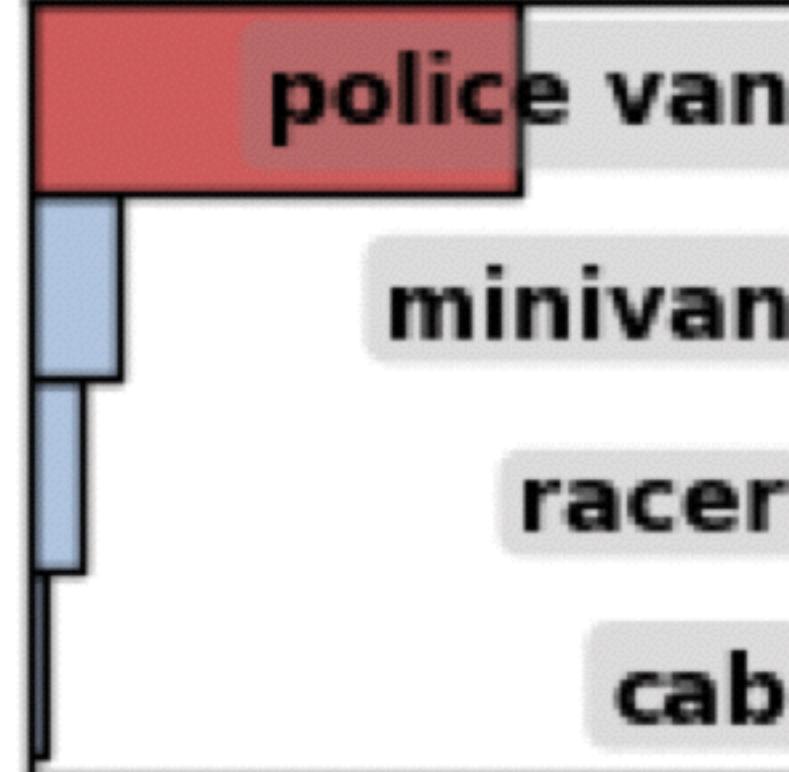
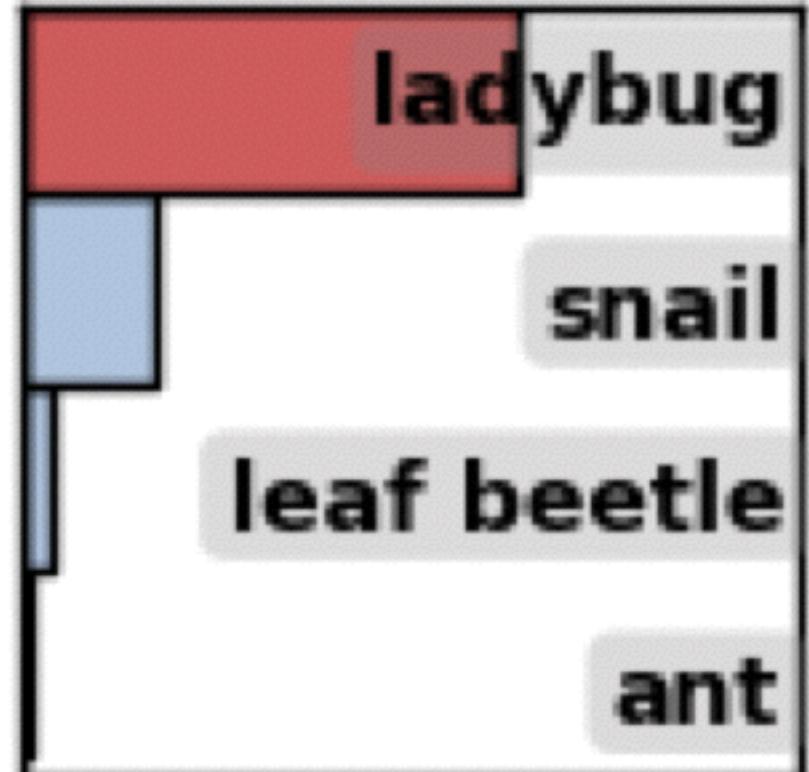
**police van**



**pumpkin**

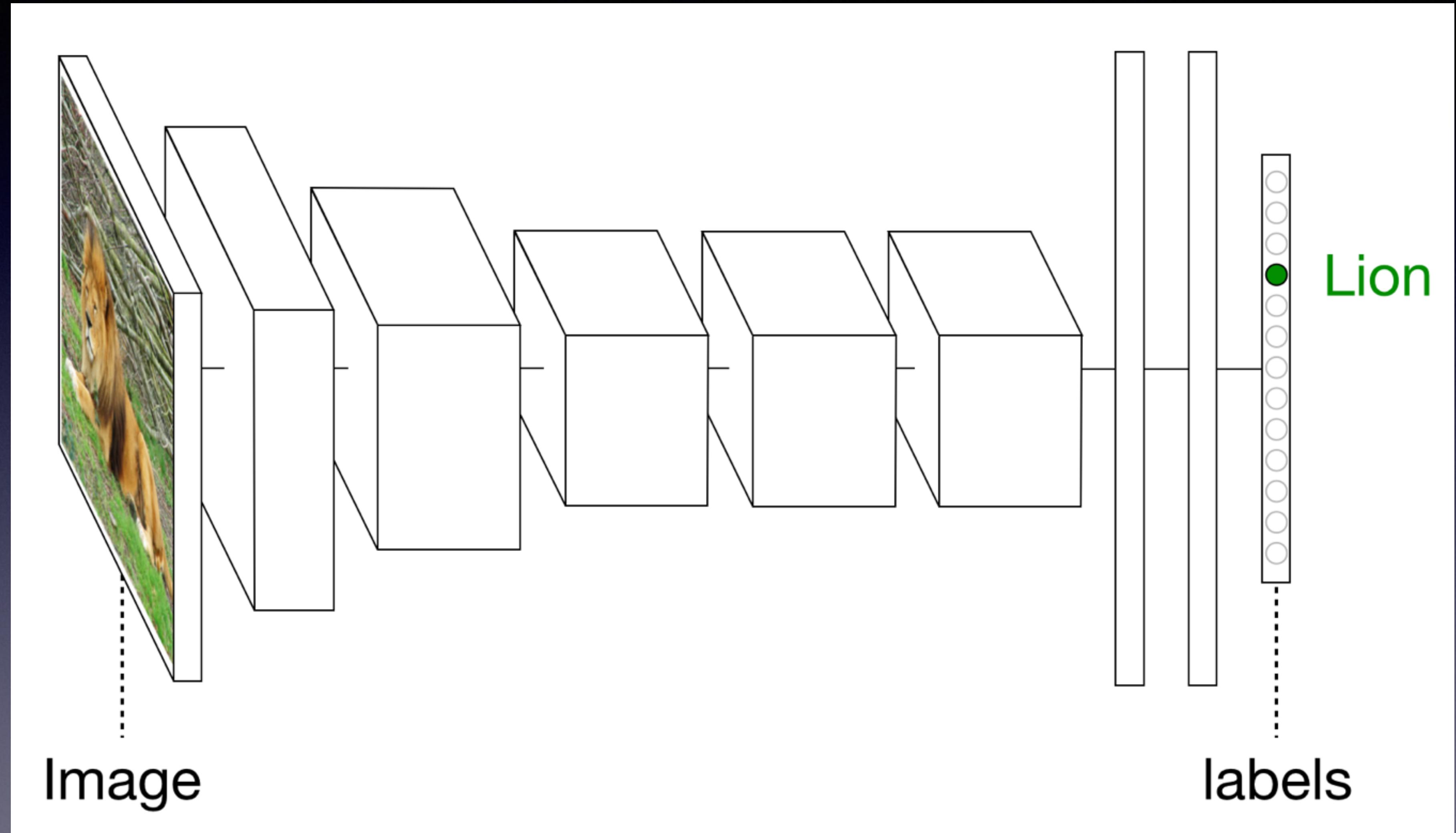


**giant panda**



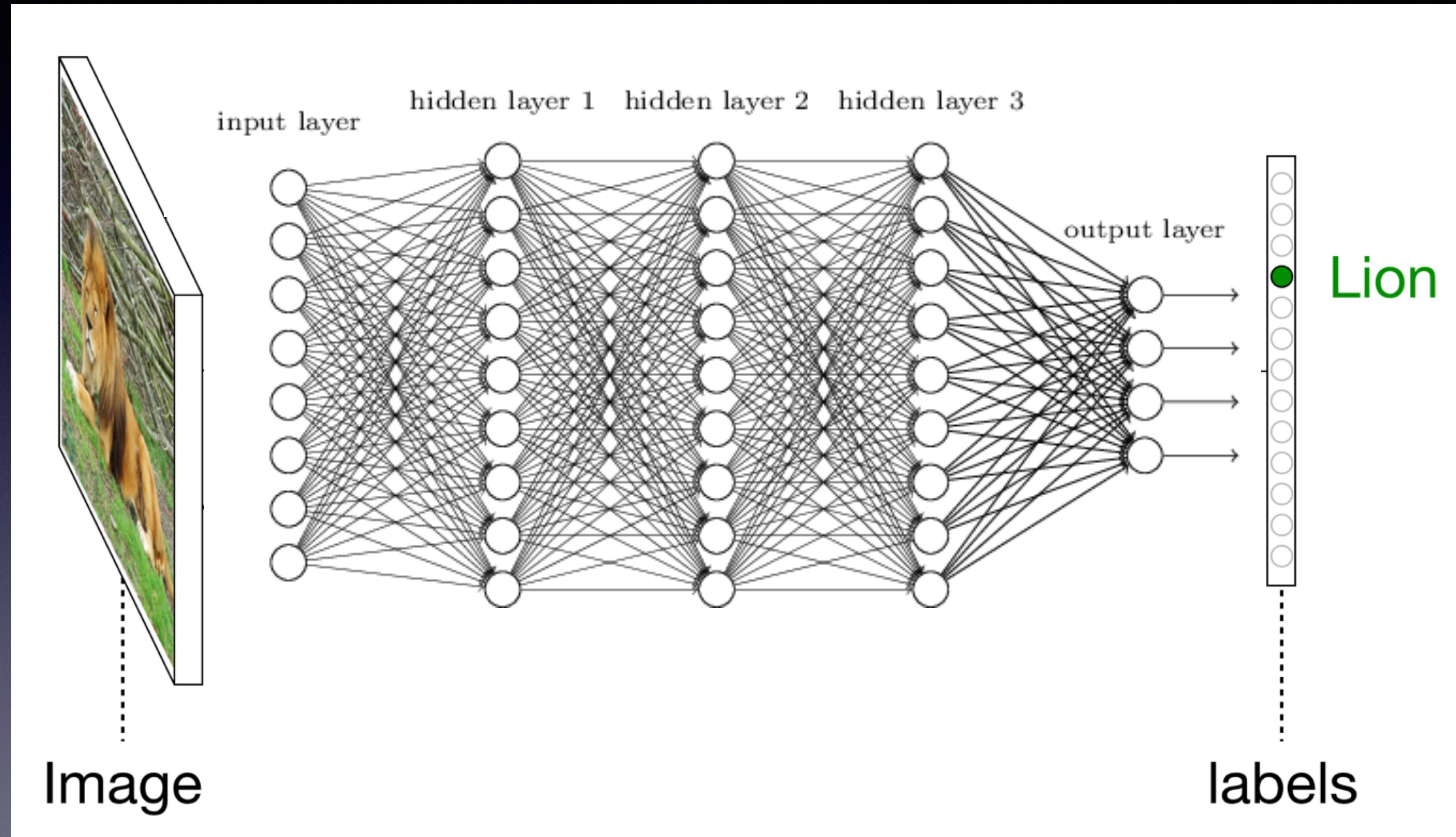
ImageNet  
1,000 Categories  
1.3 M Images  
Human error: 5%  
DNN: 3%

# Deep Neural Networks/Deep Learning



~1M neurons  
~100M weights

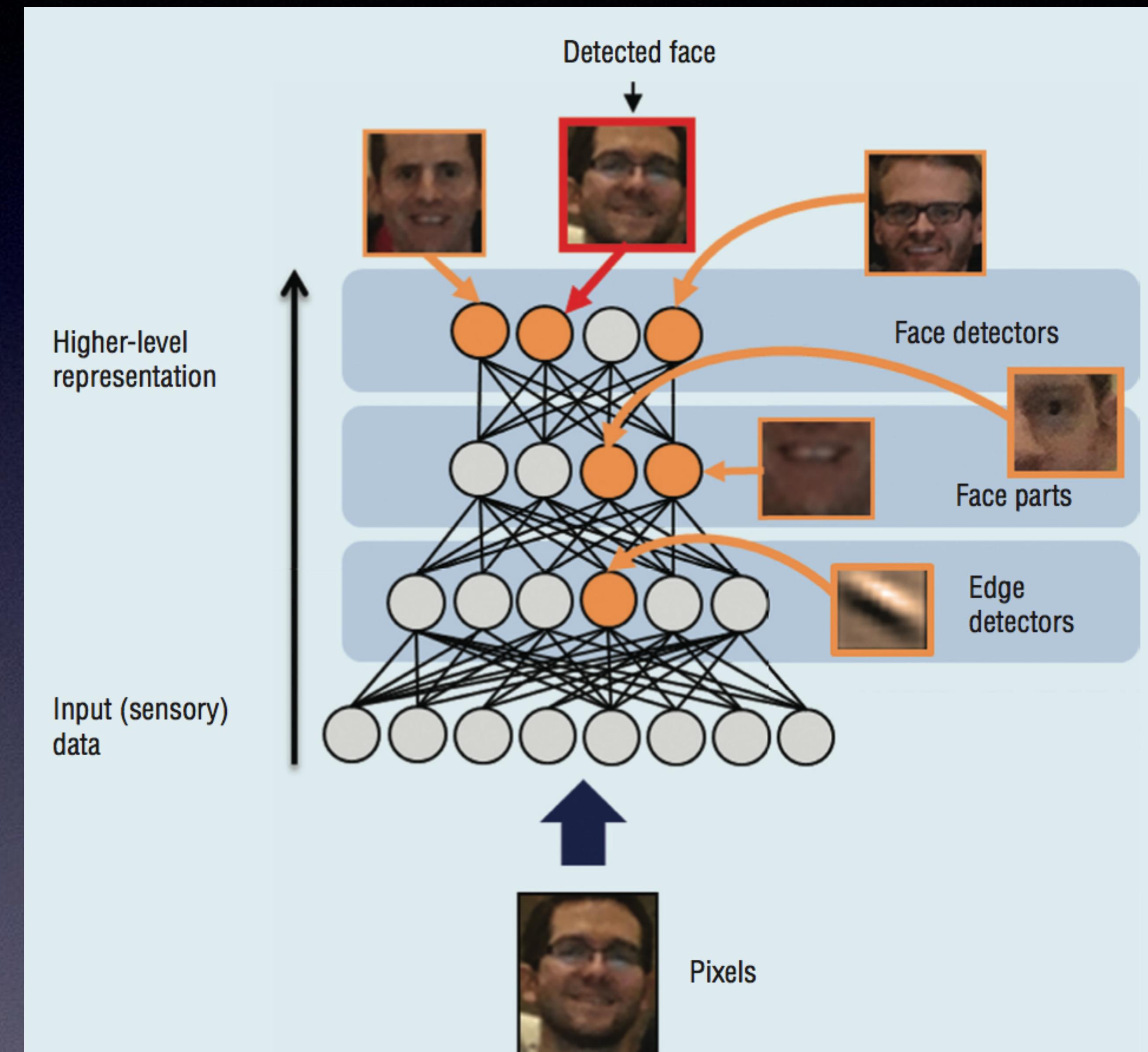
# Deep Neural Networks/Deep Learning



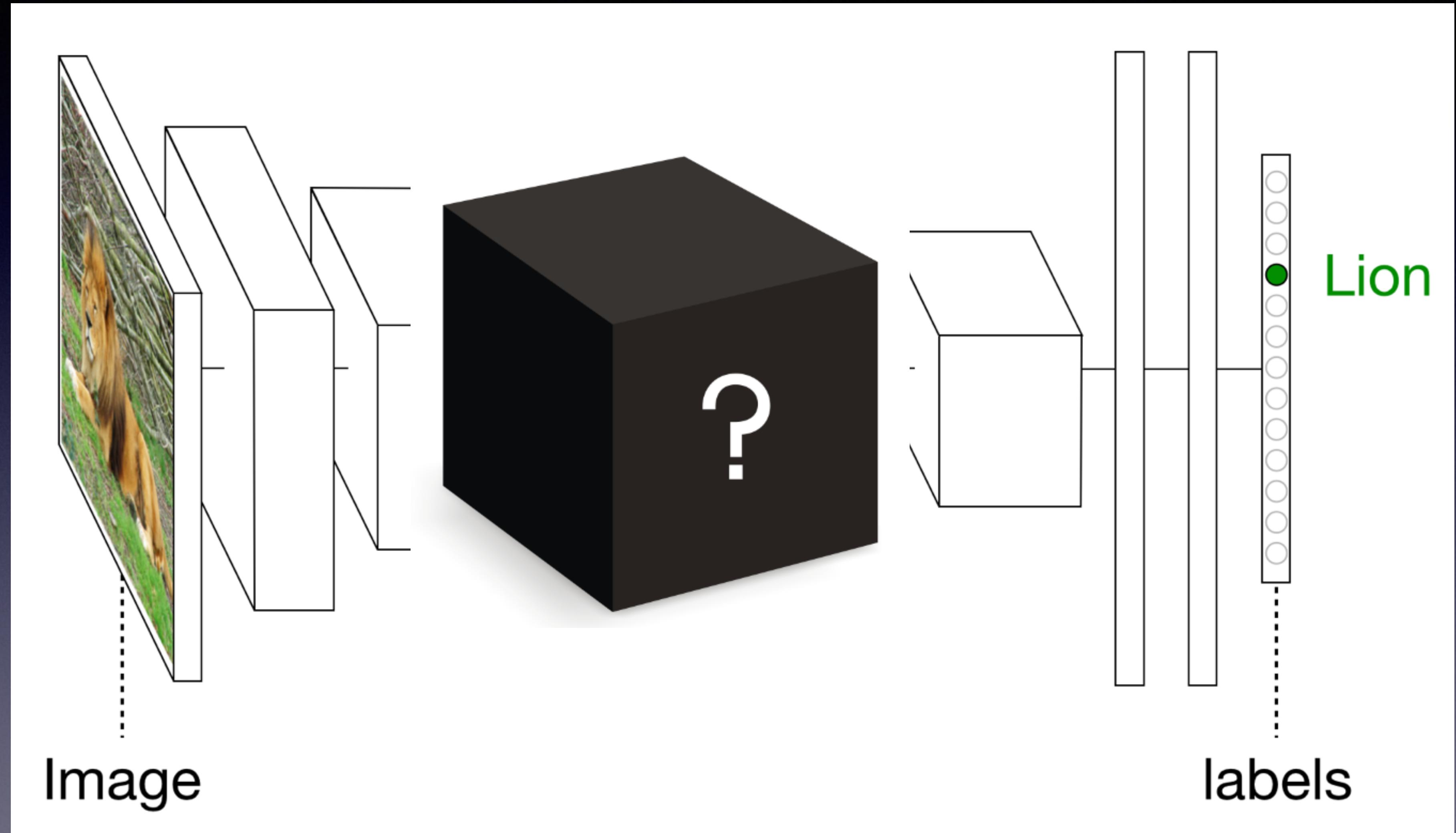
# DNNs: The Hope

- Each layer transforms data
- Recognizes increasingly abstract “features”

Hierarchically composed feature representations

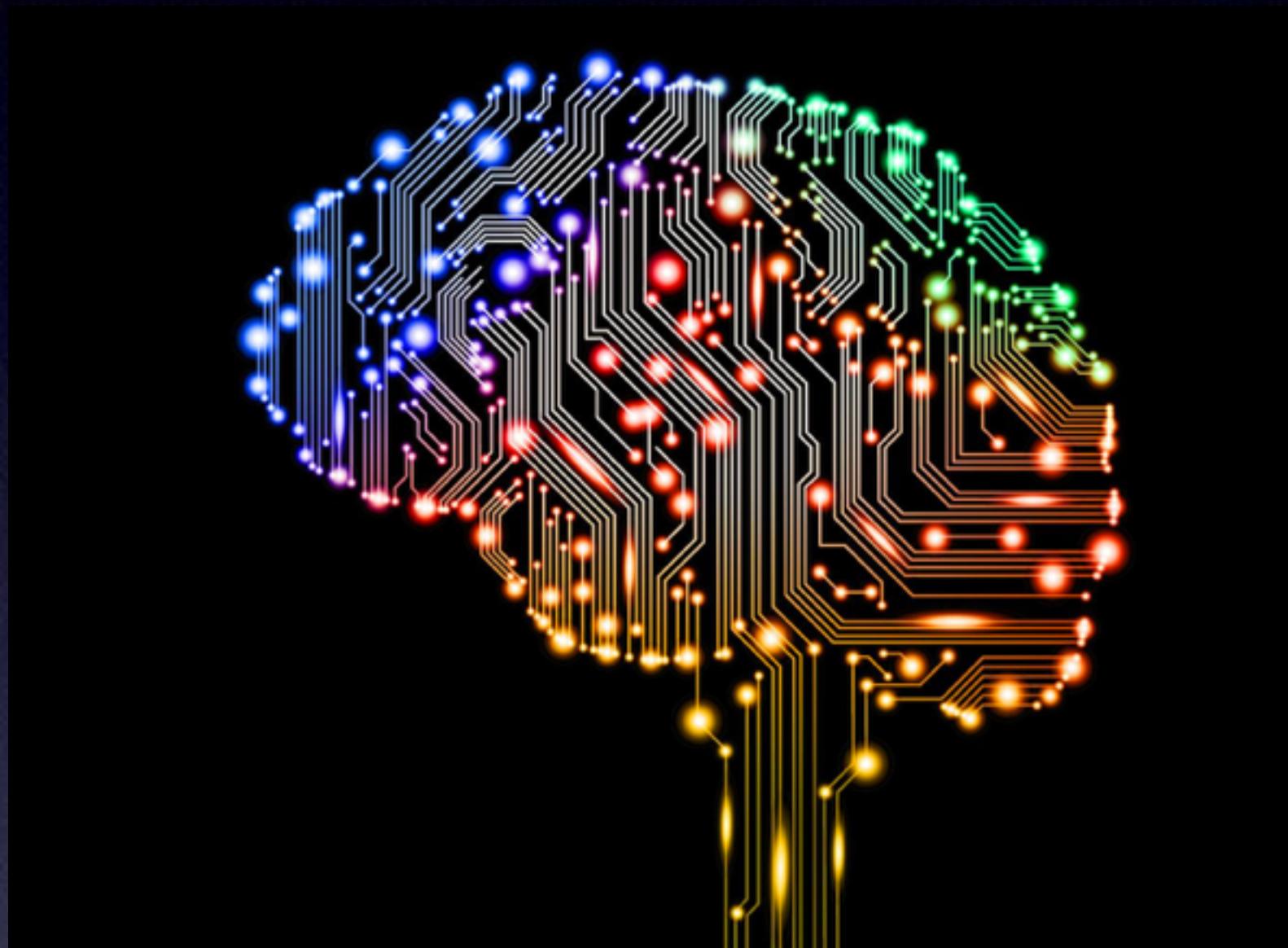


# Deep Neural Networks/Deep Learning



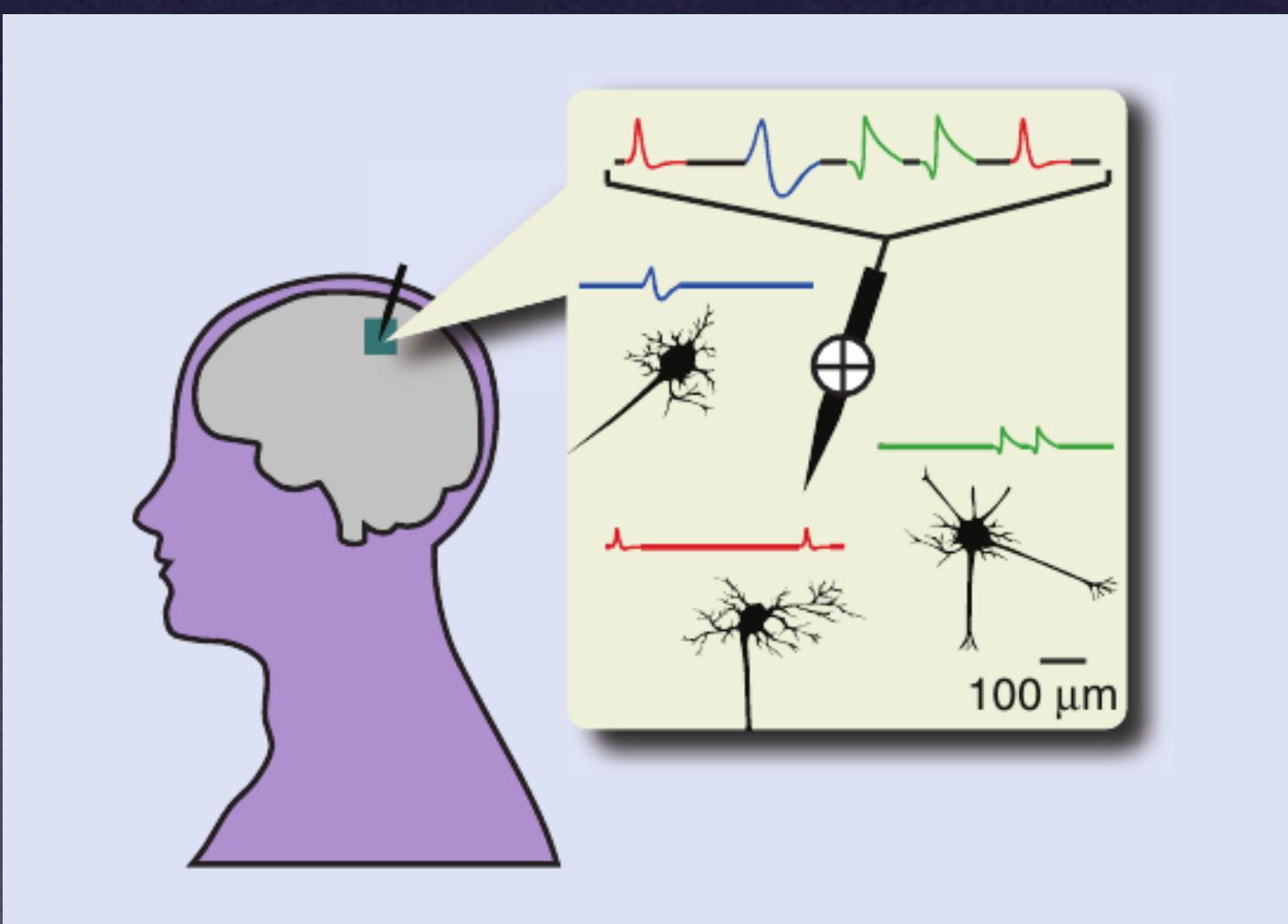
~1M neurons  
~100M weights

# Thus begins the field of “AI Neuroscience”



# One way do neuroscientists do it (of many)

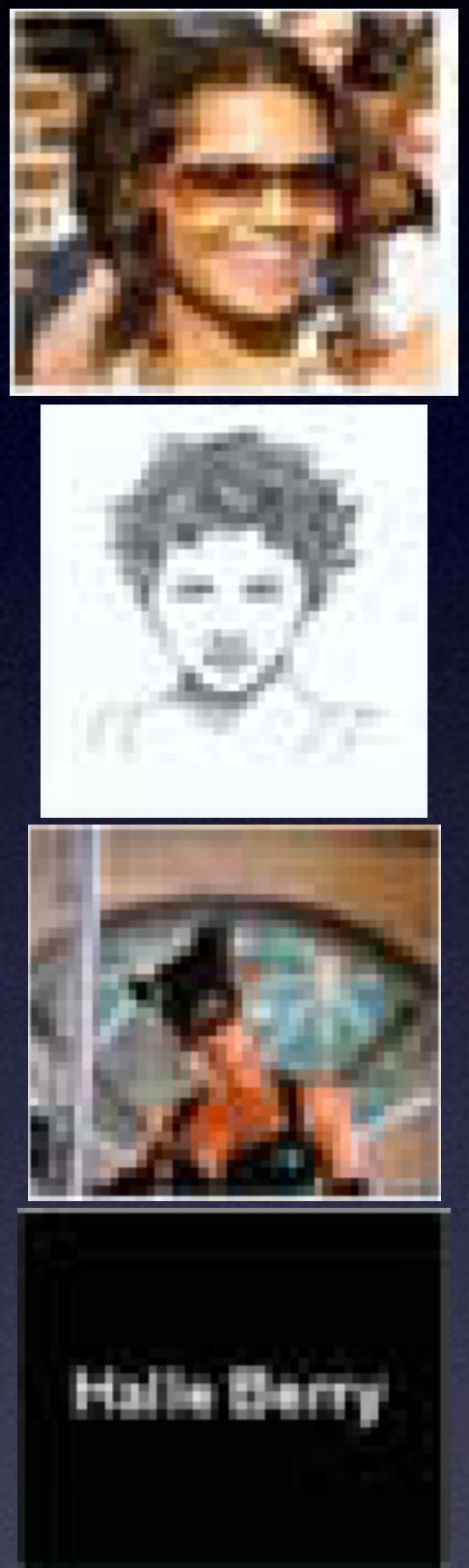
- Record a single neuron
- Show it pictures
- See what it responds to



etc....

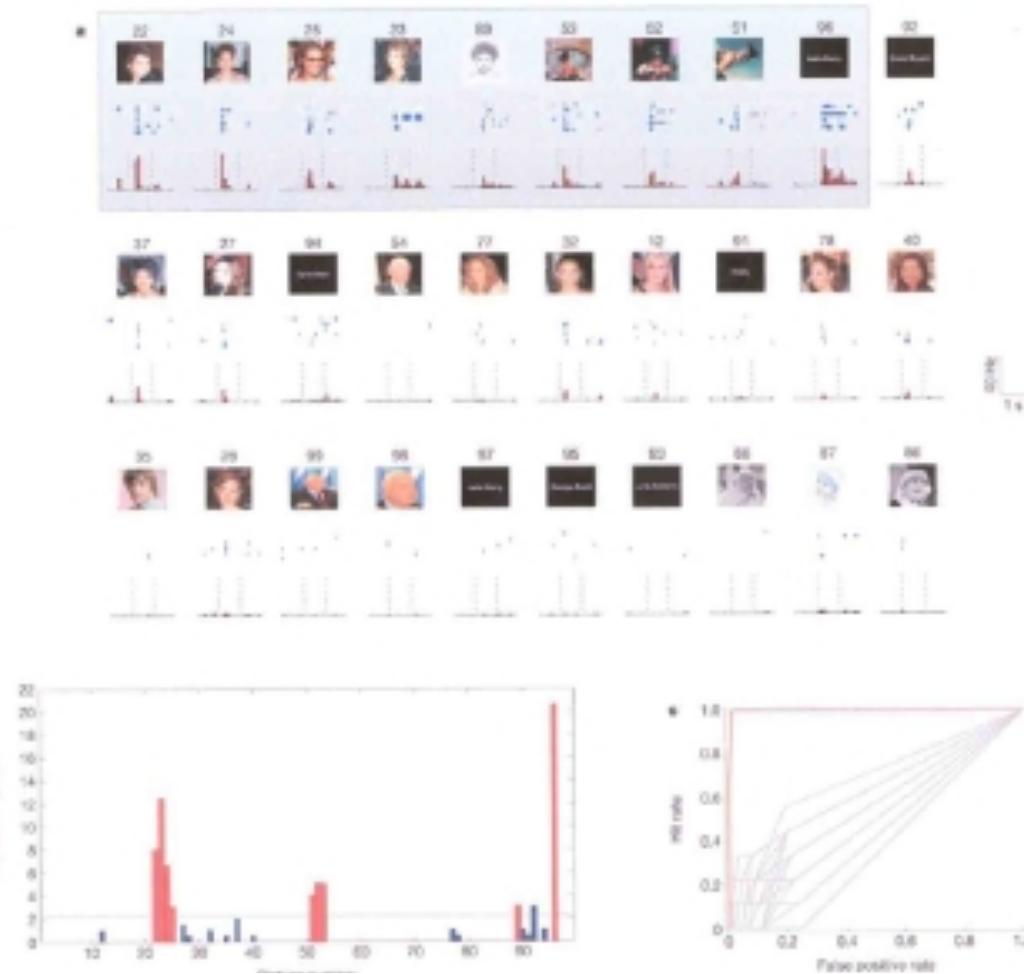
# One way do neuroscientists do it (of many)

- “Halle Berry” Neuron
  - responded to
    - her picture (but not other pictures)
    - a line drawing of her
    - her picture in her Catwoman costume
    - her name (typed out)
- different neurons for
  - Jennifer Aniston
  - Bill Clinton
  - etc.



The “Halle Berry” Neuron?

Quiroga et al. (2005)

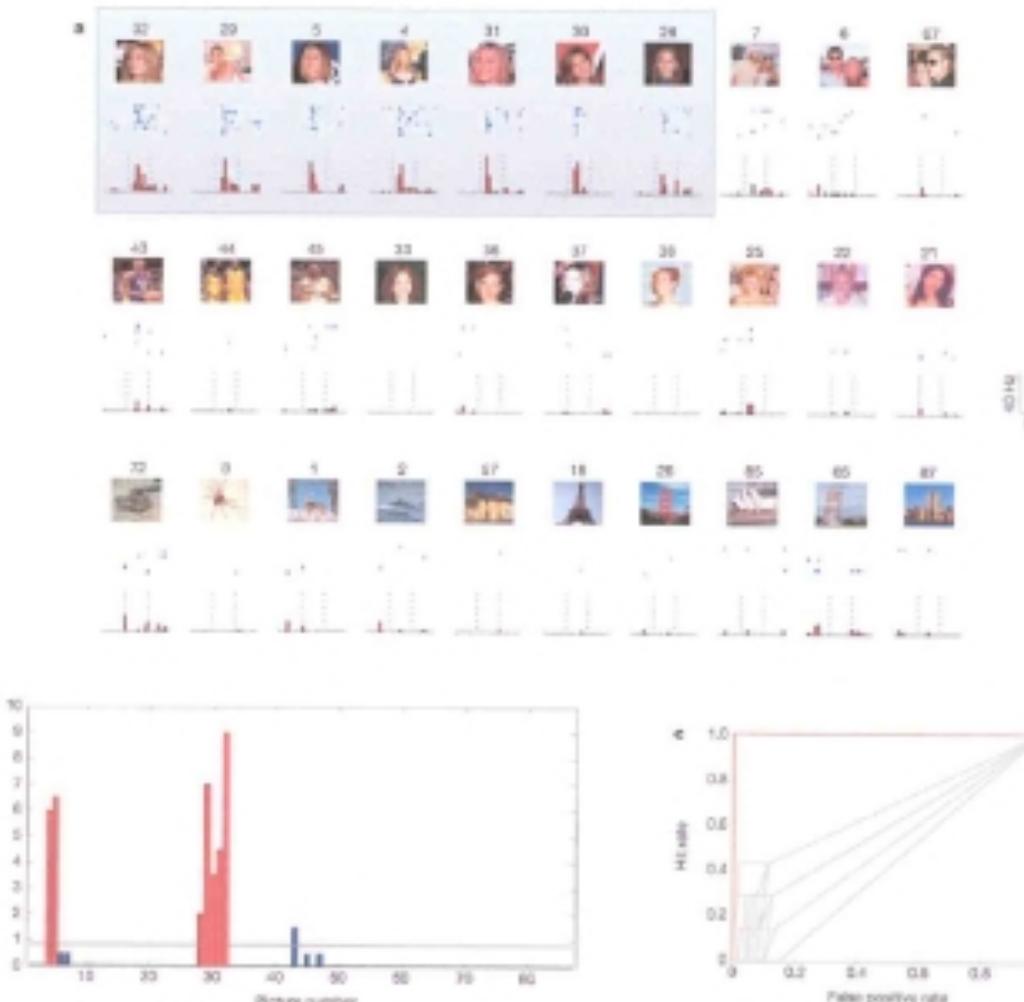


Single Unit  
in right  
Anterior  
Hippocampus

77

The “Jennifer Aniston” Neuron?

Quiroga et al. (2005)

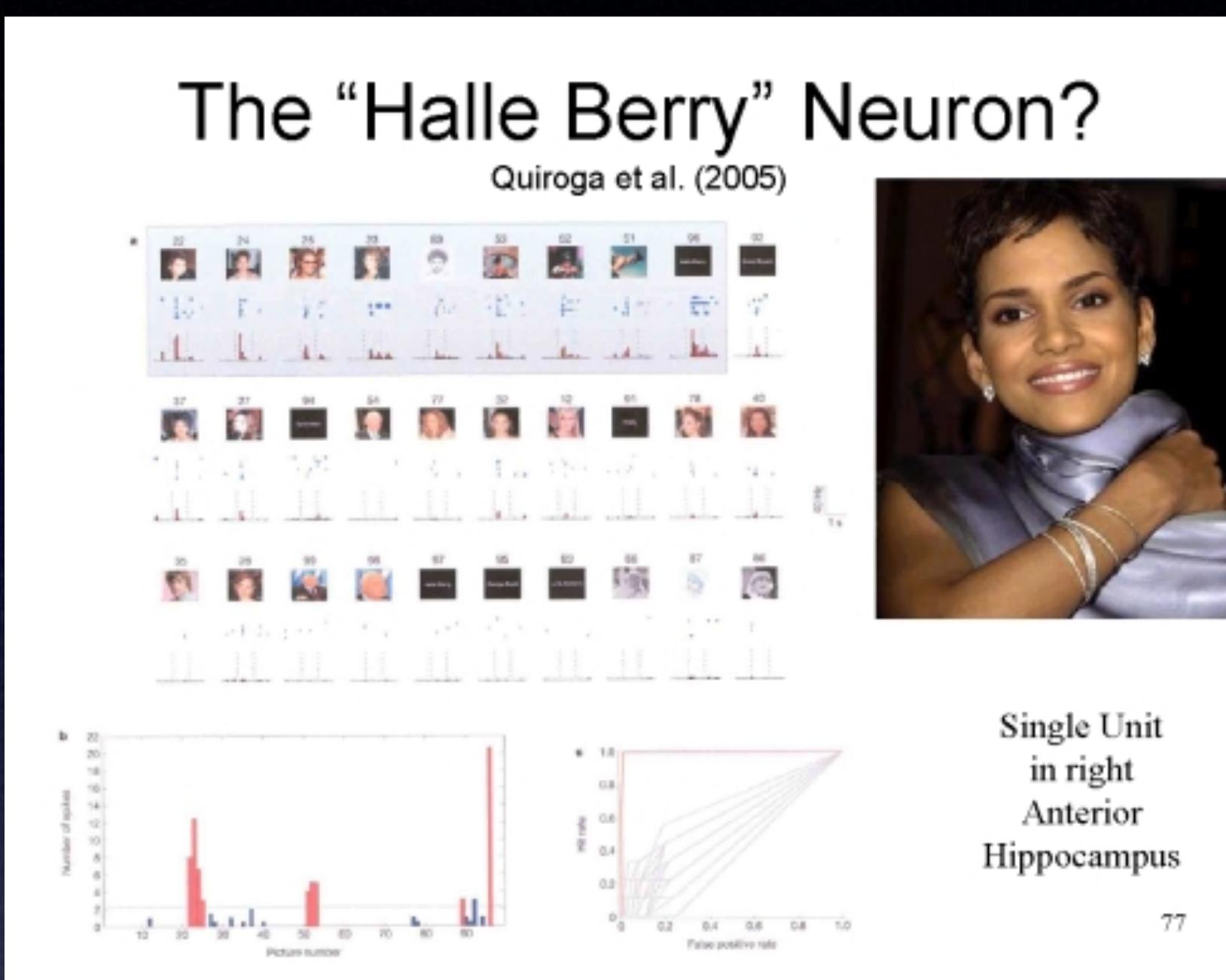


Single Unit  
in Left  
Posterior  
Hippocampus

76

# Open Questions

- Neuroscience
  - Is it really a Halle Berry neuron?
  - or an african-american actress neuron?
  - or?



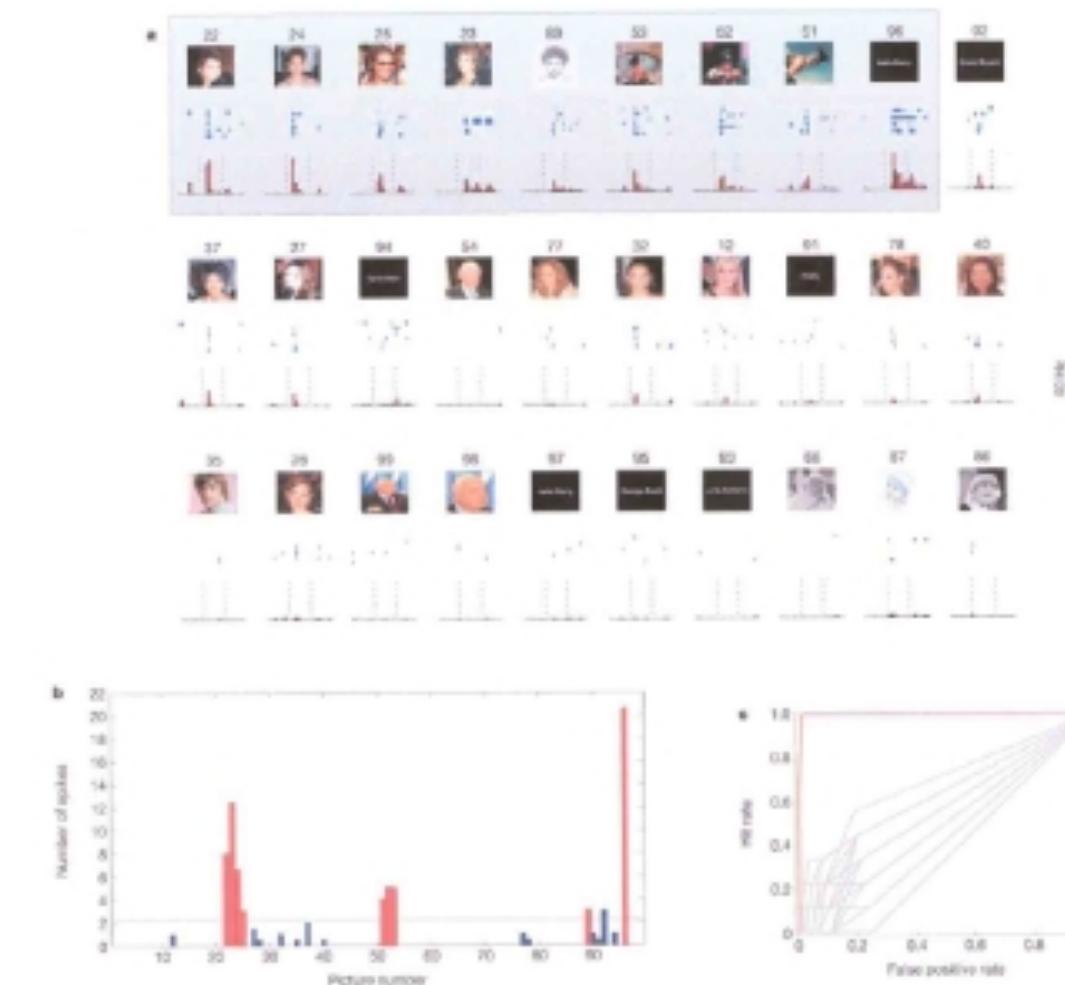
# Open Questions

- Neuroscience
  - Is it really a Halle Berry neuron?
  - or an african-american actress neuron?
  - or?



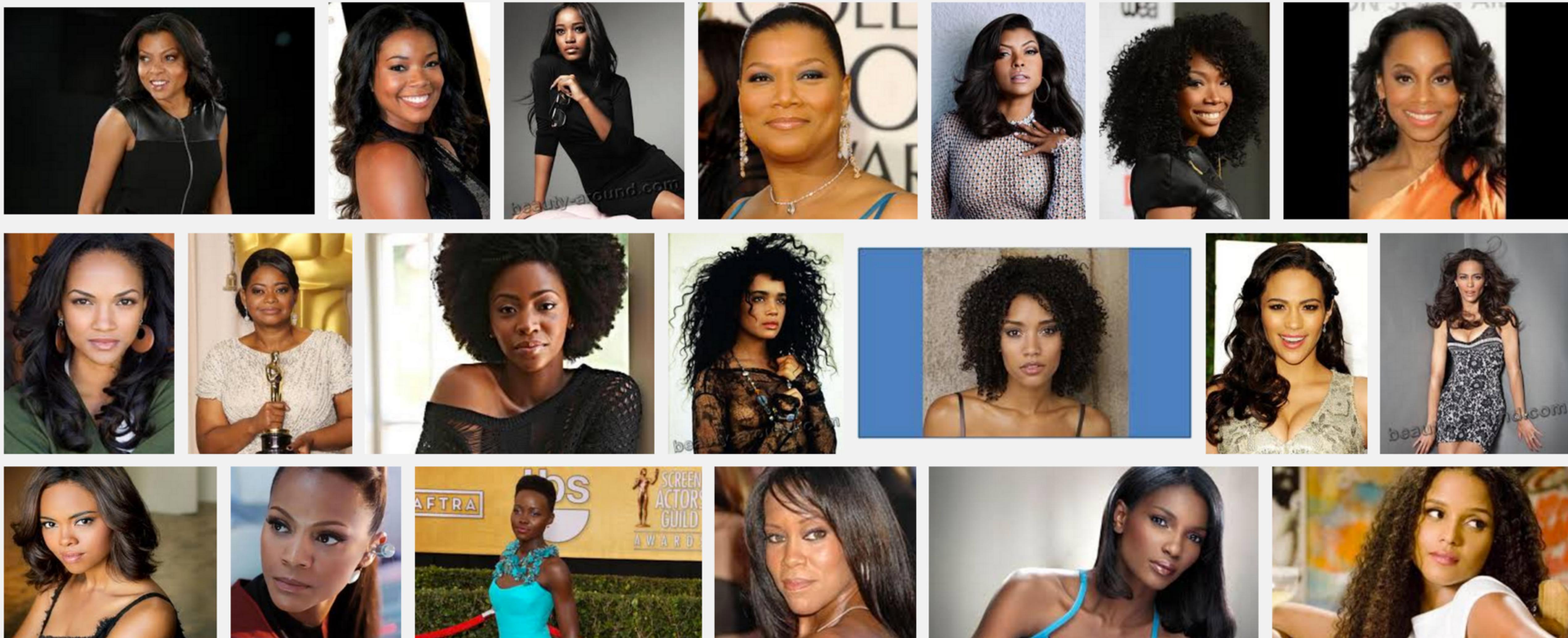
## The “Halle Berry” Neuron?

Quiroga et al. (2005)



Single Unit  
in right  
Anterior  
Hippocampus

# Ideal Test: Synthesize Preferred Inputs



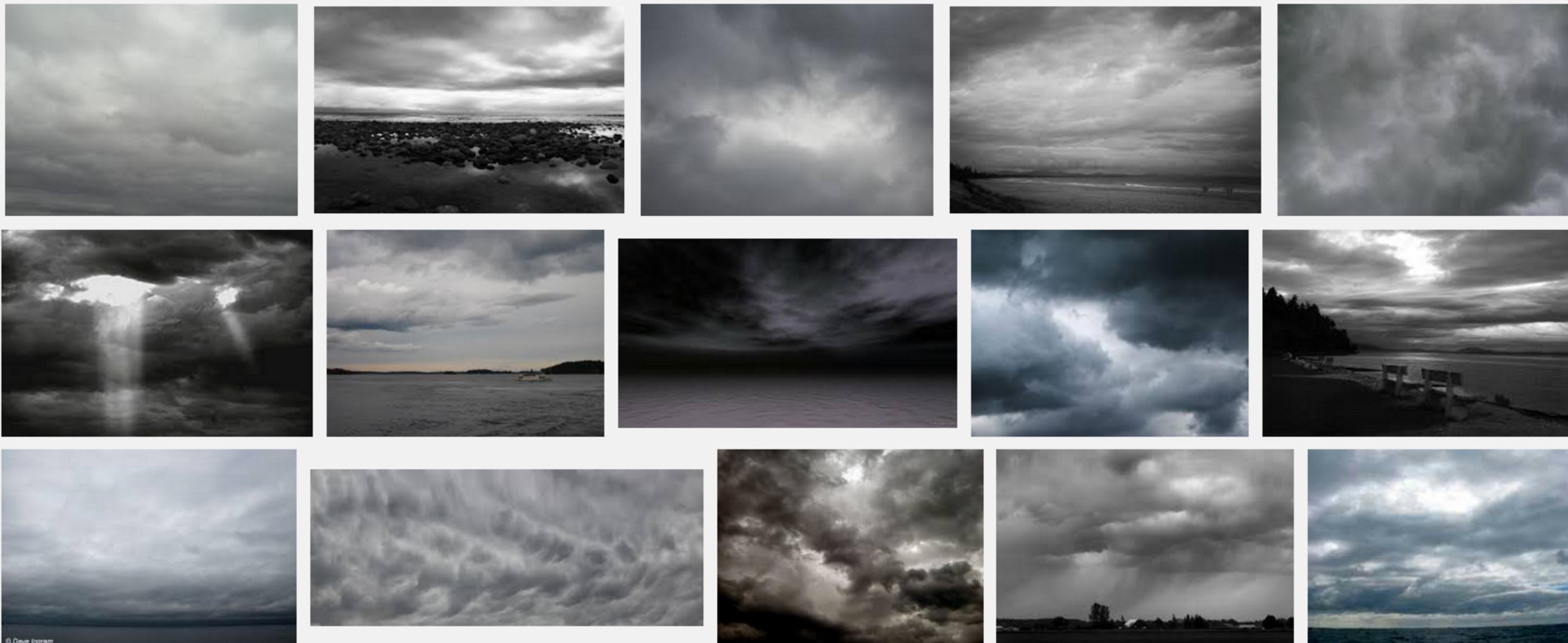
African-american actress neuron

# Ideal Test: Synthesize Preferred Inputs



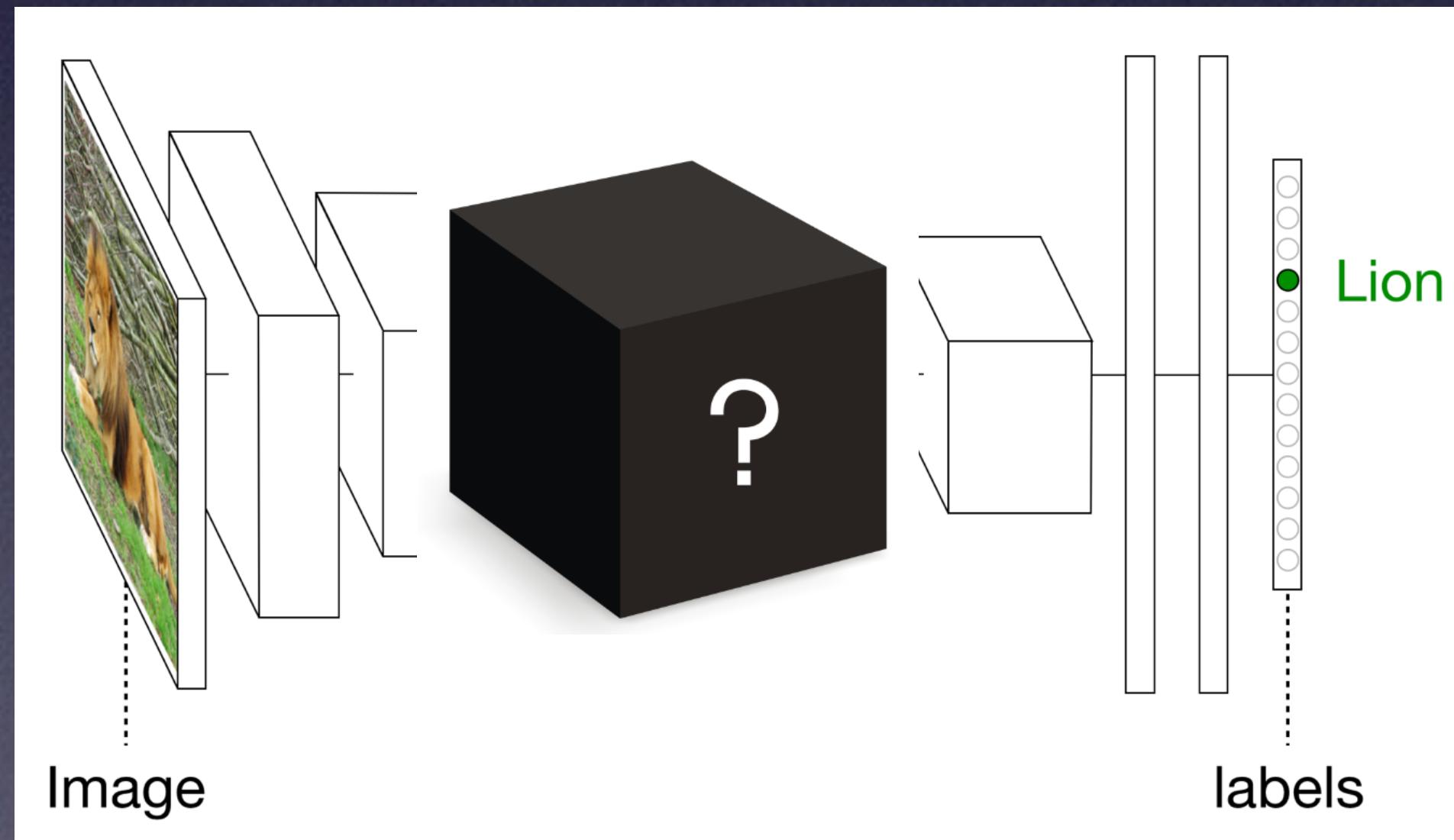
Halle Berry Neuron

# Ideal Test: Synthesize Preferred Inputs

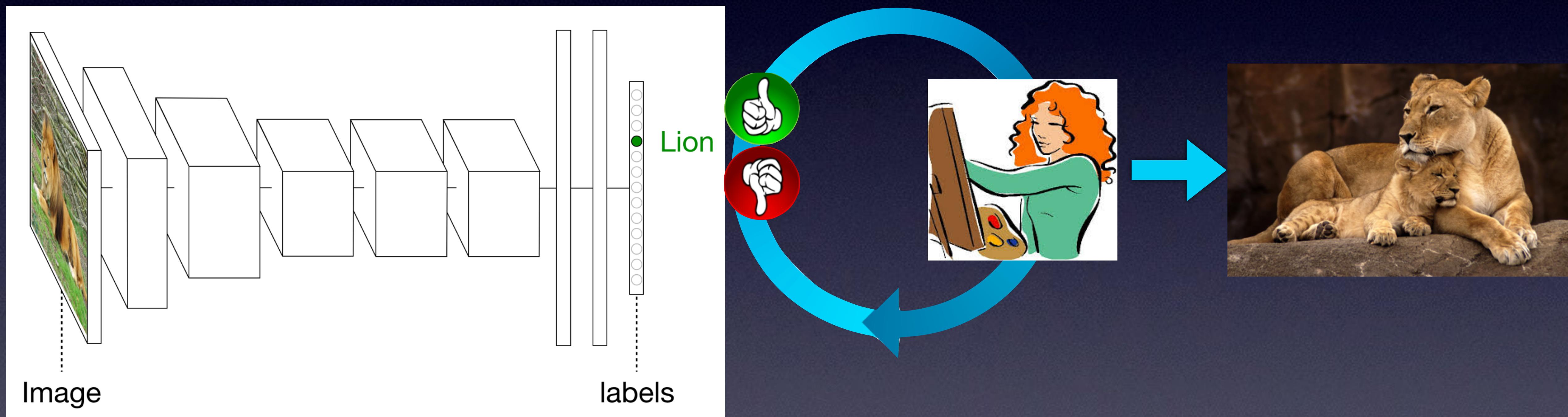


# Possible with Artificial Neural Networks!

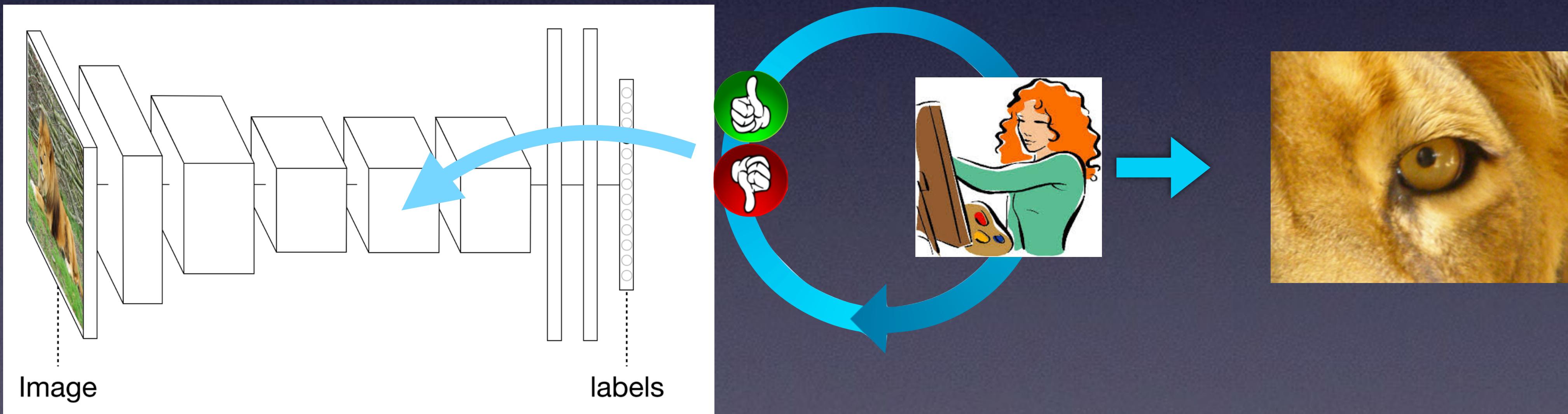
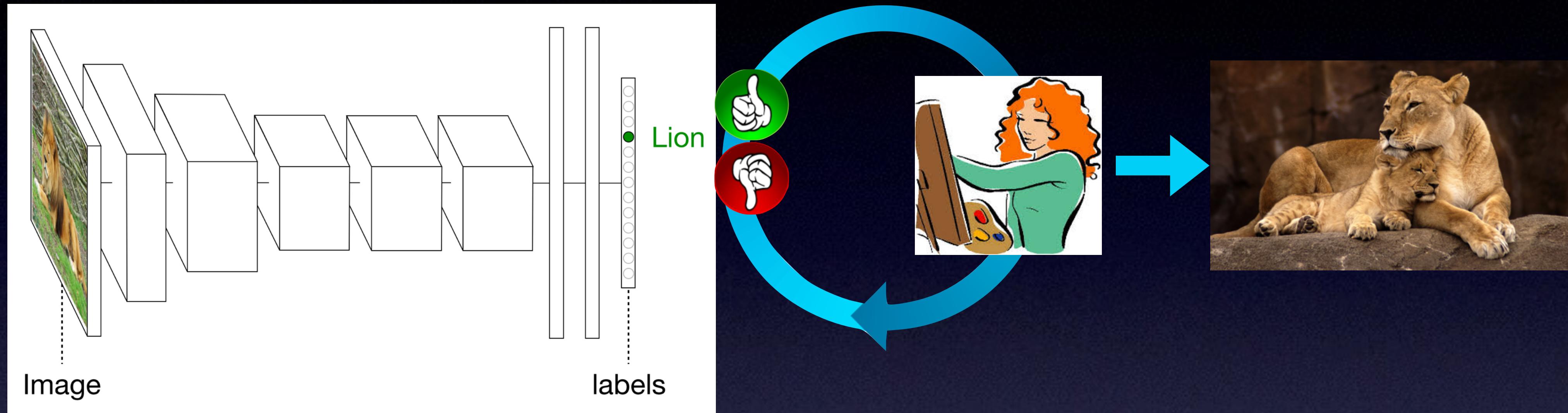
- How much do DNNs understand about the objects they recognize?
- Do they have increasingly abstract feature detectors?
  - Are they “multifaceted”
- Are representations “local” or “distributed”?
- What sorts of errors are they likely to make?



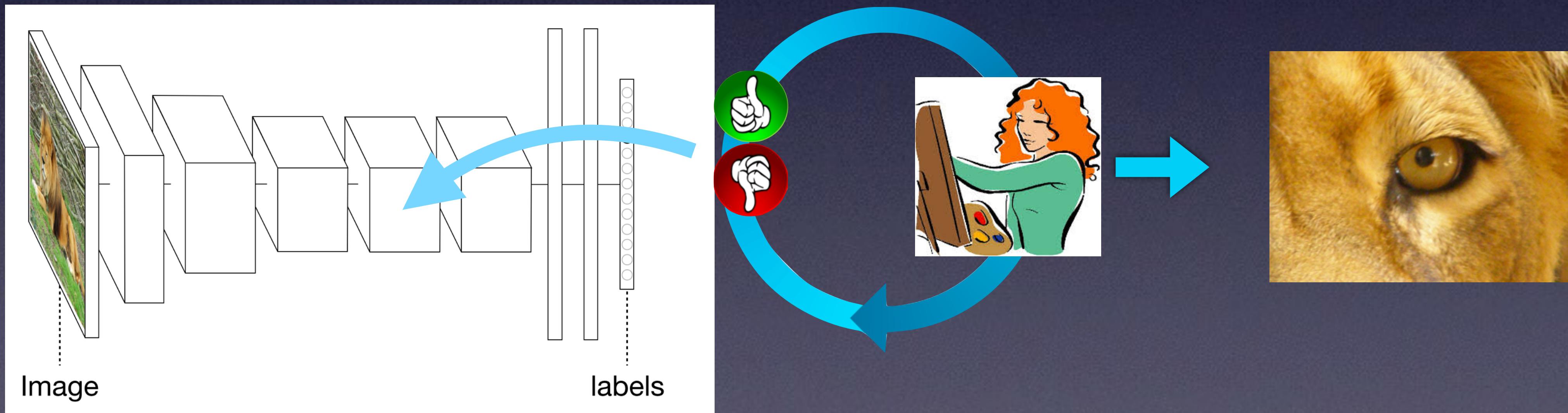
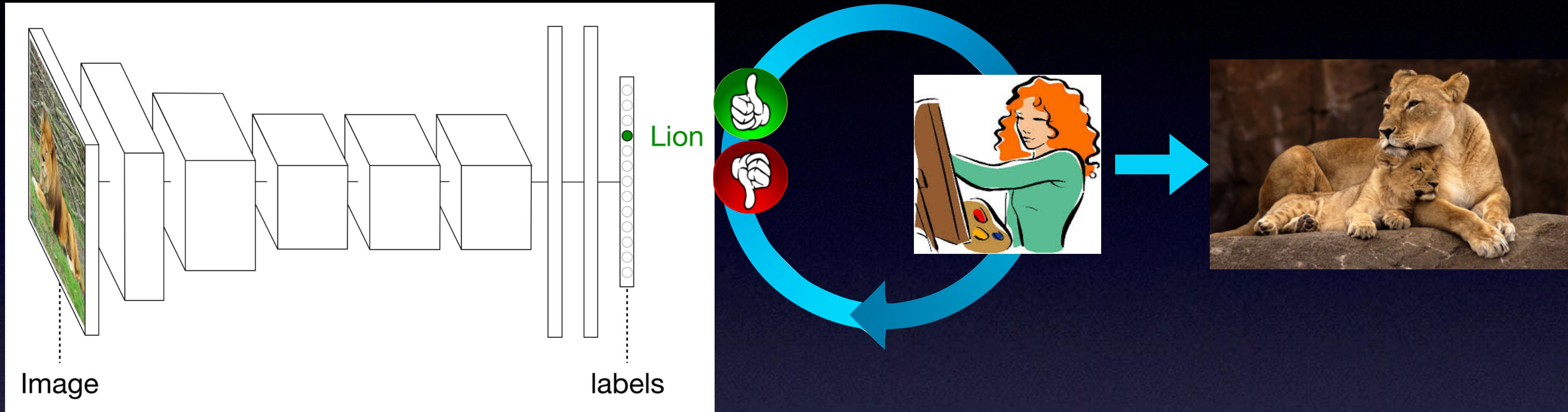
# Investigating What Each Neuron Does



# Investigating What Each Neuron Does

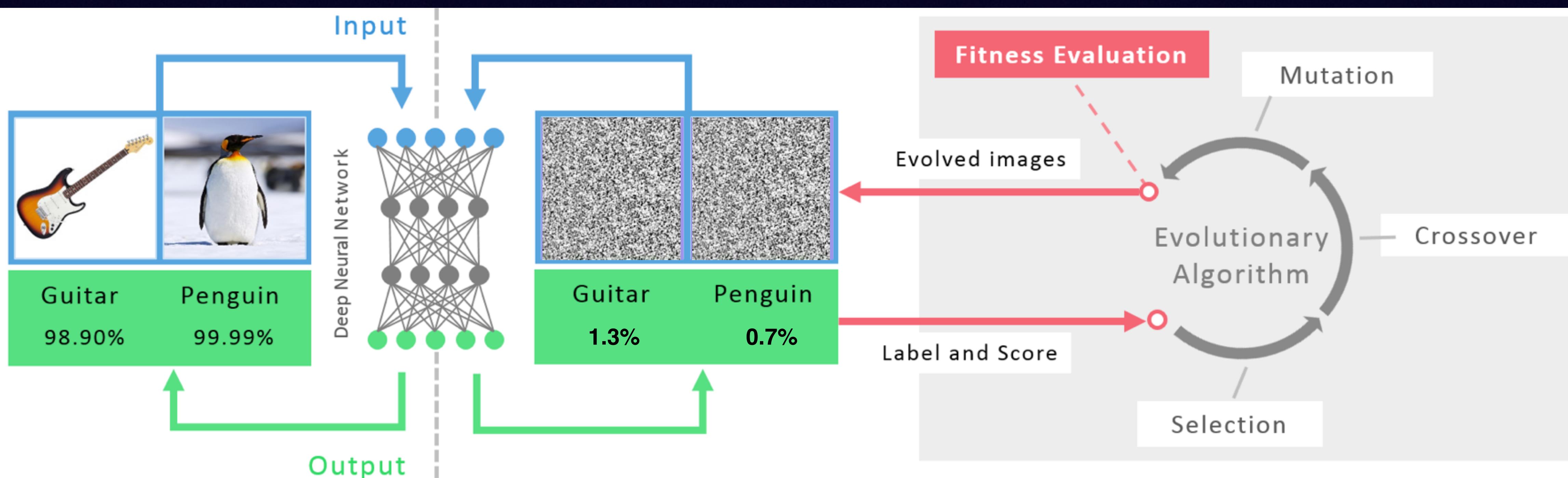


# “Deep Visualization”



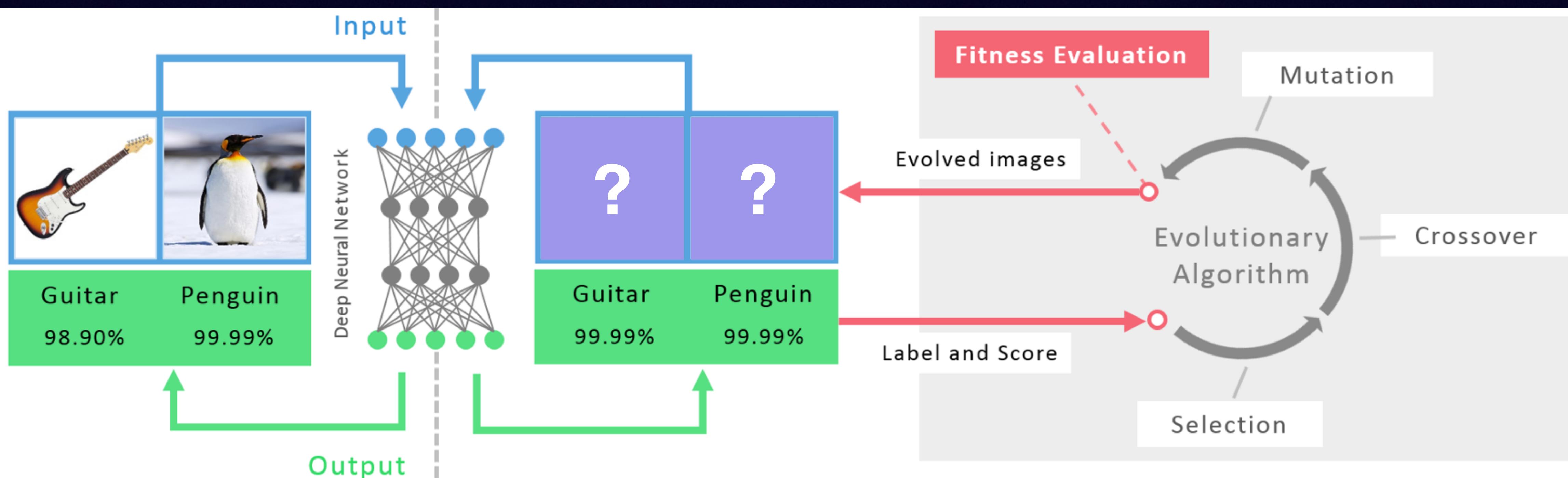
# Deep Visualization Take 1

Nguyen, Yosinski, Clune, 2015, CVPR



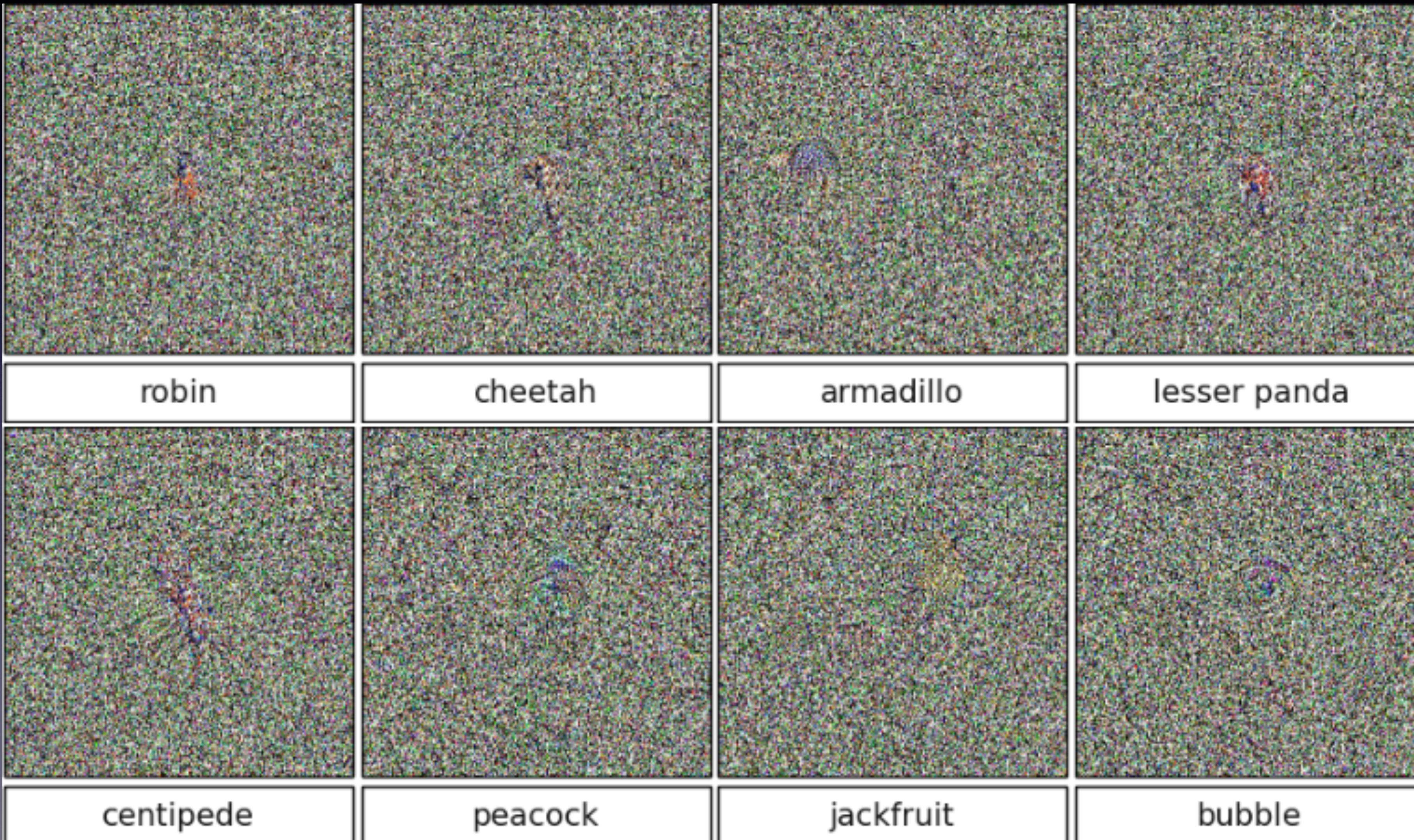
# Deep Visualization Take 1

Nguyen, Yosinski, Clune, 2015, CVPR



# Deep Visualization Take 1

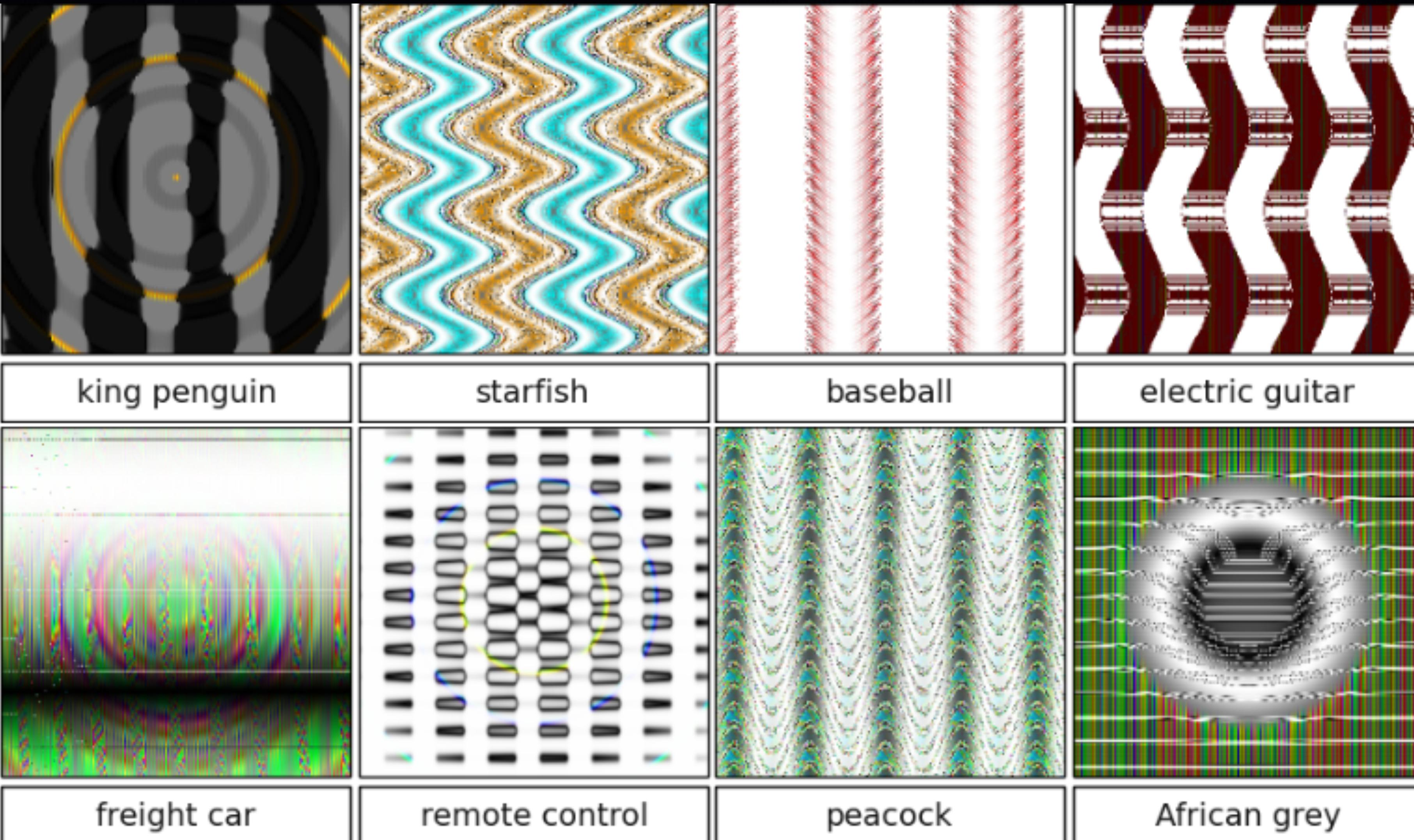
Nguyen, Yosinski, Clune, 2015, CVPR



DNN Confidence: > 99.6 % for all

# Deep Visualization Take 1

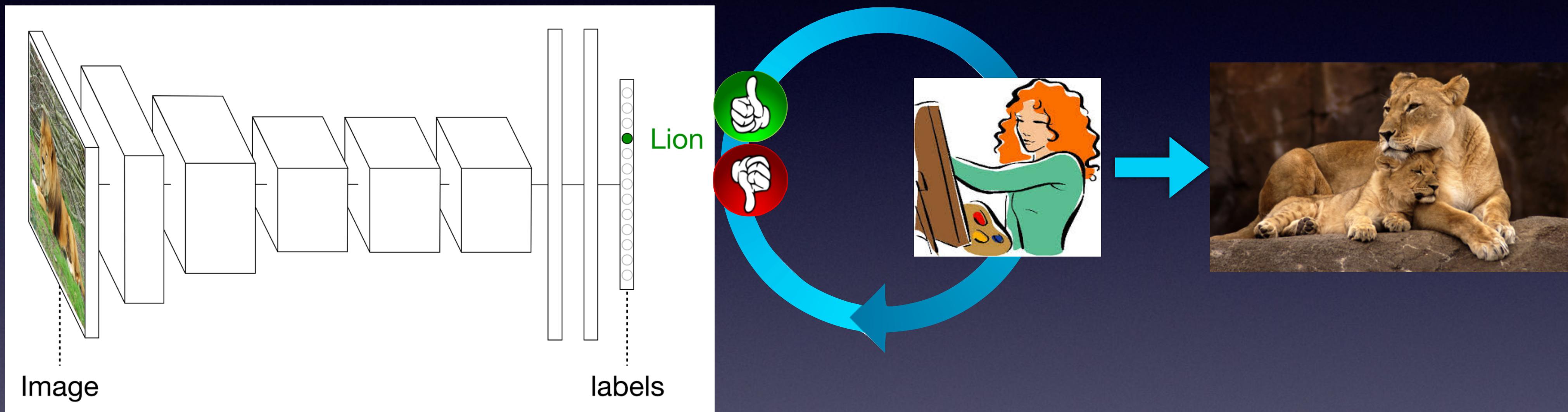
Nguyen, Yosinski, Clune, 2015, CVPR



DNN Confidence: > 99.6 % for all

# Deep Visualization Take 1

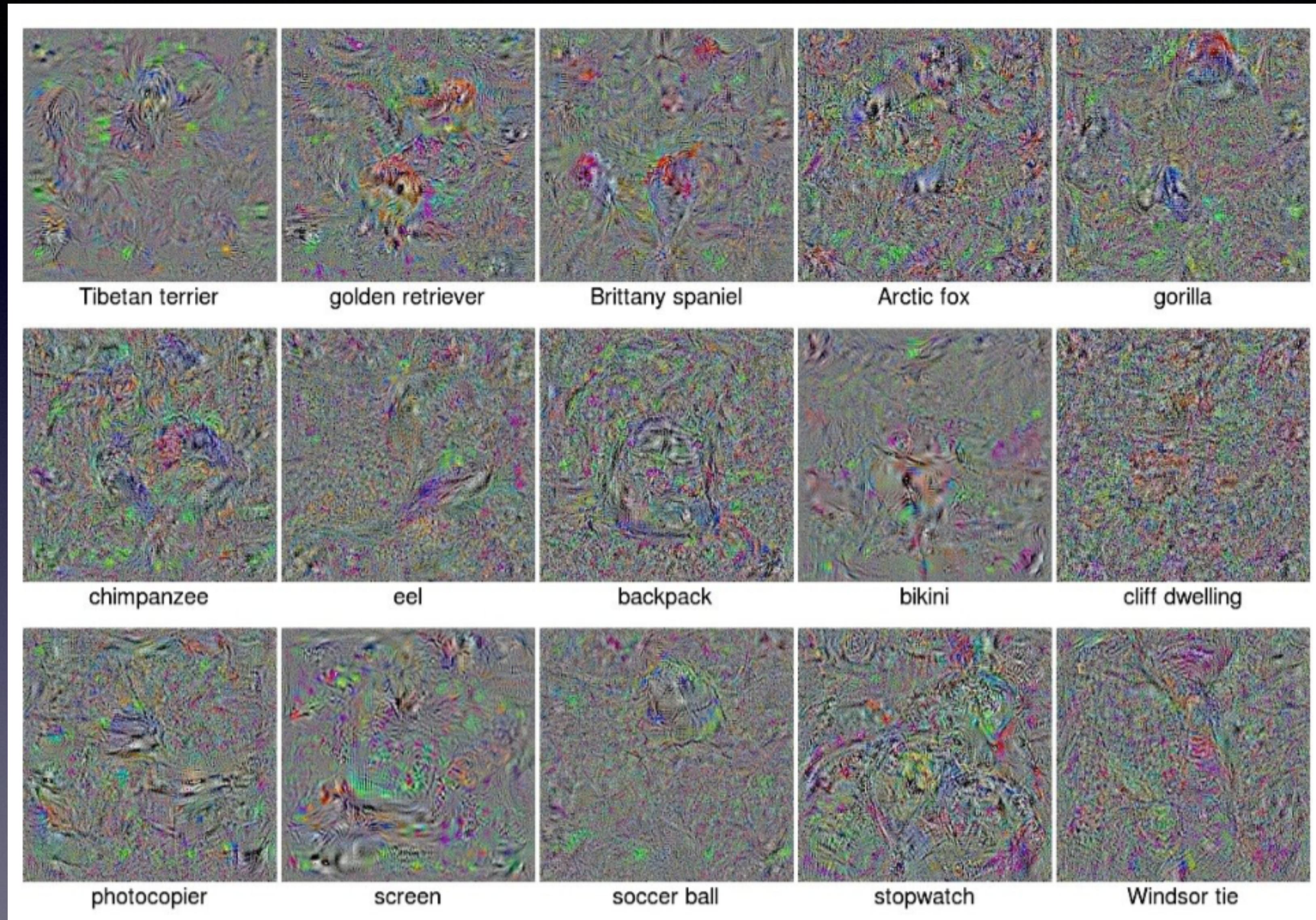
Nguyen, Yosinski, Clune, 2015, CVPR



Backprop to  
Pixel Changes

# Deep Visualization Take 1

Nguyen, Yosinski, Clune, 2015, CVPR



DNN Confidence: > 99.9 % for all

# Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images

Anh Nguyen  
University of Wyoming  
anguyen8@uwyko.edu

Jason Yosinski  
Cornell University  
yosinski@cs.cornell.edu

Jeff Clune  
University of Wyoming  
jeffclune@uwyko.edu

## Abstract

Deep neural networks (DNNs) have recently been achieving state-of-the-art performance on a variety of pattern-recognition tasks, most notably visual classification problems. Given that DNNs are now able to classify objects in images with near-human-level performance, questions naturally arise as to what differences remain between computer and human vision. A recent study [30] revealed that changing an image (e.g. of a lion) in a way imperceptible to humans can cause a DNN to label the image as something else entirely (e.g. mislabeling a lion a library). Here we show a related result: it is easy to produce images that are completely unrecognizable to humans, but that state-of-the-art DNNs believe to be recognizable objects with 99.99% confidence (e.g. labeling with certainty that white noise static is a lion). Specifically, we take convolutional neural networks trained to perform well on either the ImageNet or MNIST datasets and then find images with evolutionary algorithms or gradient ascent that DNNs label with high confidence as belonging to each dataset class. It is possible to produce images totally unrecognizable to human eyes that DNNs believe with near certainty are familiar objects, which we call “fooling images” (more generally, fooling examples). Our results shed light on interesting differences between human vision and current DNNs, and raise questions about the generality of DNN computer vision.

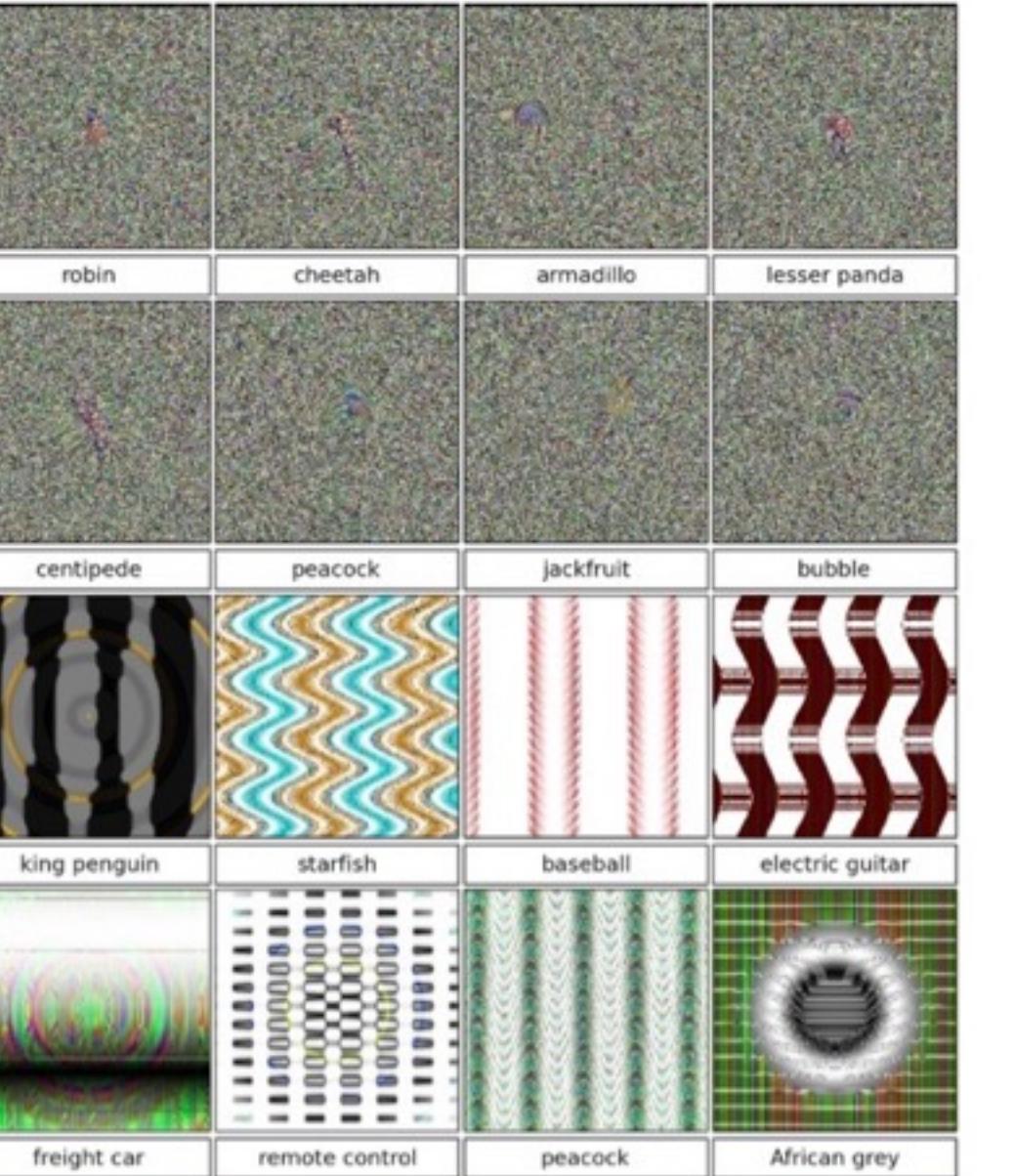


Figure 1. Evolved images that are unrecognizable to humans, but that state-of-the-art DNNs trained on ImageNet believe with  $\geq 99.6\%$  certainty to be a familiar object. This result highlights differences between how DNNs and humans recognize objects. Images are either directly (*top*) or indirectly (*bottom*) encoded.

## 1. Introduction

Deep neural networks (DNNs) learn hierarchical layers of representation from sensory input in order to perform pattern recognition [2, 14]. Recently, these deep architectures have demonstrated impressive, state-of-the-art, and sometimes human-competitive results on many pattern recognition tasks, especially vision classification problems [16, 7, 31, 17]. Given the near-human ability of DNNs to classify visual objects, questions arise as to what differences remain between computer and human vision.

A recent study revealed a major difference between DNN and human vision [30]. Changing an image, originally correctly classified (e.g. as a lion), in a way imperceptible to human eyes, can cause a DNN to label the image as something else entirely (e.g. mislabeling a lion a library).

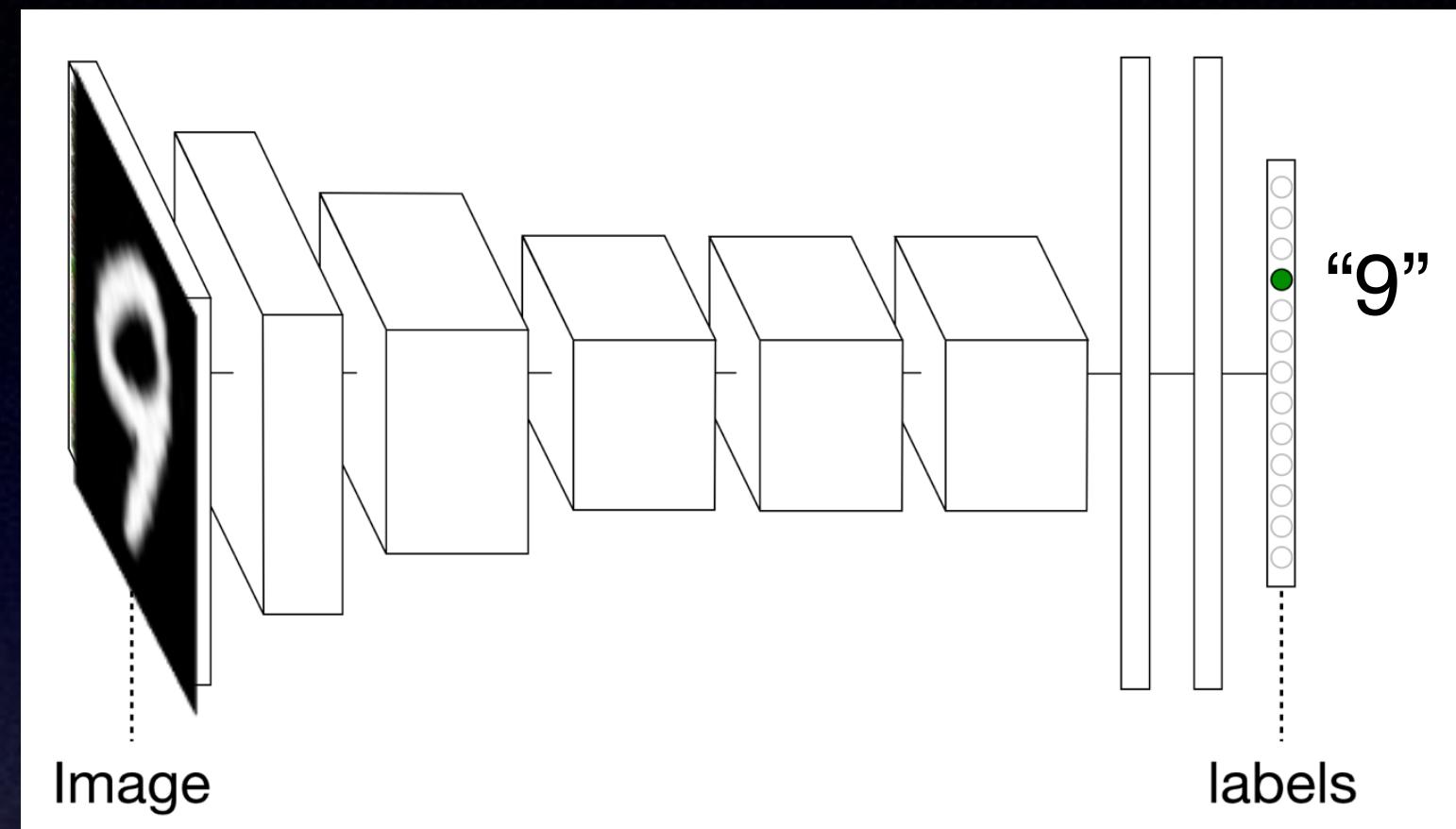
In this paper, we show another way that DNN and human vision differ: It is easy to produce images that are completely unrecognizable to humans (Fig. 1), but that state-of-the-art DNNs believe to be recognizable objects with over 99% confidence (e.g. labeling with certainty that TV static

- AI “sees” the world differently
- Clever Hans?
- Huge security concern



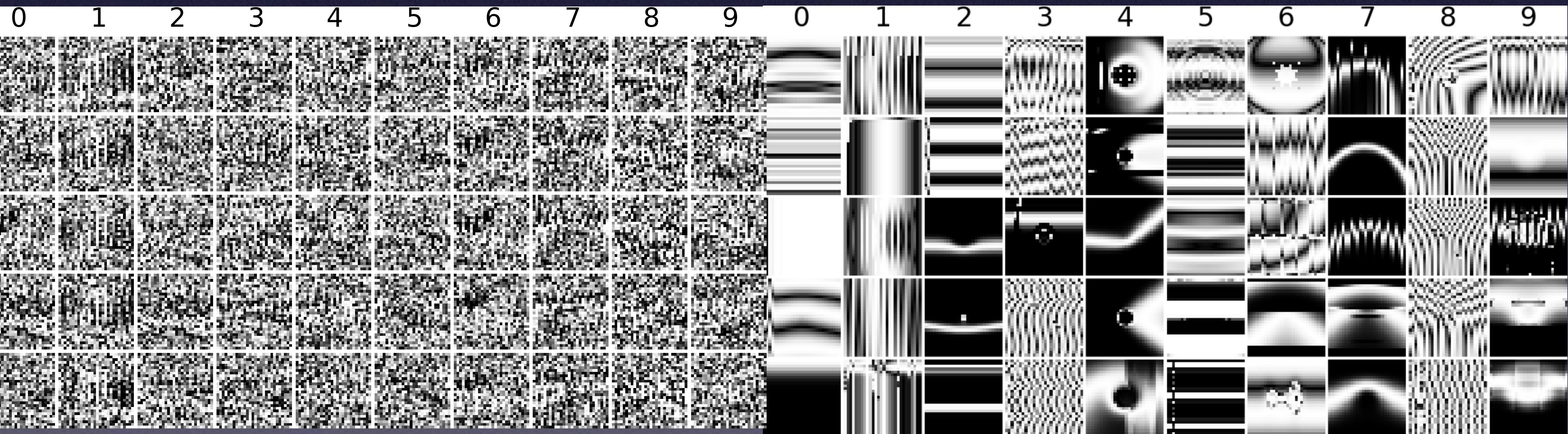
School bus

# Digits

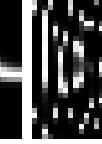
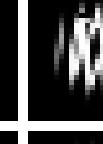
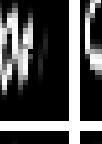
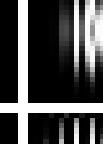
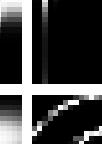
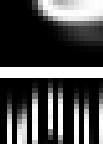
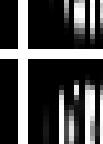
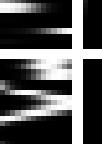
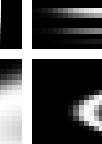
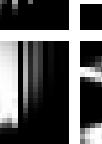
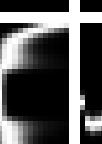
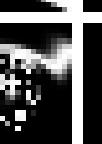
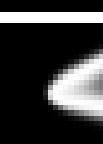
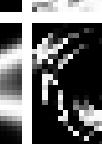
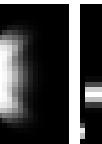
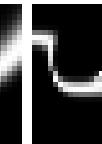
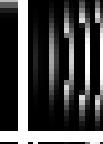
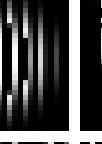
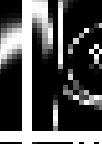
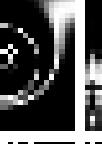
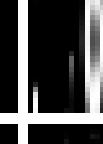
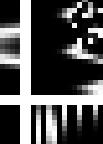
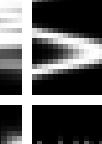
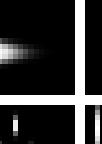
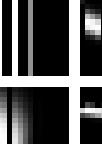


0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9  
0 1 2 3 4 5 6 7 8 9

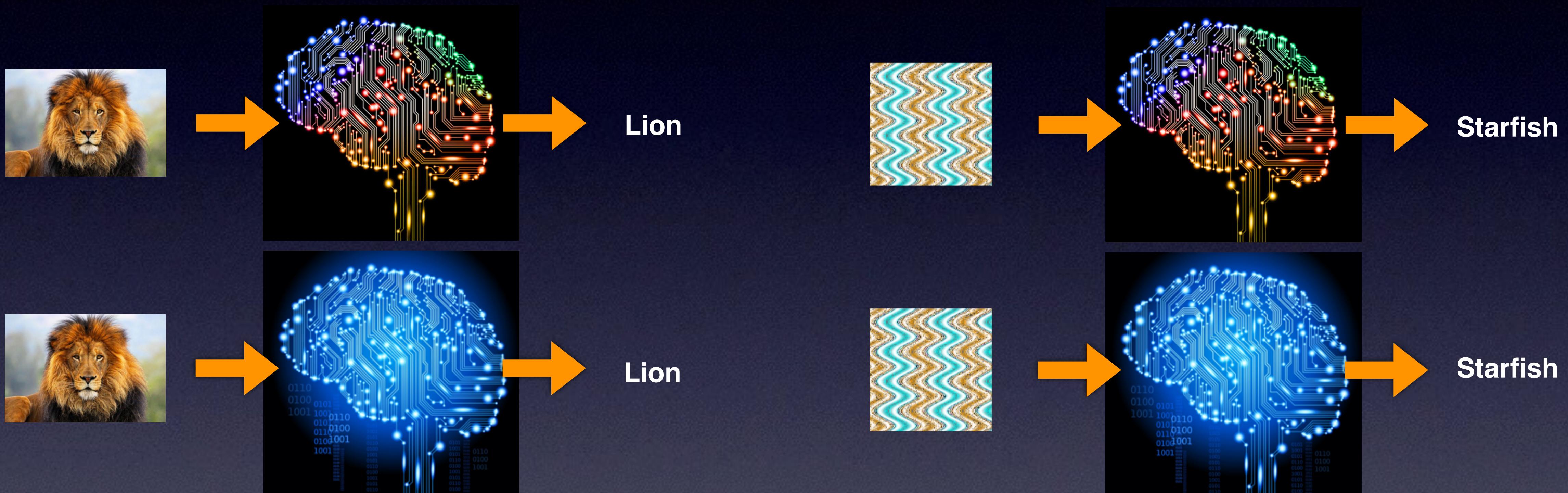
>99% accurate



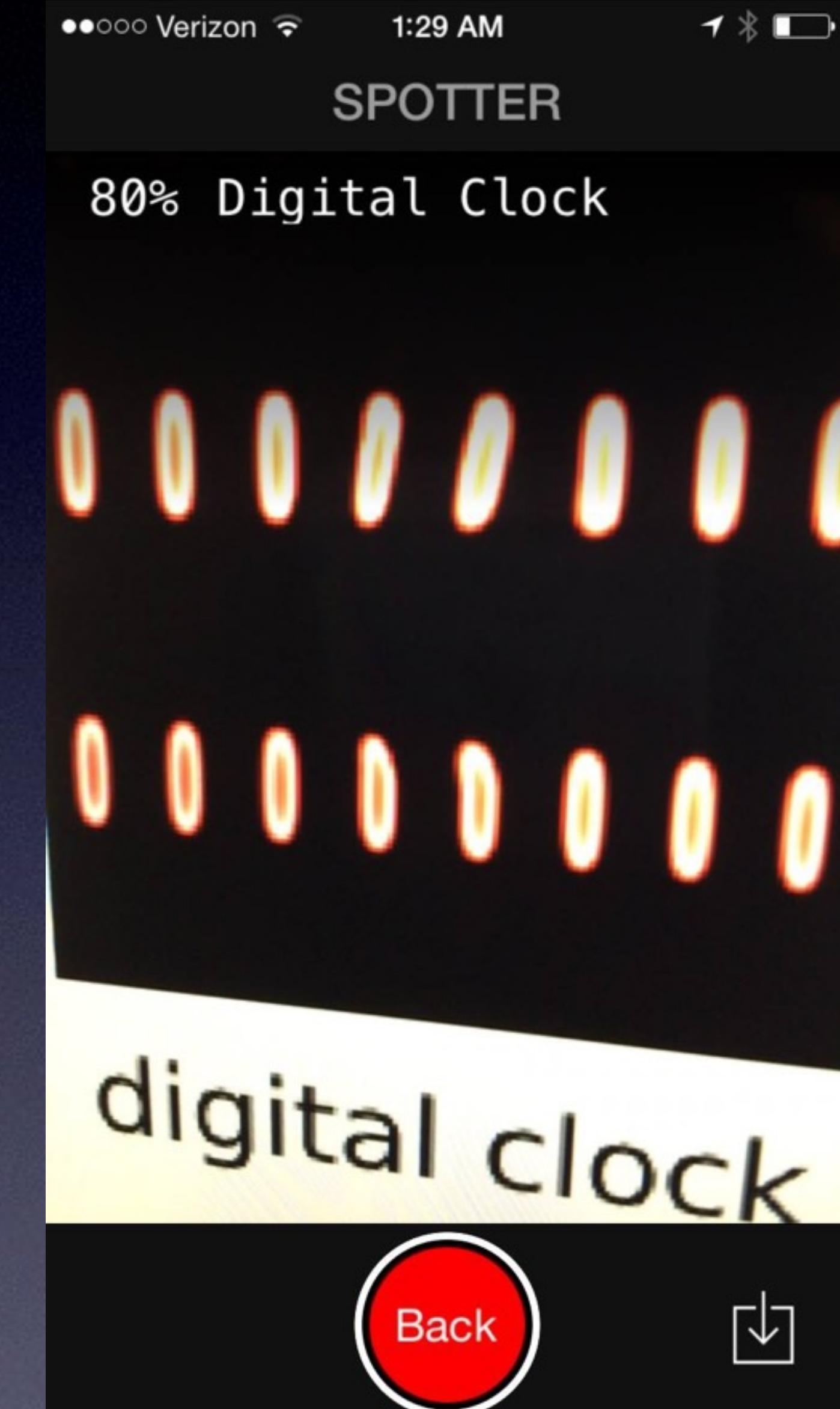
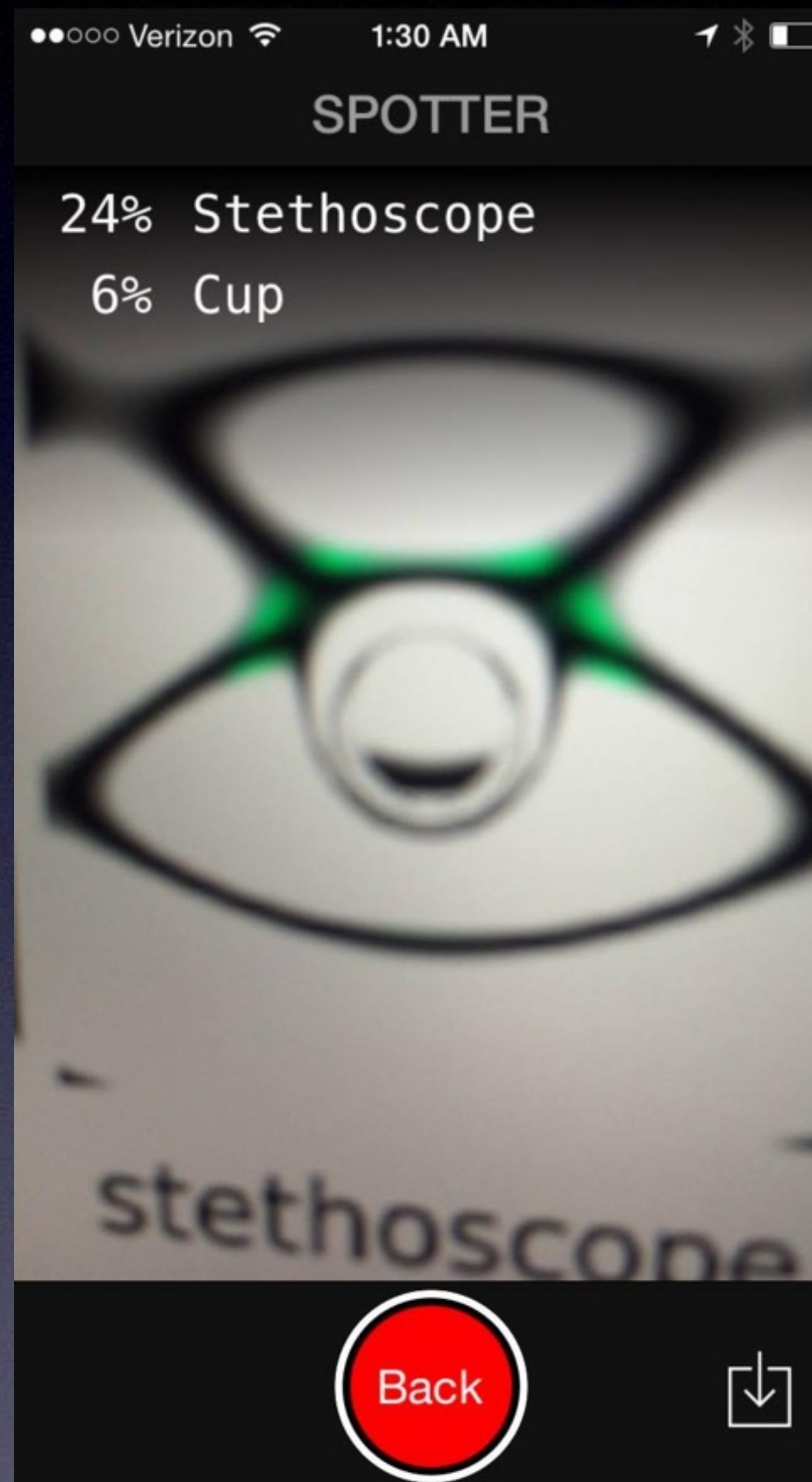
# Training on these examples does not help

	0	1	2	3	4	5	6	7	8	9	Median confidence
1											99.99
2											97.42
3											99.83
4											72.52
5											97.55
6											99.68
7											76.13
8											99.96
9											99.51
10											99.48
11											98.62
12											99.97
13											99.93
14											99.15
15											99.15

# Images that fool one network fool others!



# Images that fool one network fool others!



# Huge reaction

63rd most talked about scientific paper worldwide in 2015 - Altmetric





Larry Page,  
Google co-Founder

Gary Marcus,  
NYU Prof & CEO



Don't worry killer robot,  
I'm really a starfish.

# Well-received Academically

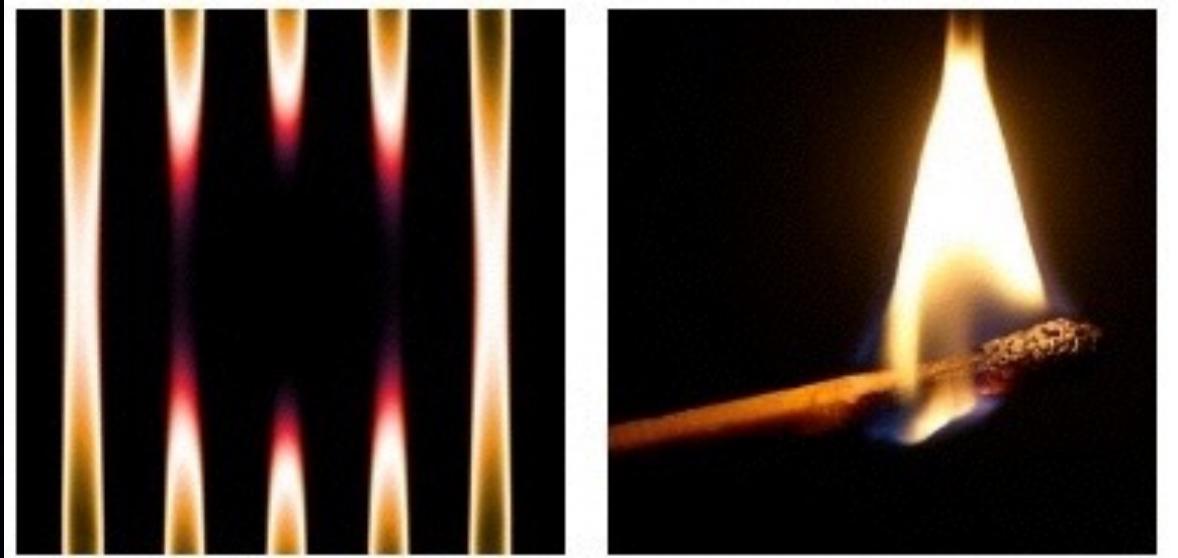
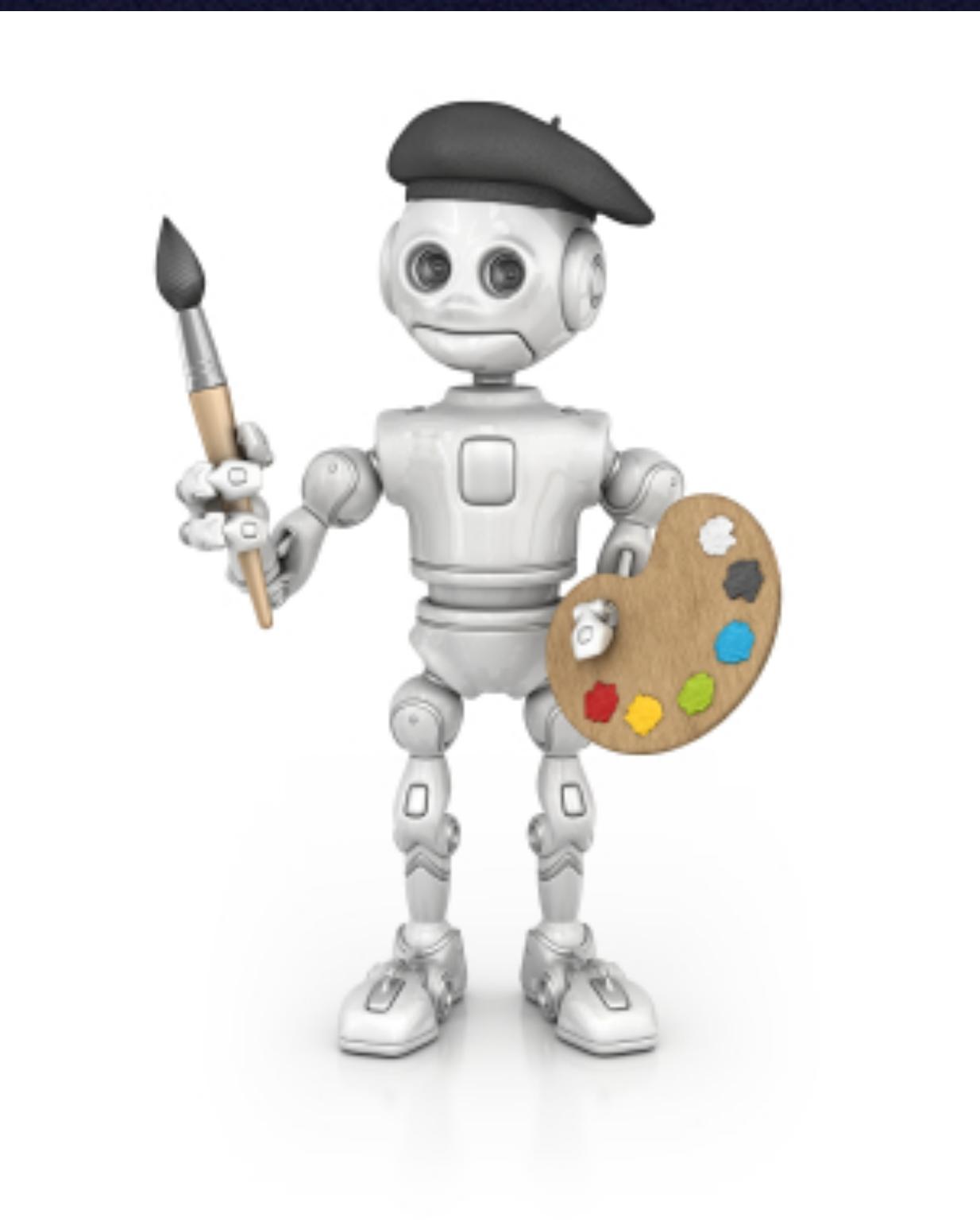
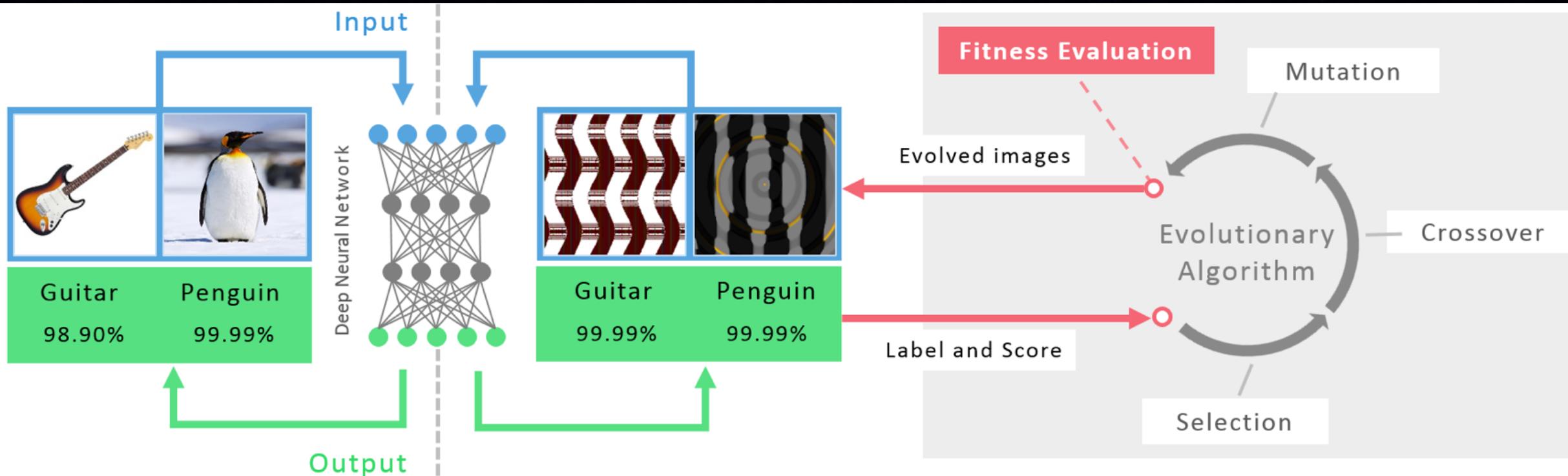
- Computer Vision and Pattern Recognition (CVPR)
  - Oral presentation (3% acceptance rate)
  - Winner: Community Top Paper Award (3rd place)



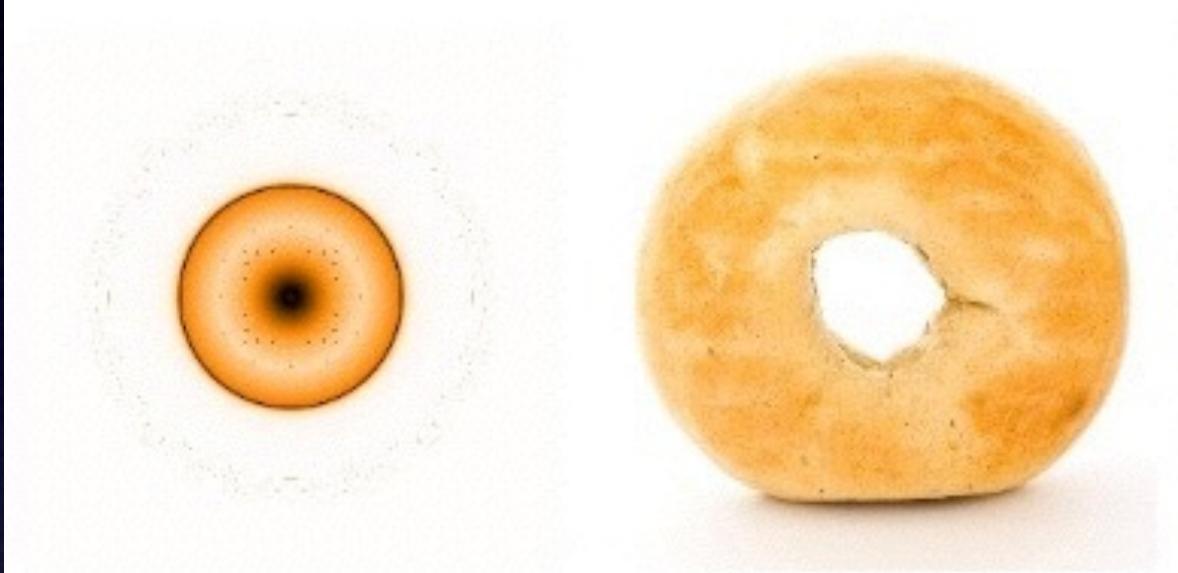
# What does the huge reaction mean?

- A lot of interest in whether Deep Learning is “real” intelligence
  - or a “Clever Hans”

# Automatic Art Generator



Matchstick



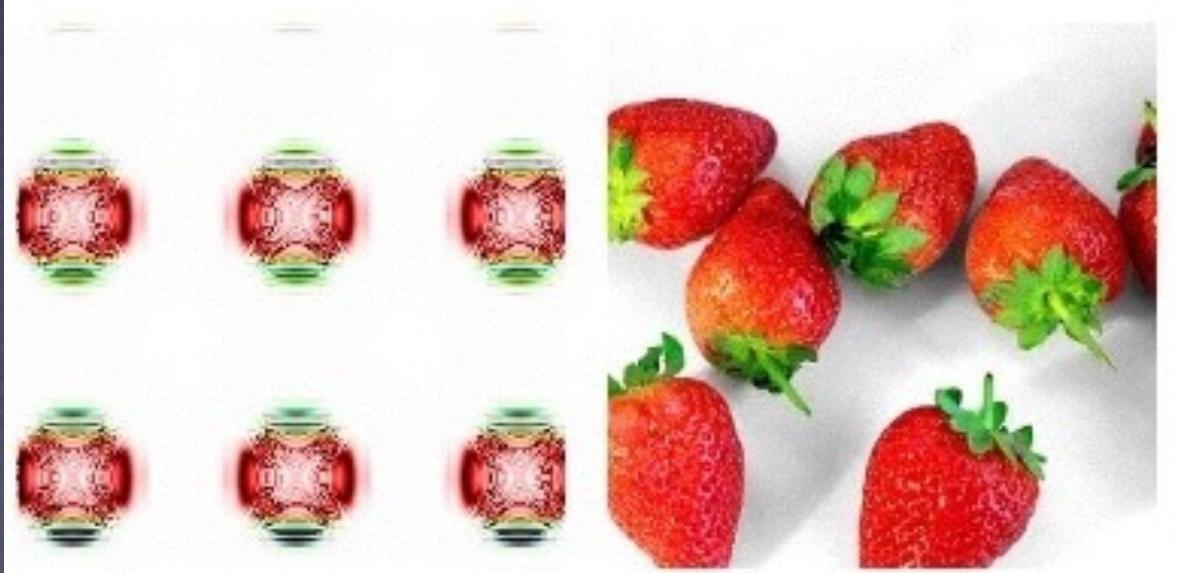
Bagel



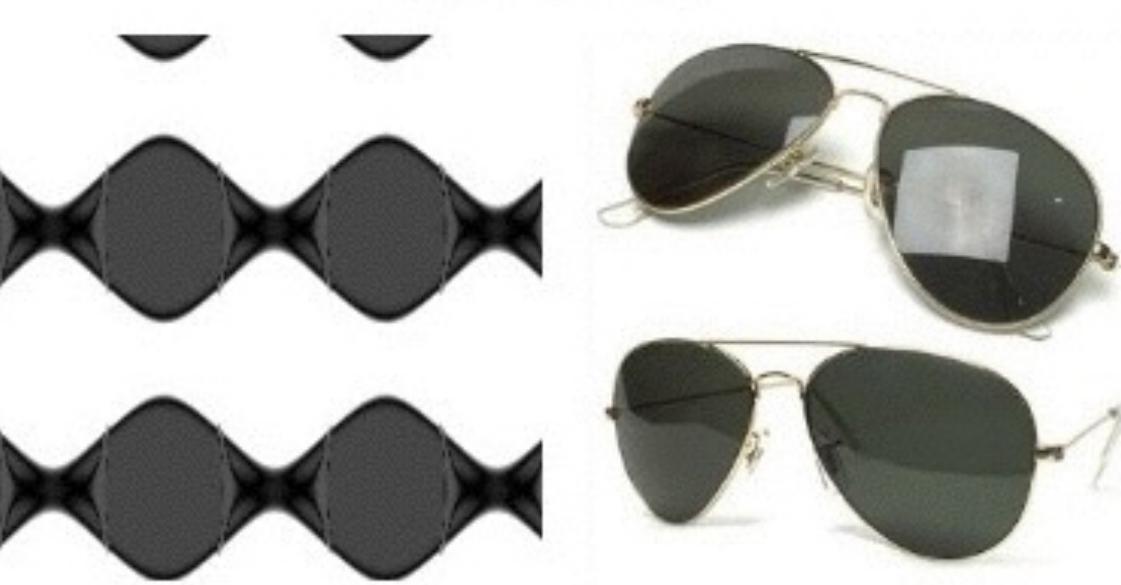
Prison



Tile roof

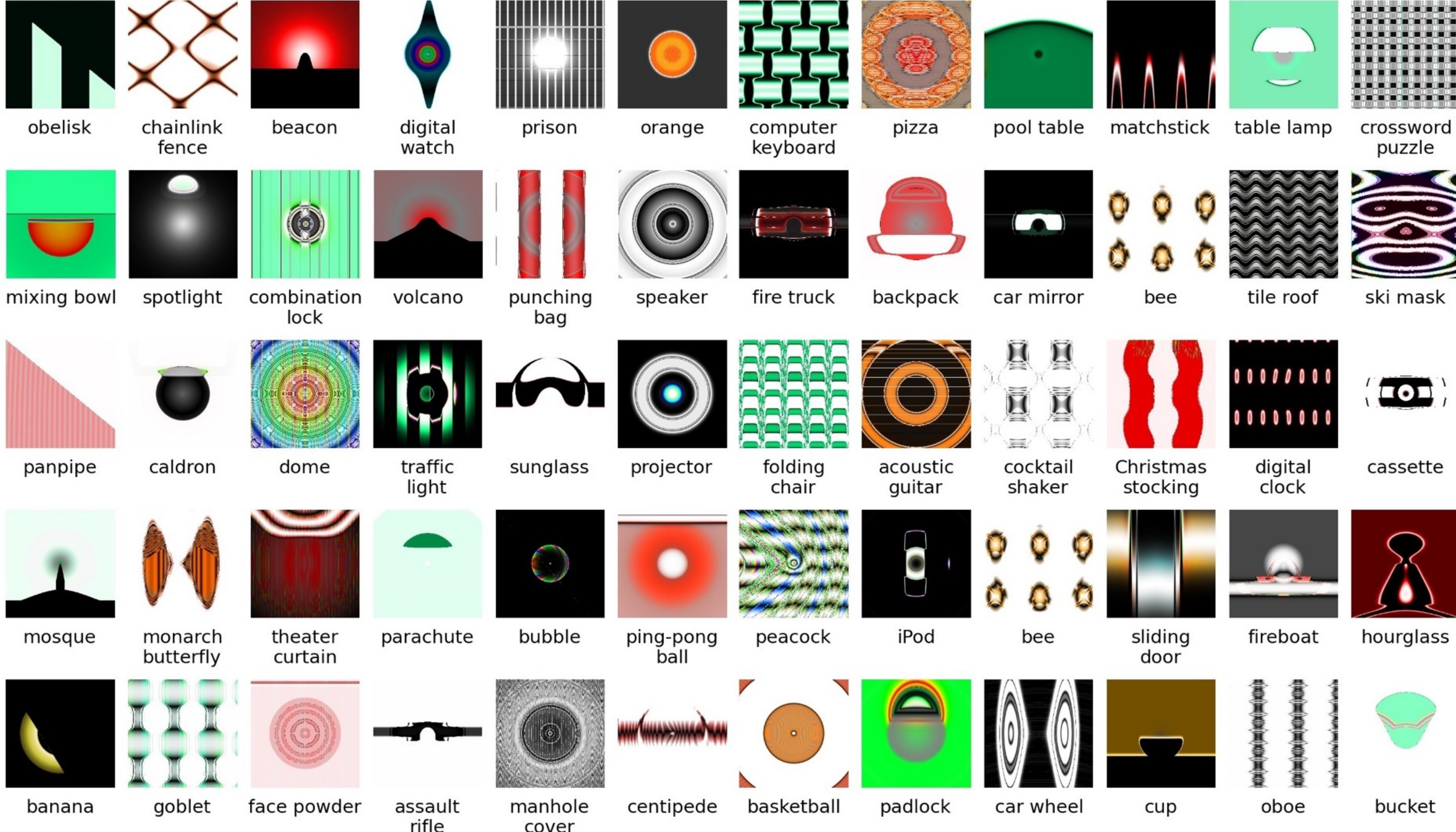


Strawberry



Sunglasses

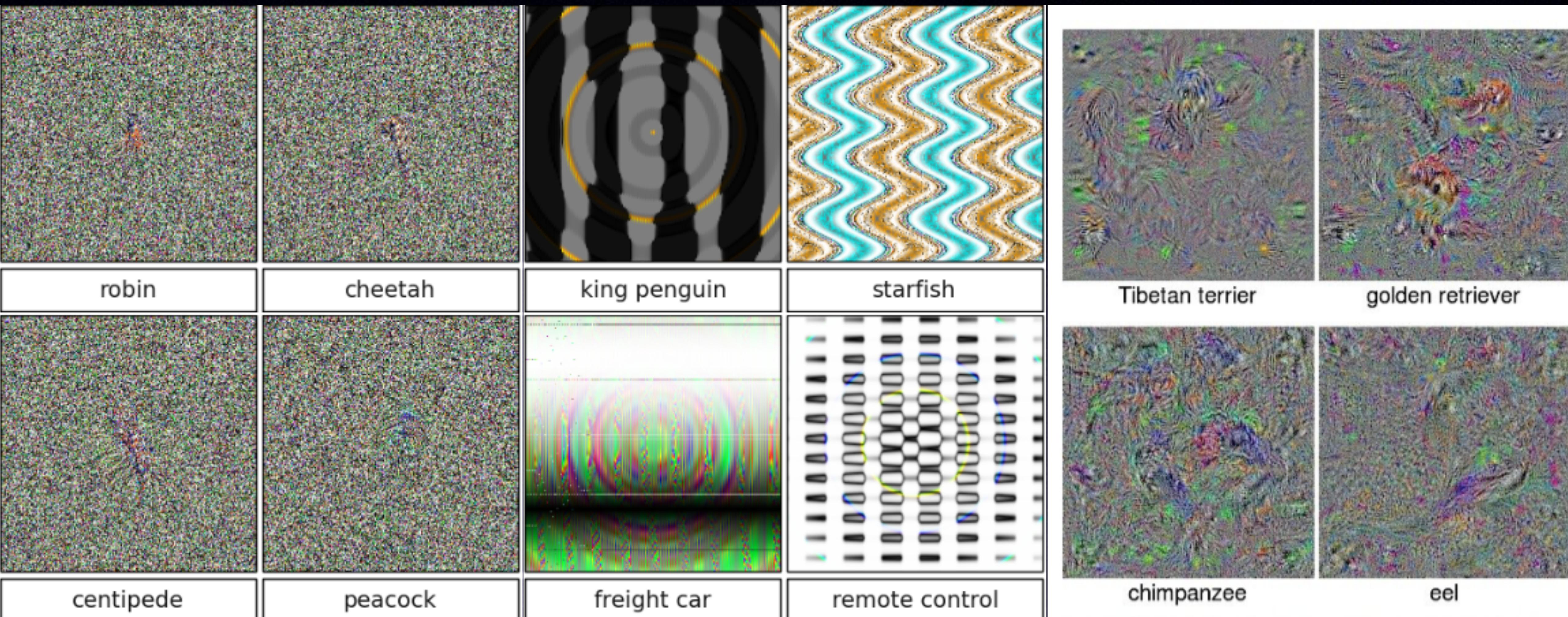






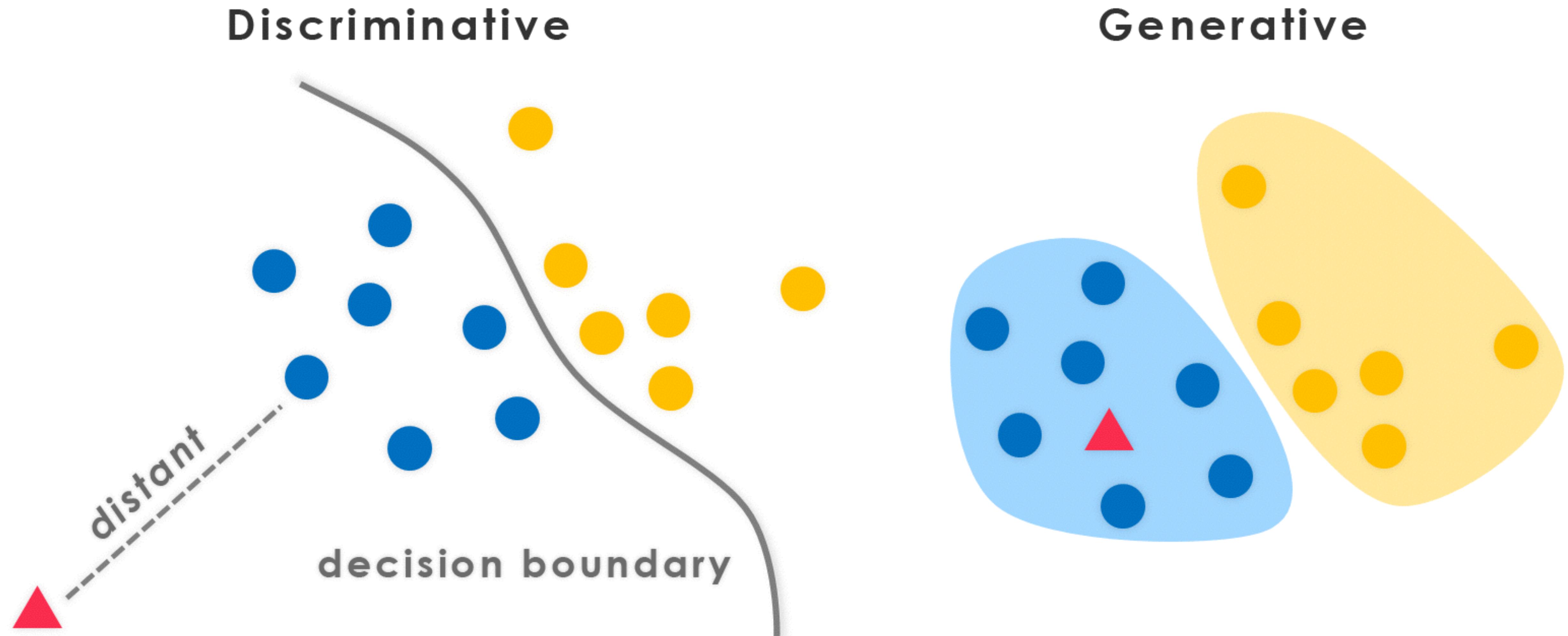
- UW Museum Student Art Competition
- Judges did not know art was AI-generated (and not human artist)
- 35% acceptance rate, and an award

# Why are networks easily fooled?



DNN Confidence: > 99.6 % for all

# Hypothesis 1: DNNs do understand, test is bad



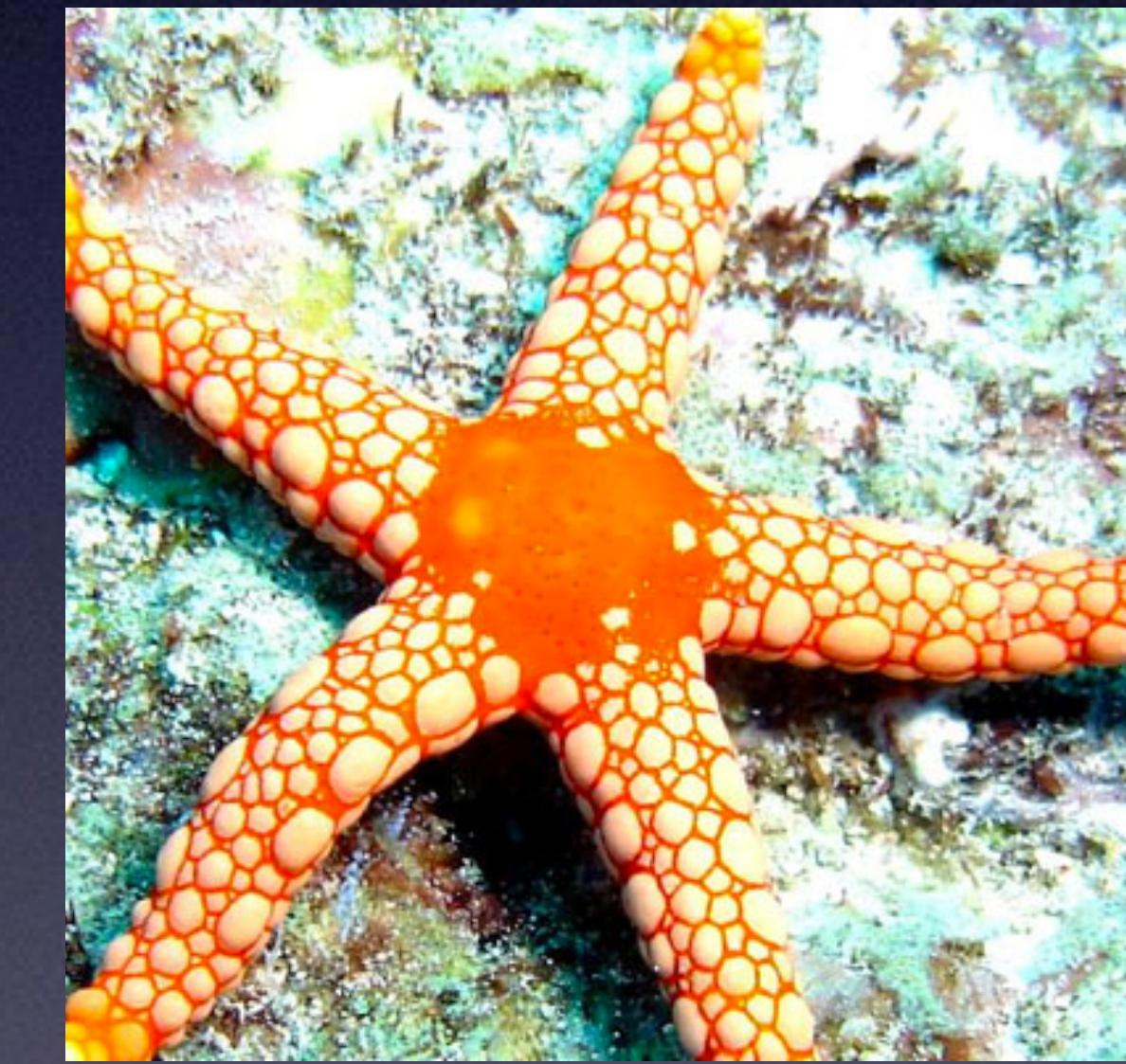
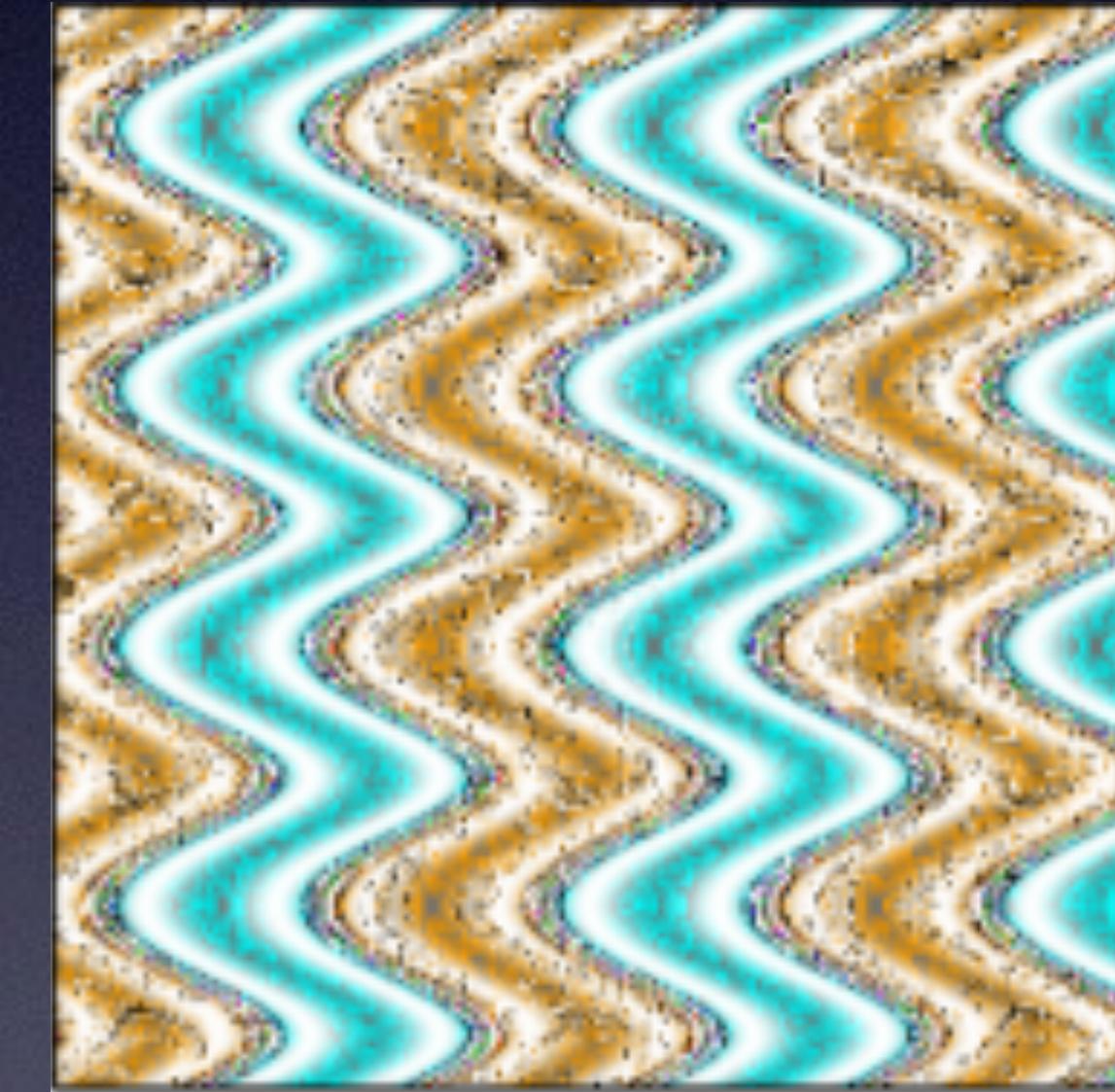
Prediction: With constraints to stay in the space of natural images, we **WOULD** get recognizable objects.

# Hypothesis 2: Clever Hans

- Only learns distinguishing features



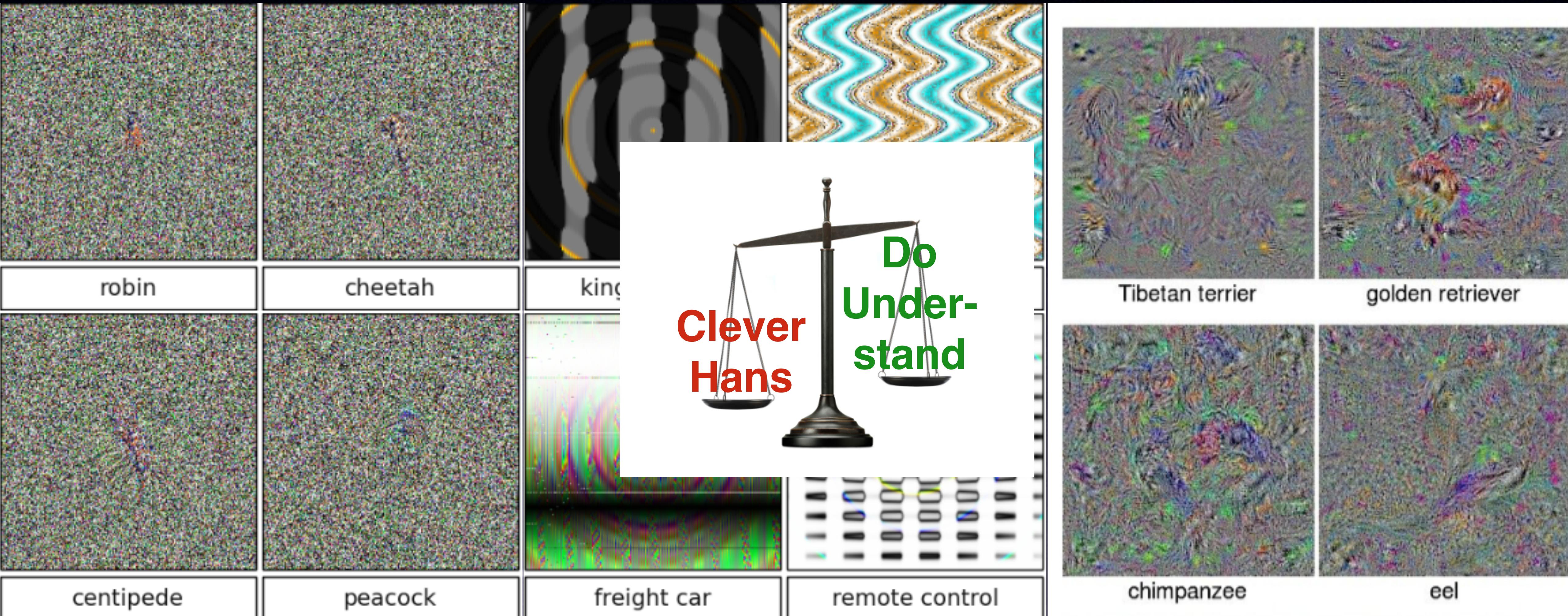
School Bus



Starfish

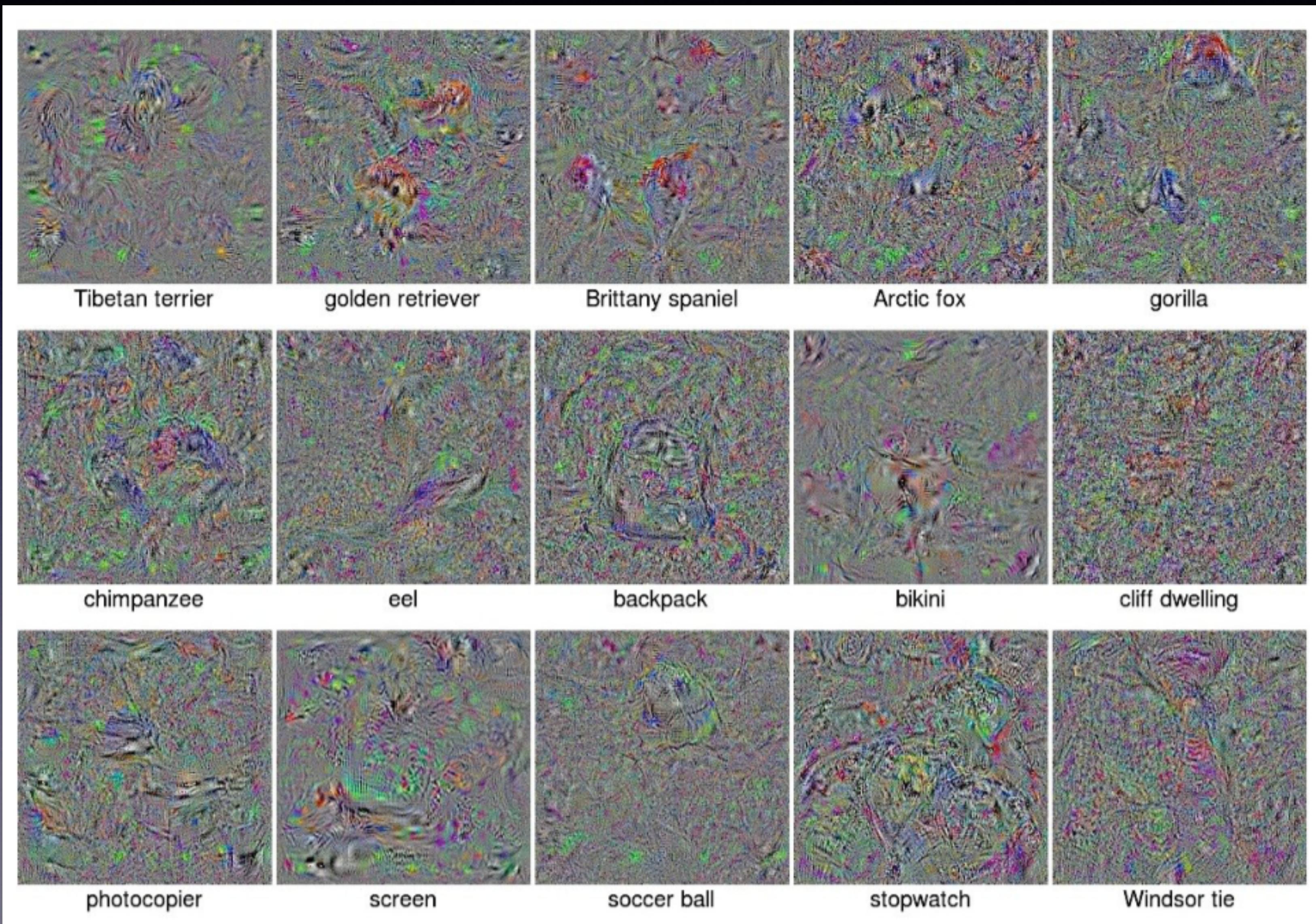
Prediction: With constraints to stay in the space of natural images, we **WOULD NOT** get recognizable objects.

# Clever Hans seems more likely...



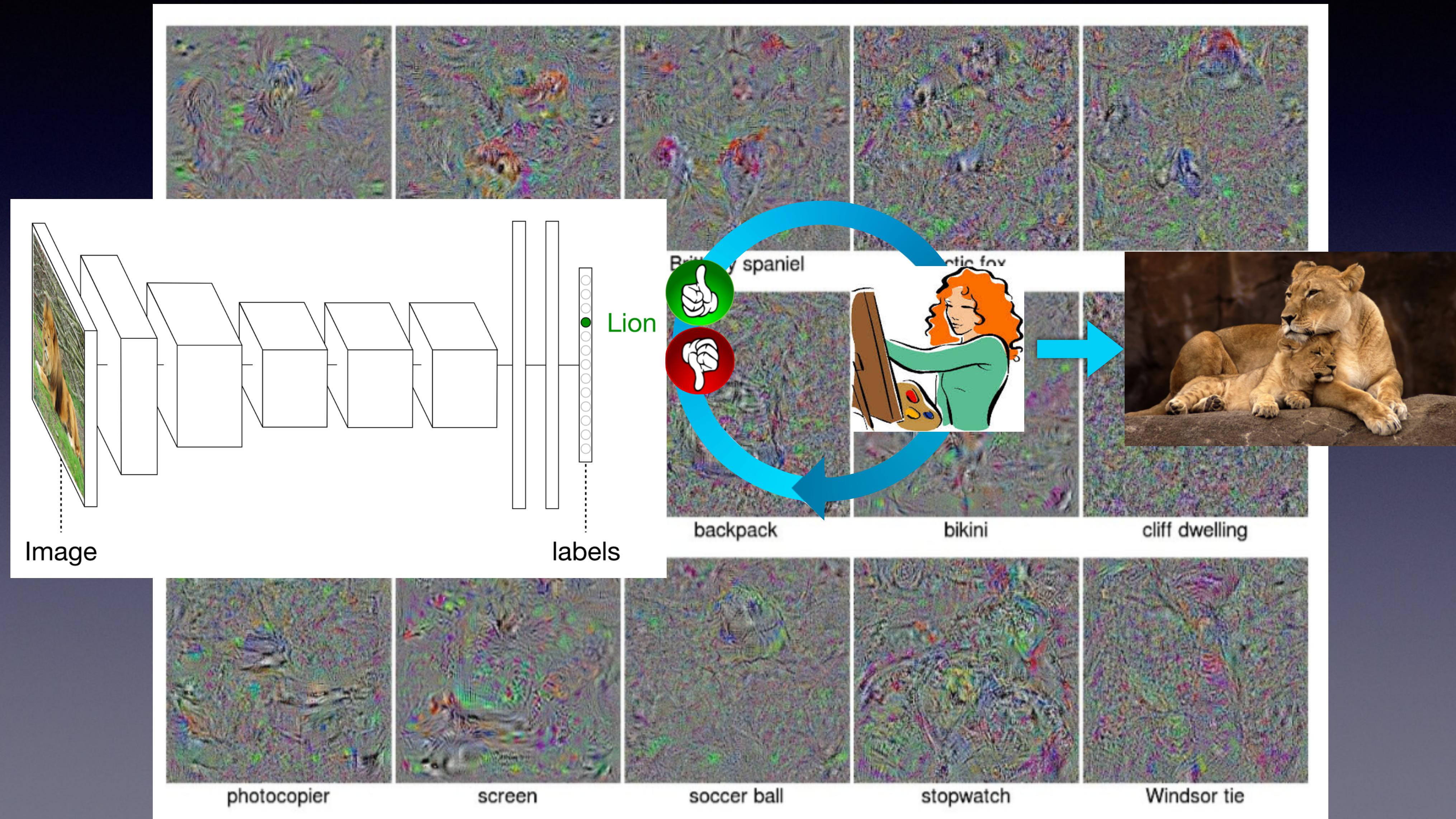
# Deep Visualization Take 1

Nguyen, Yosinski, Clune, 2015, CVPR



# Deep Visualization Take 1

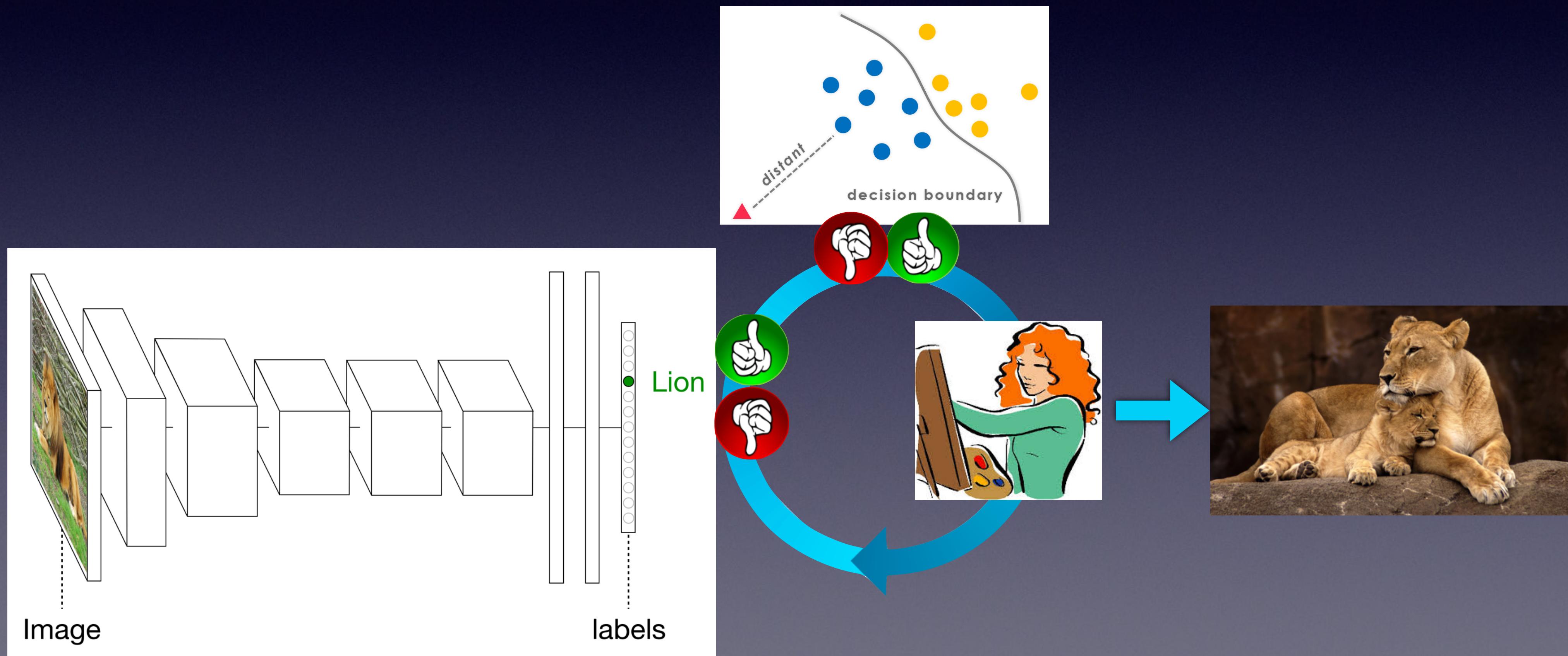
Nguyen, Yosinski, Clune, 2015, CVPR



# Deep Visualization Take 2

Yosinski, Clune, Nguyen, Lipson, 2015, ICML Deep Learning Workshop

## Manually Engineered Natural Image Priors

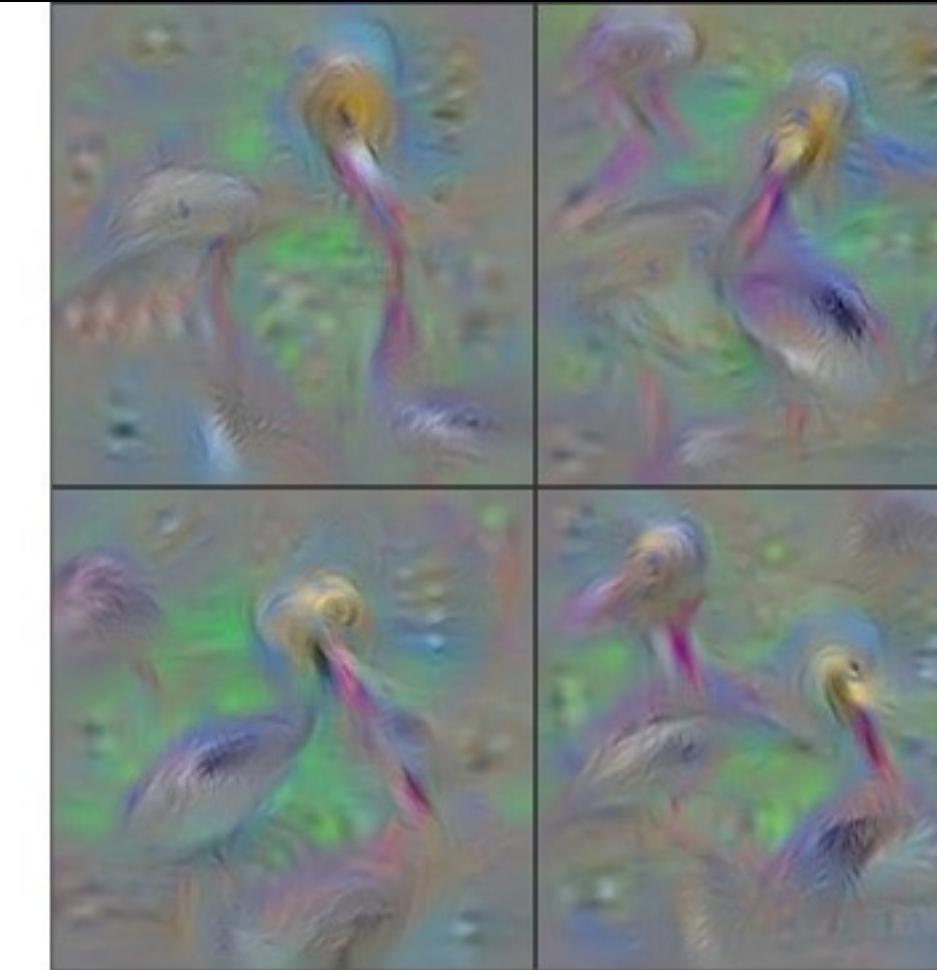


# Deep Visualization Take 2

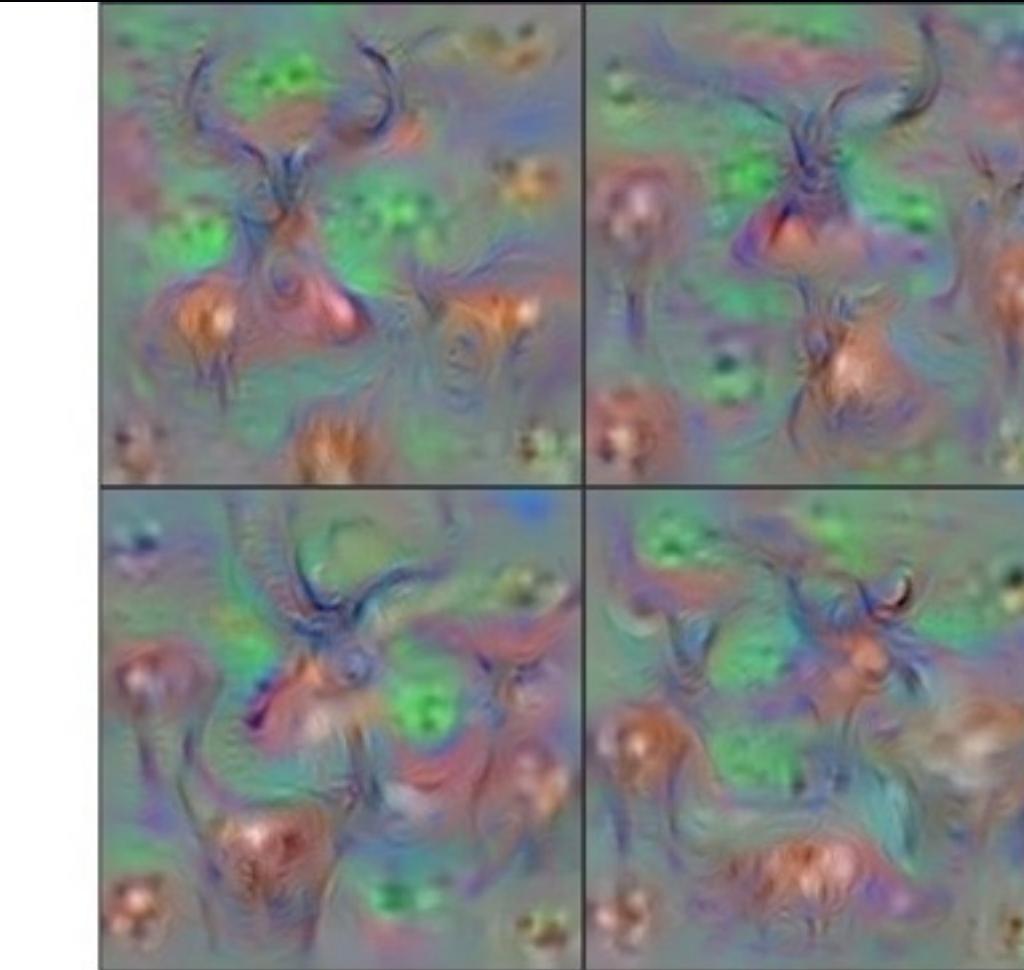
Yosinski, Clune, Nguyen, Lipson, 2015, ICML Deep Learning Workshop



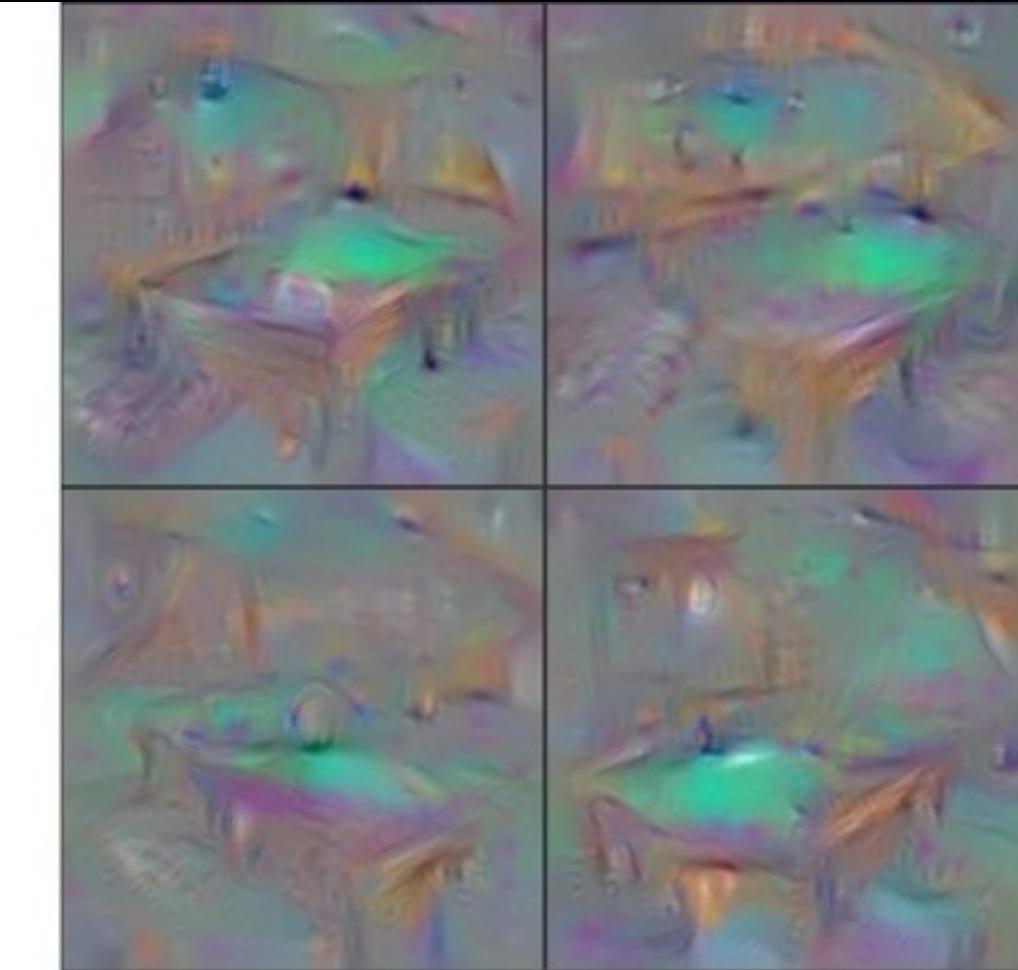
Flamingo



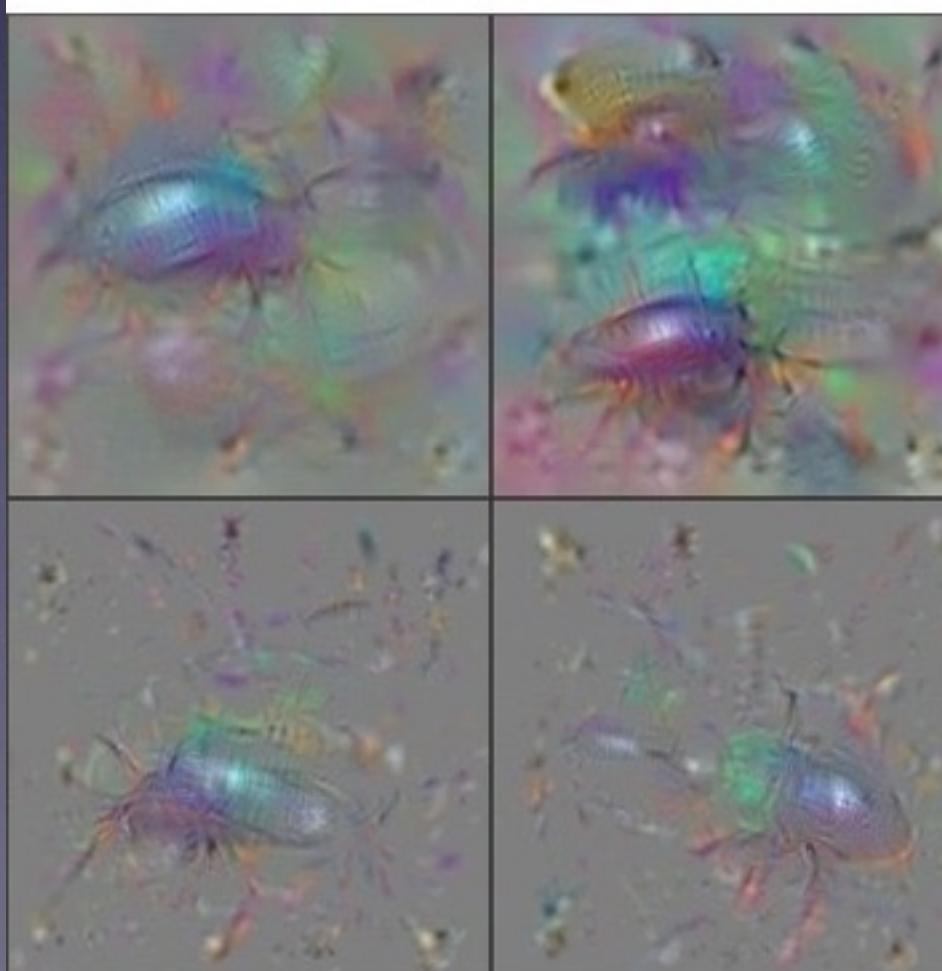
Pelican



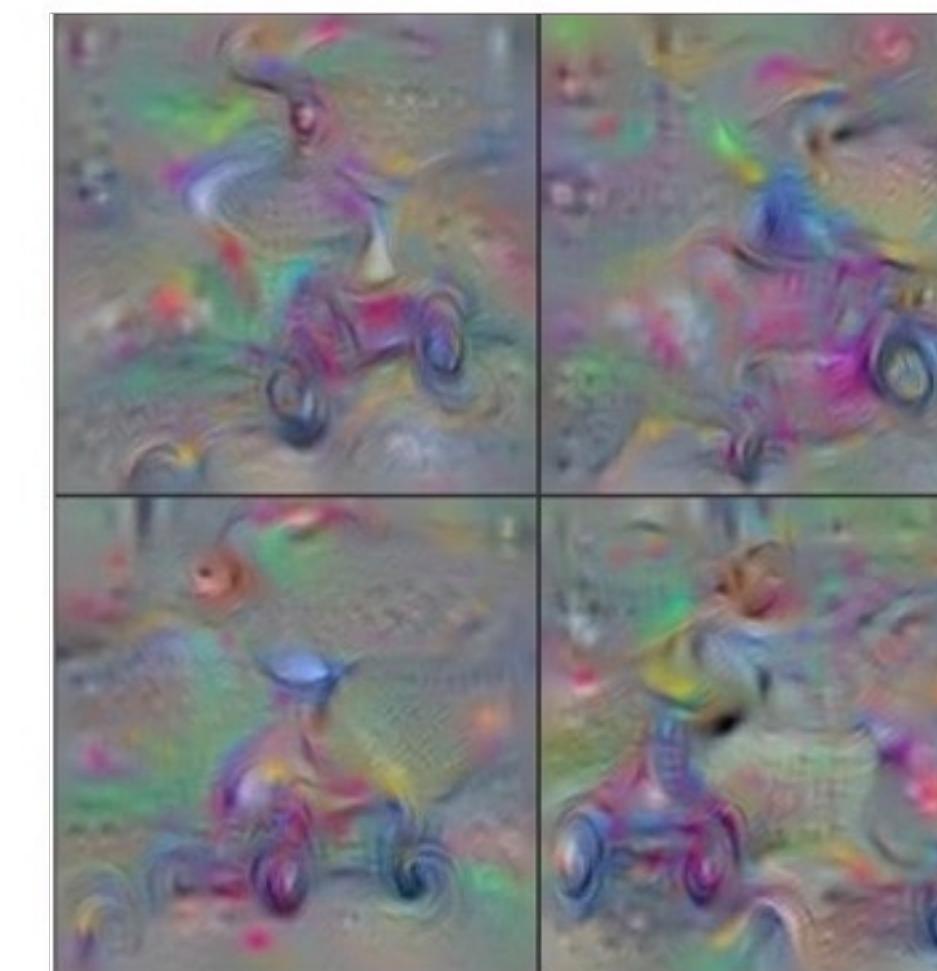
Hartebeest



Billiard Table



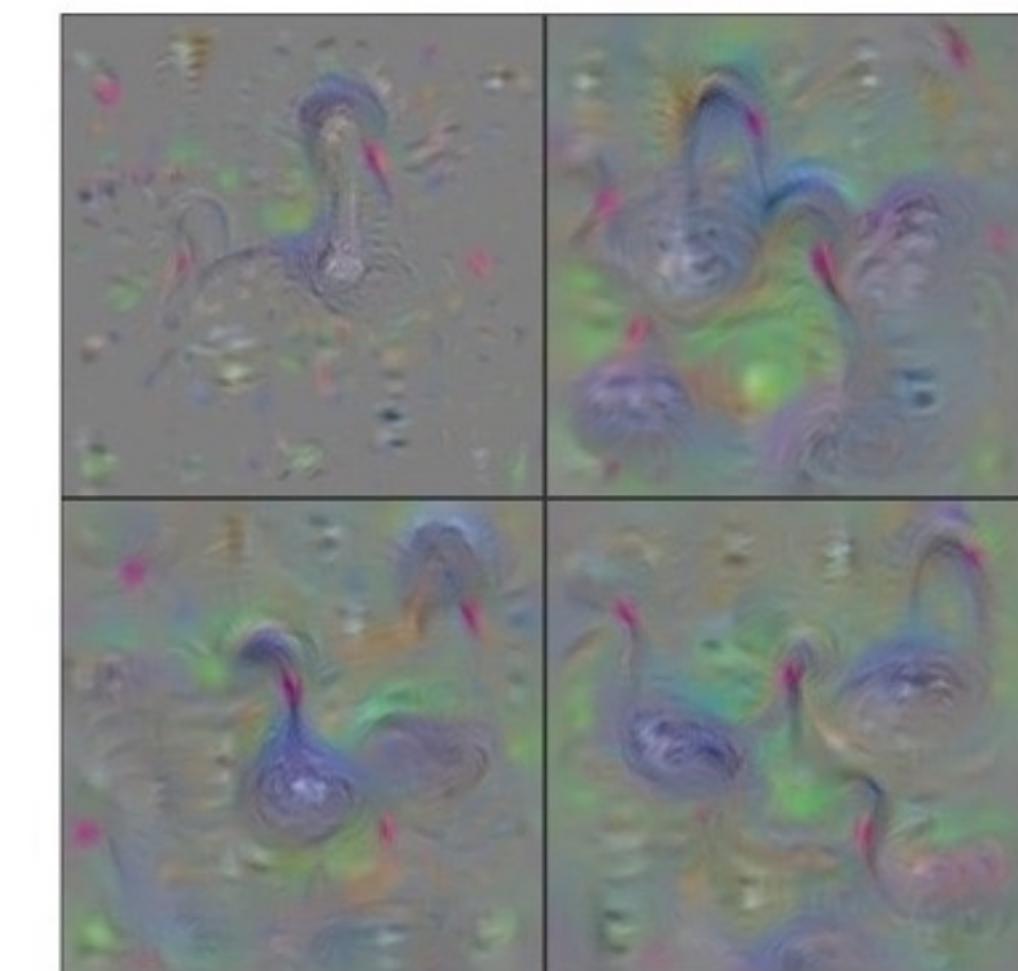
Ground Beetle



Tricycle



School Bus



Black Swan

# Deep Visualization Take 2

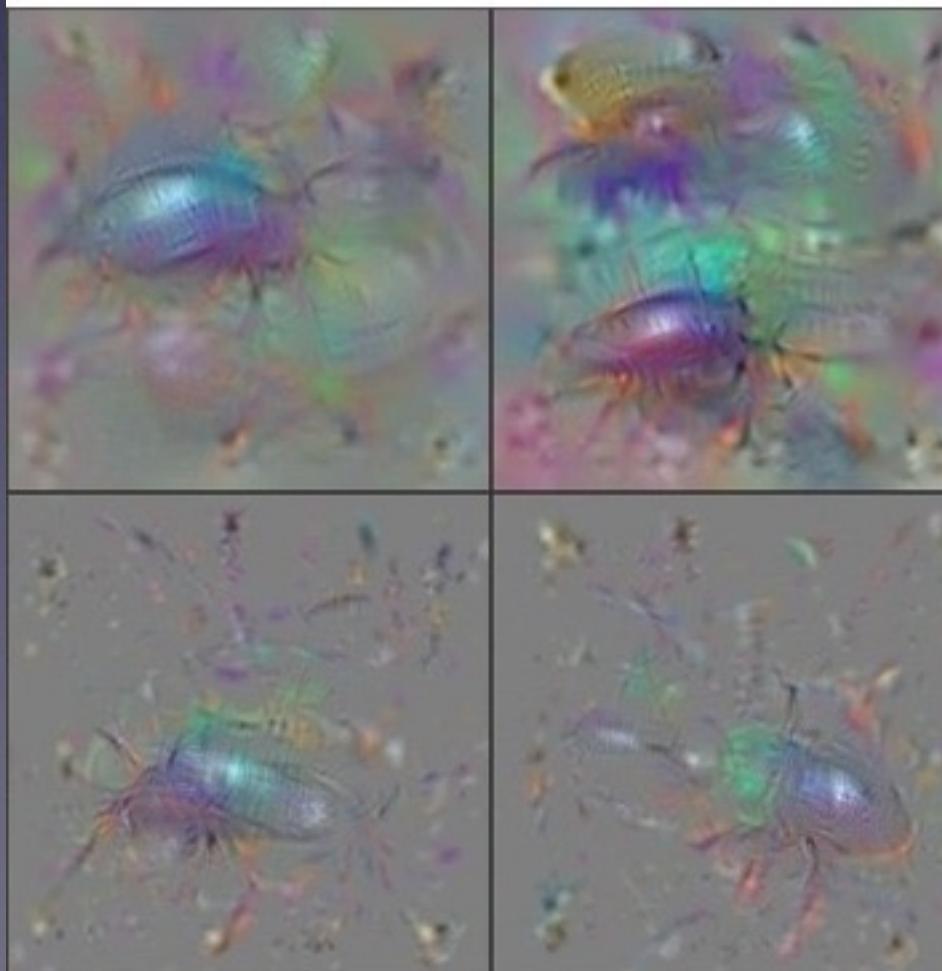
Yosinski, Clune, Nguyen, Lipson, 2015, ICML Deep Learning Workshop



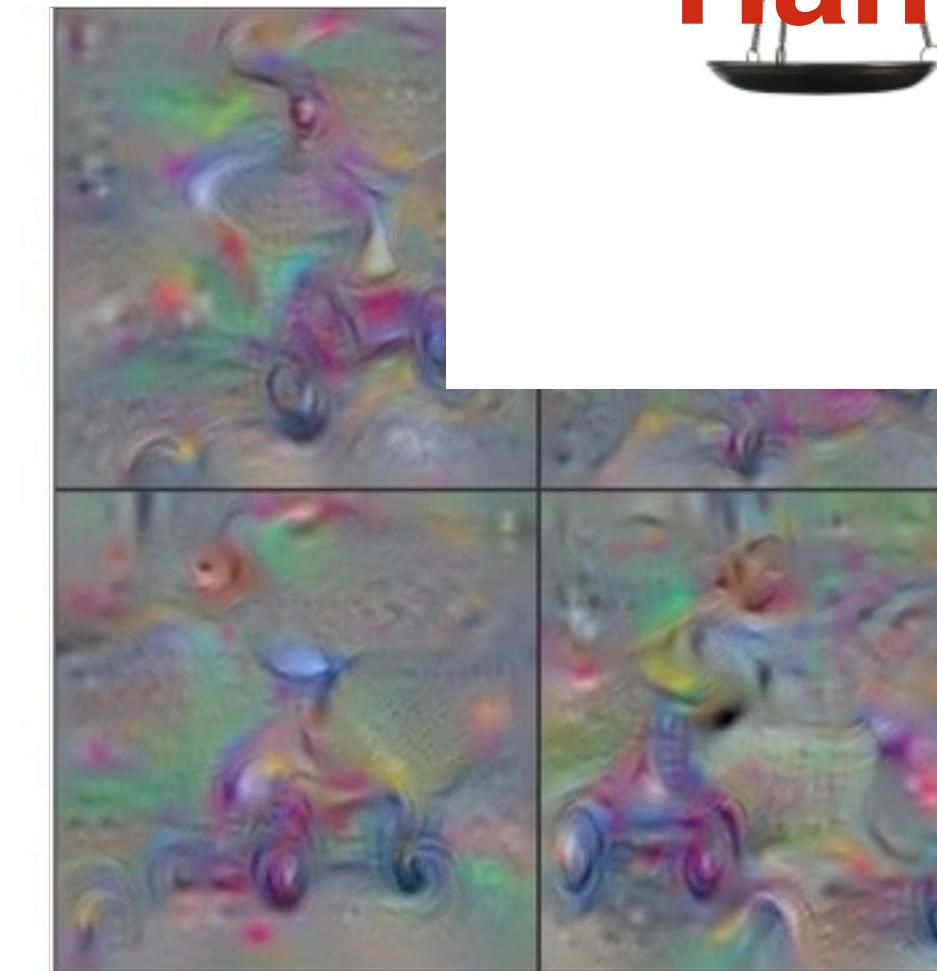
Flamingo



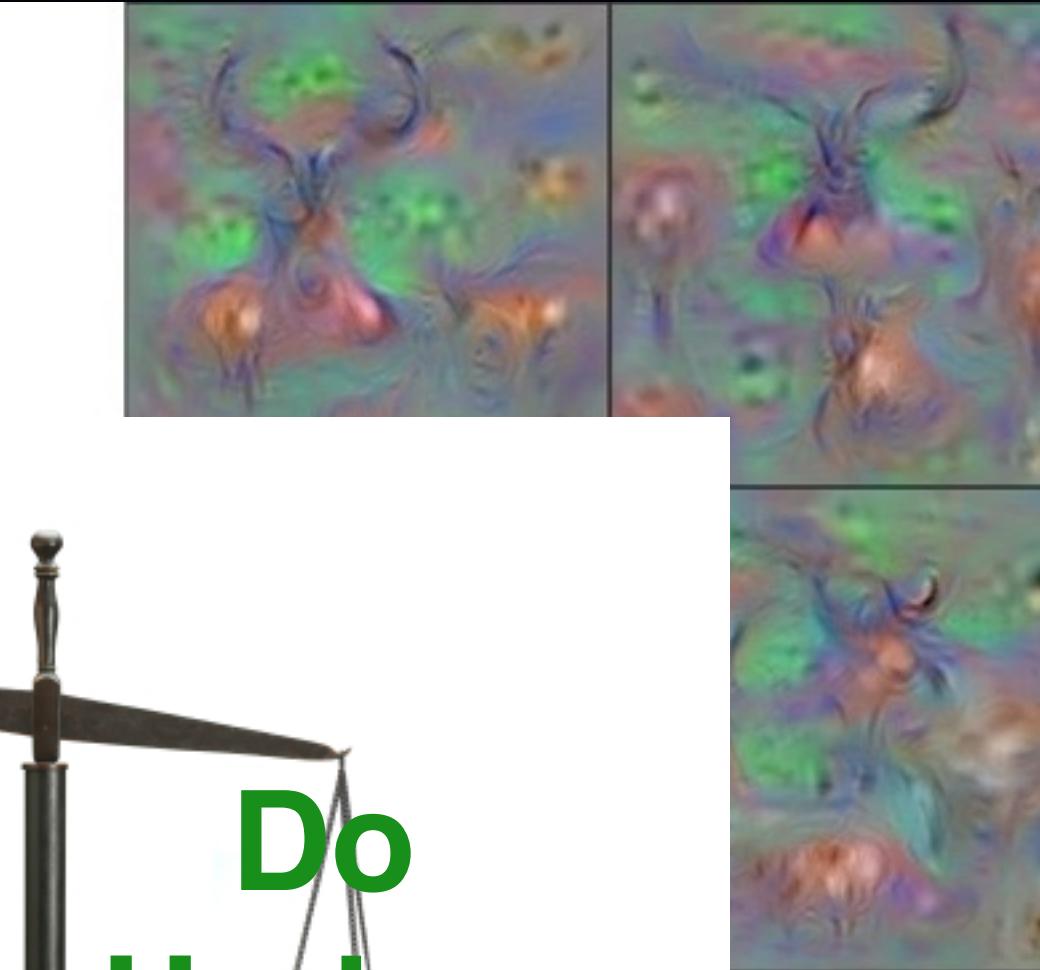
P



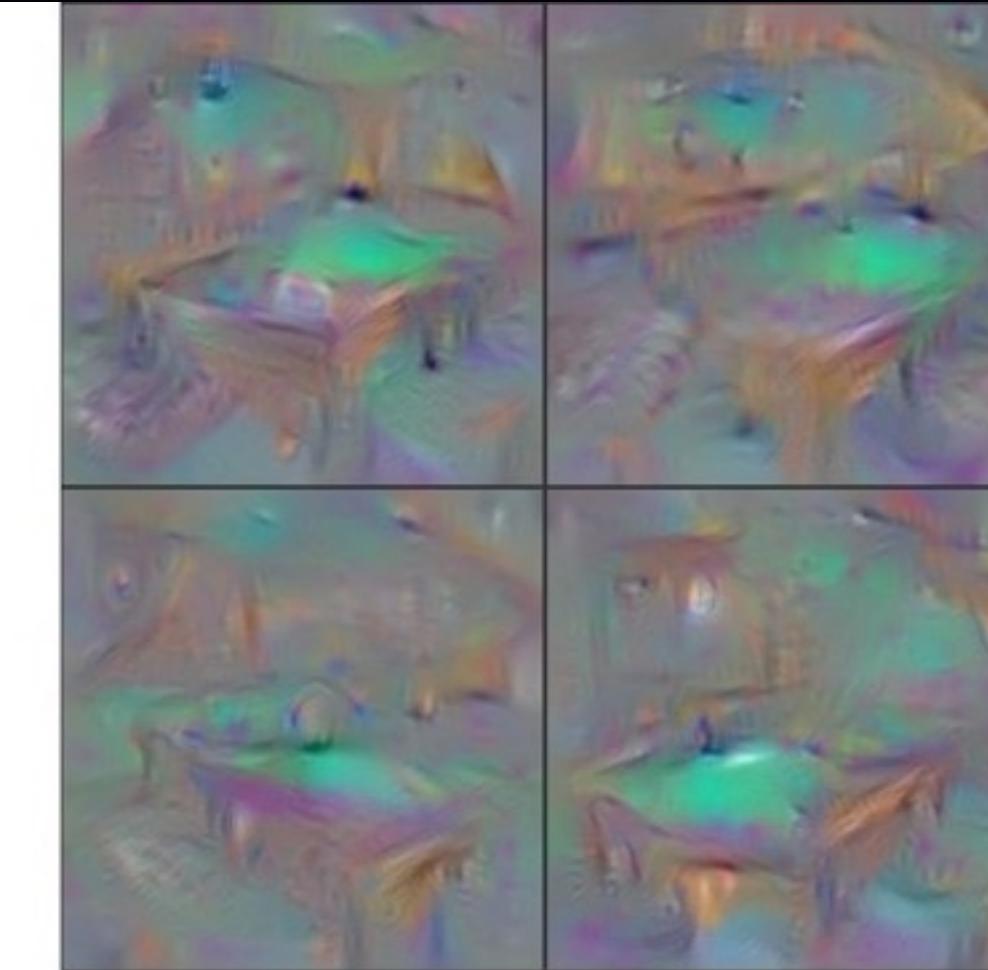
Ground Beetle



Tricycle



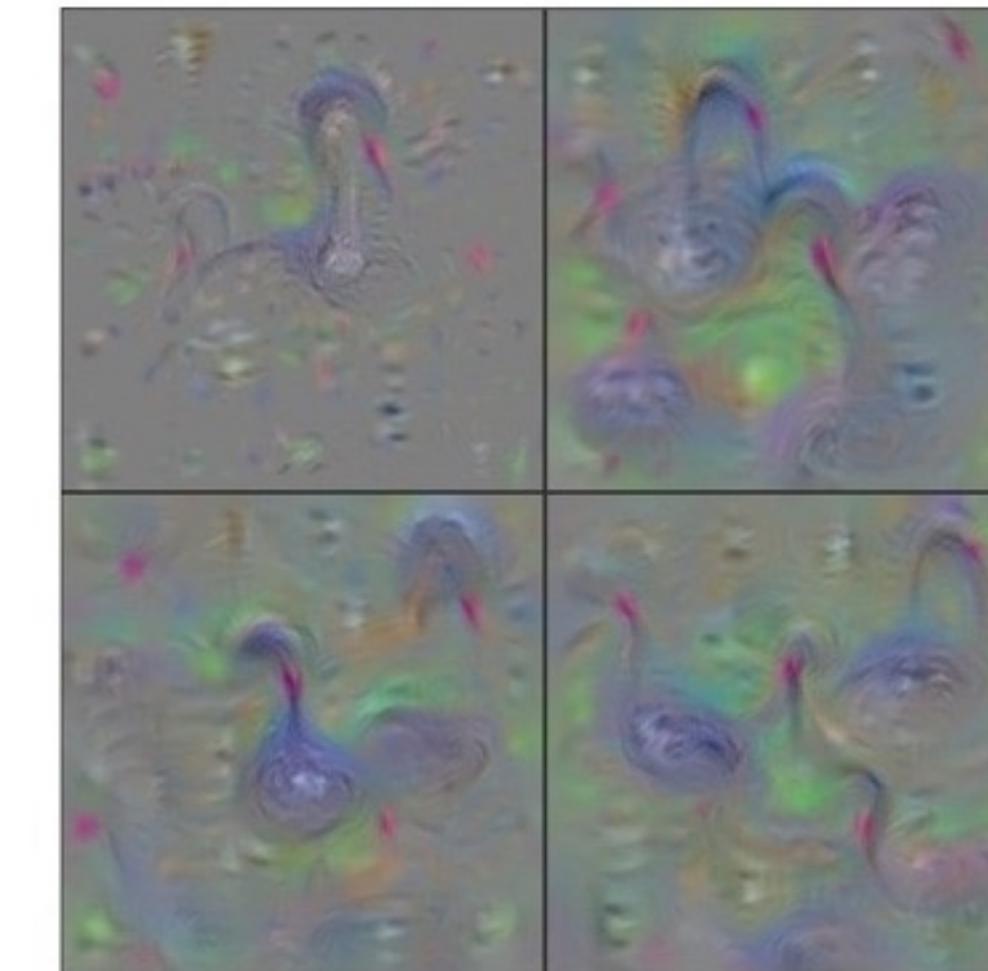
est



Billiard Table



School Bus

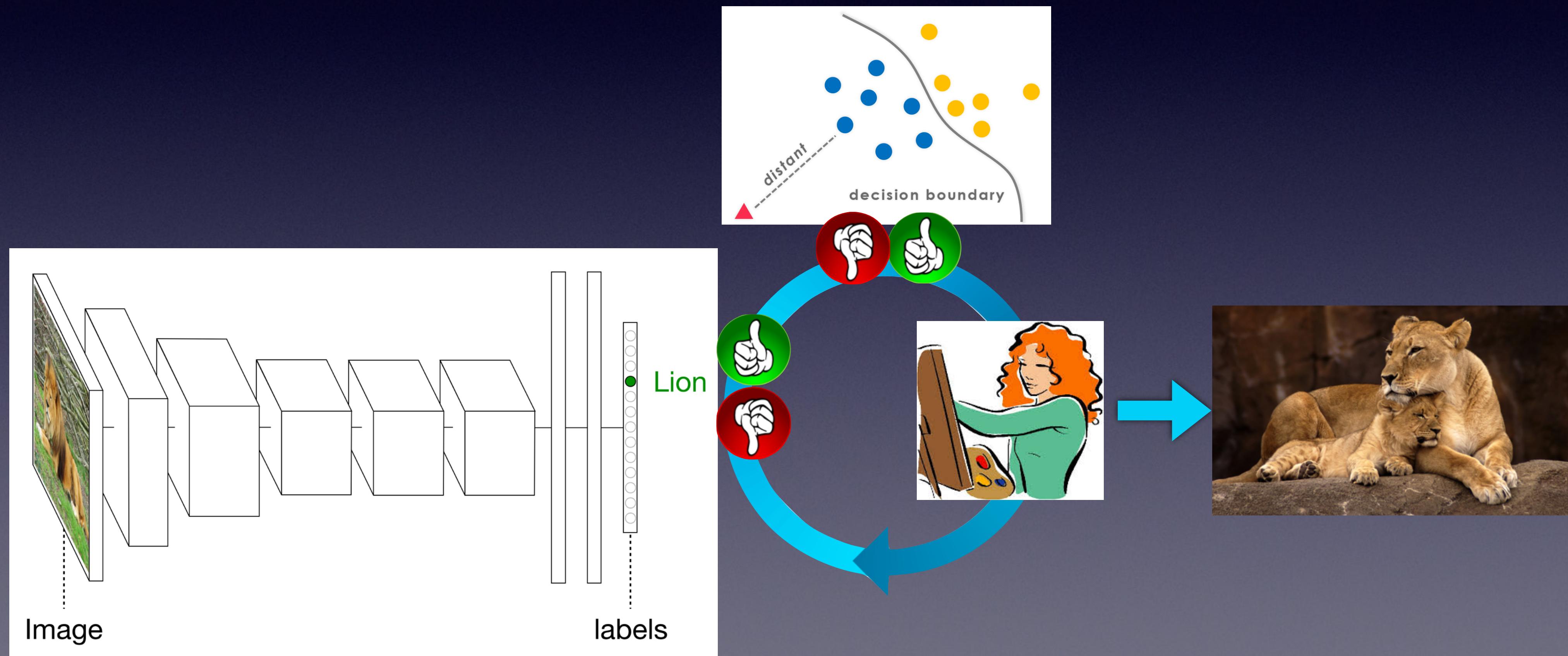


Black Swan

# Deep Visualization Take 3

Multifaceted Feature Visualization. Nguyen, Yosinski, Clune 2016, ICML Workshop

Different Manually Engineered Natural Image Priors  
(+ mean seeds from data)



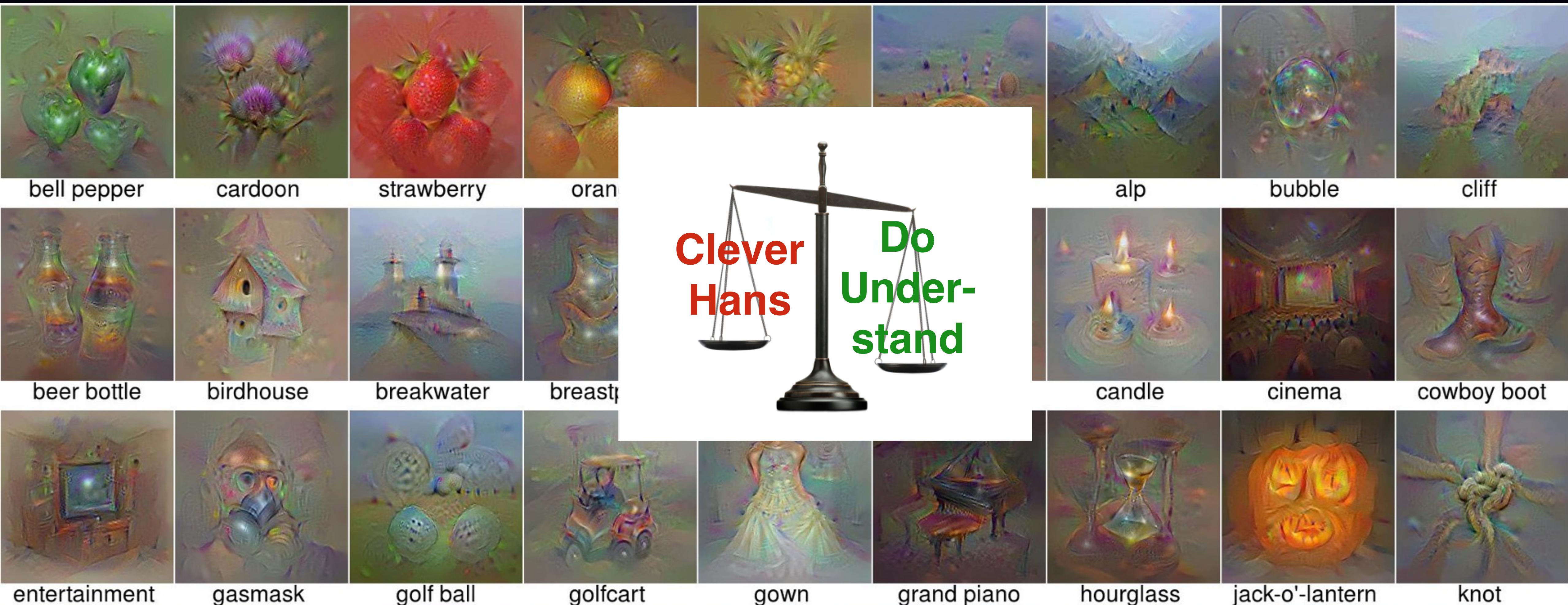
# Deep Visualization Take 3

Multifaceted Feature Visualization. Nguyen, Yosinski, Clune 2016, ICML Workshop



# Deep Visualization Take 3

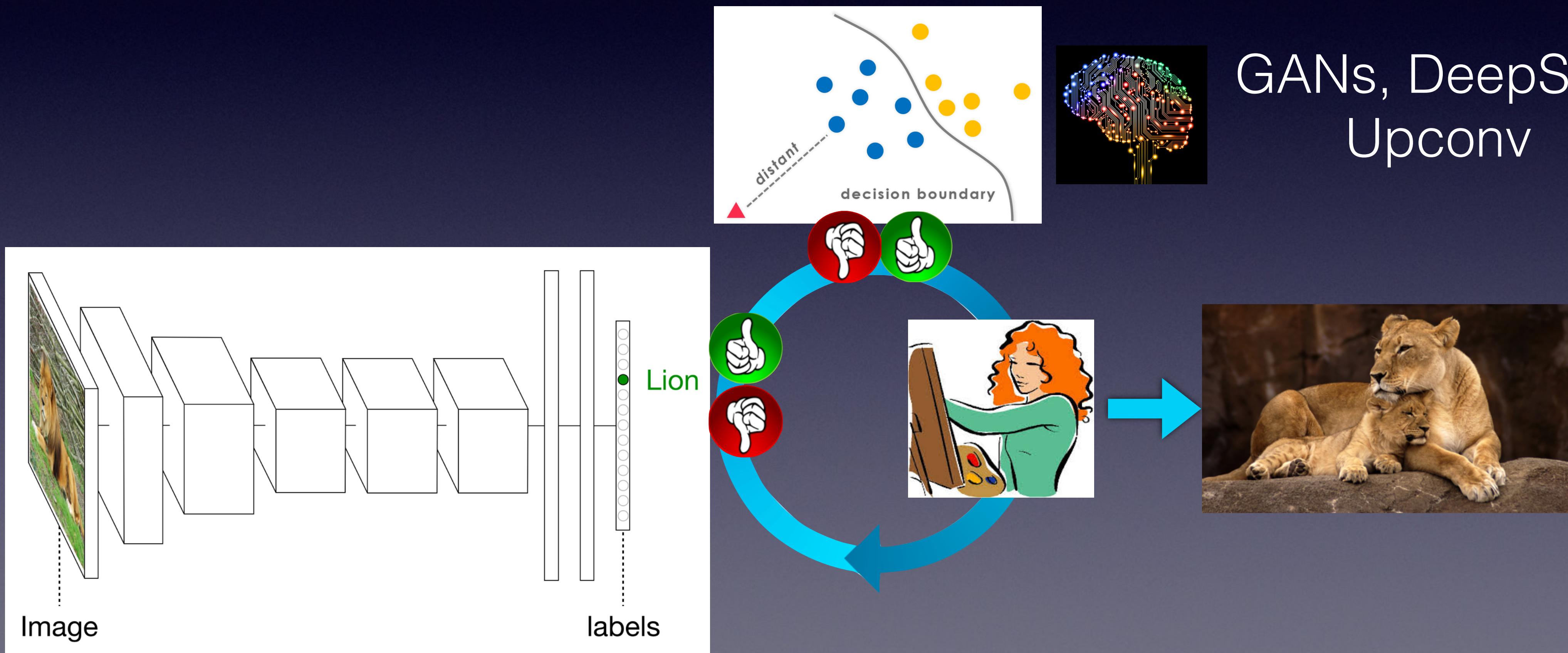
Multifaceted Feature Visualization. Nguyen, Yosinski, Clune 2016, ICML Workshop



# Deep Visualization Take 4

Nguyen, Dosovitskiy, Yosinski, Brox, Clune. 2016. Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. NIPS

## Learned Natural Image Priors



# Deep Visualization Take 4

Nguyen, Dosovitskiy, Yosinski, Brox, Clune. 2016. NIPS



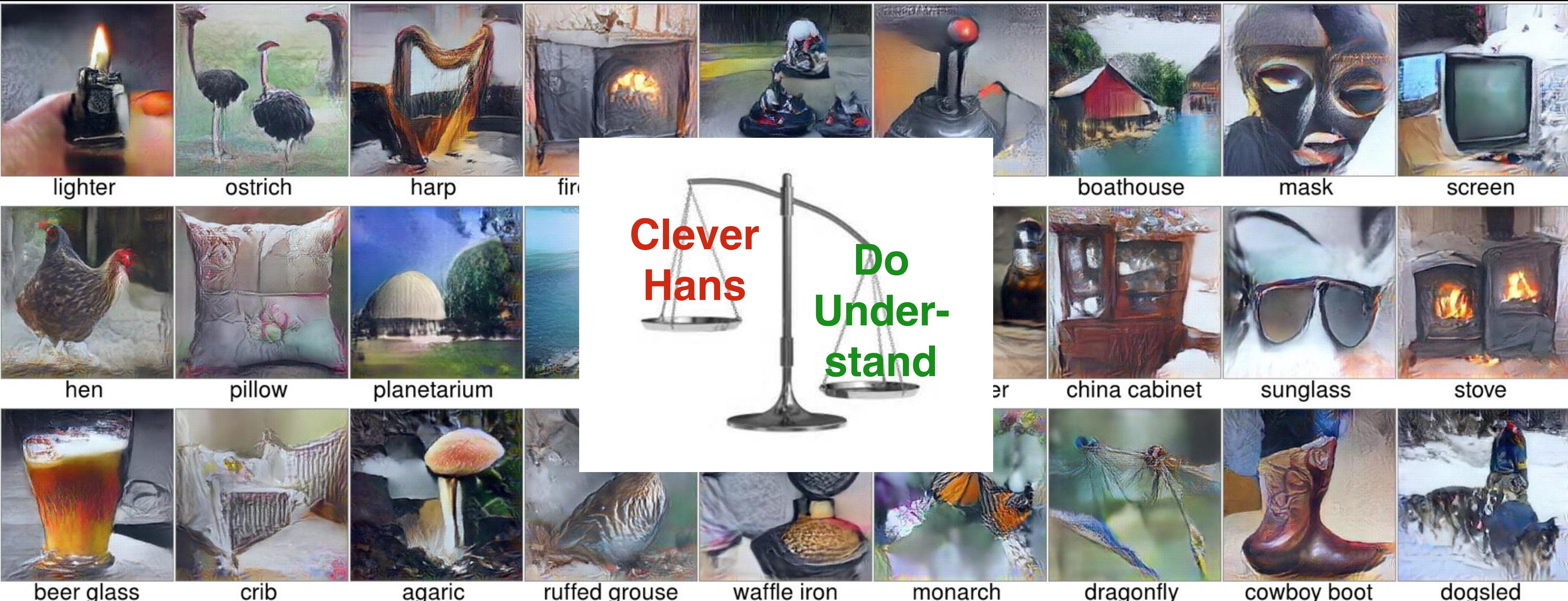
# Deep Visualization Take 4

Nguyen, Dosovitskiy, Yosinski, Clune. Work in progress.

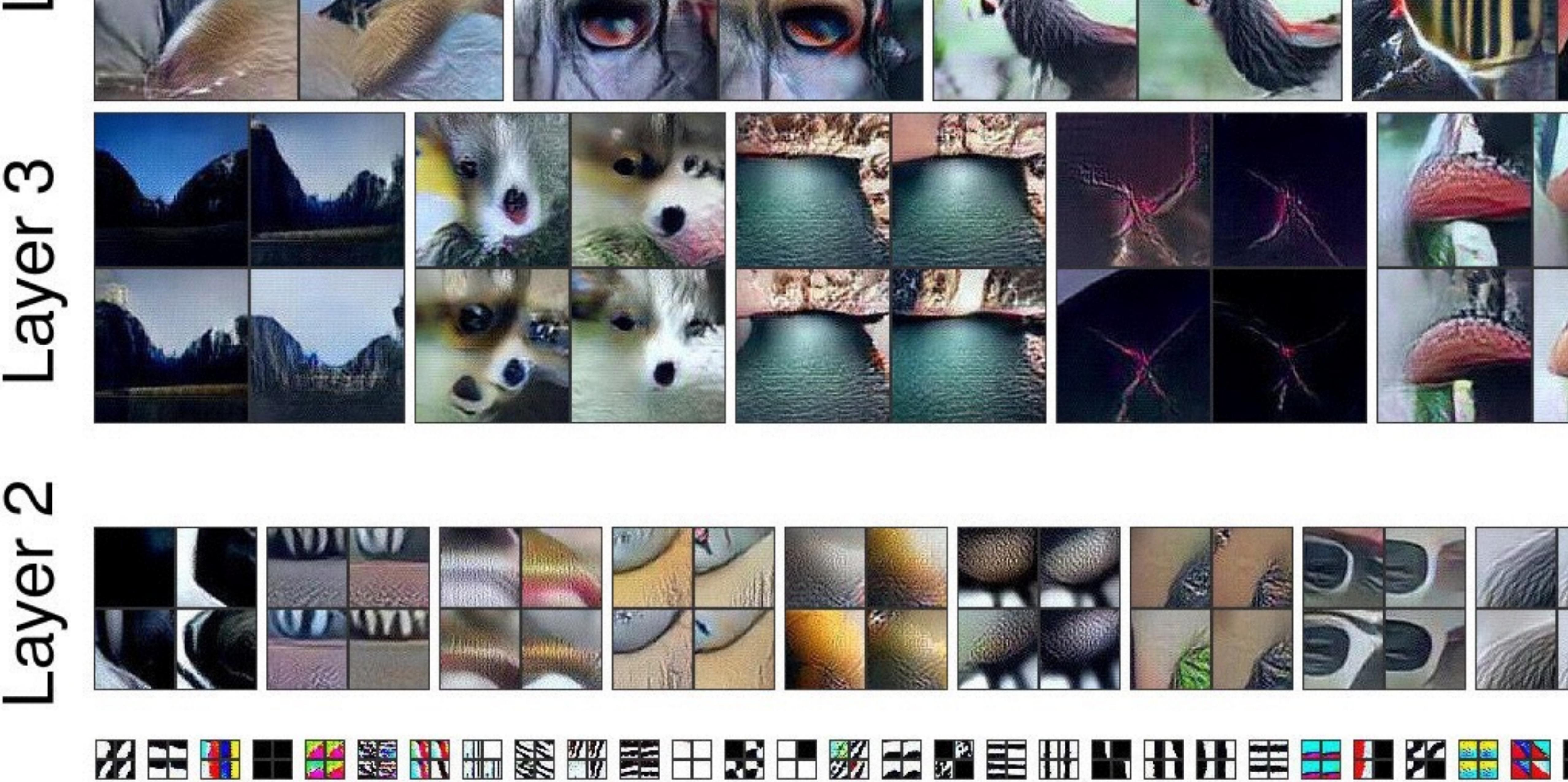


# Deep Visualization Take 4

Nguyen, Dosovitskiy, Yosinski, Clune. Work in progress.







Layer 1

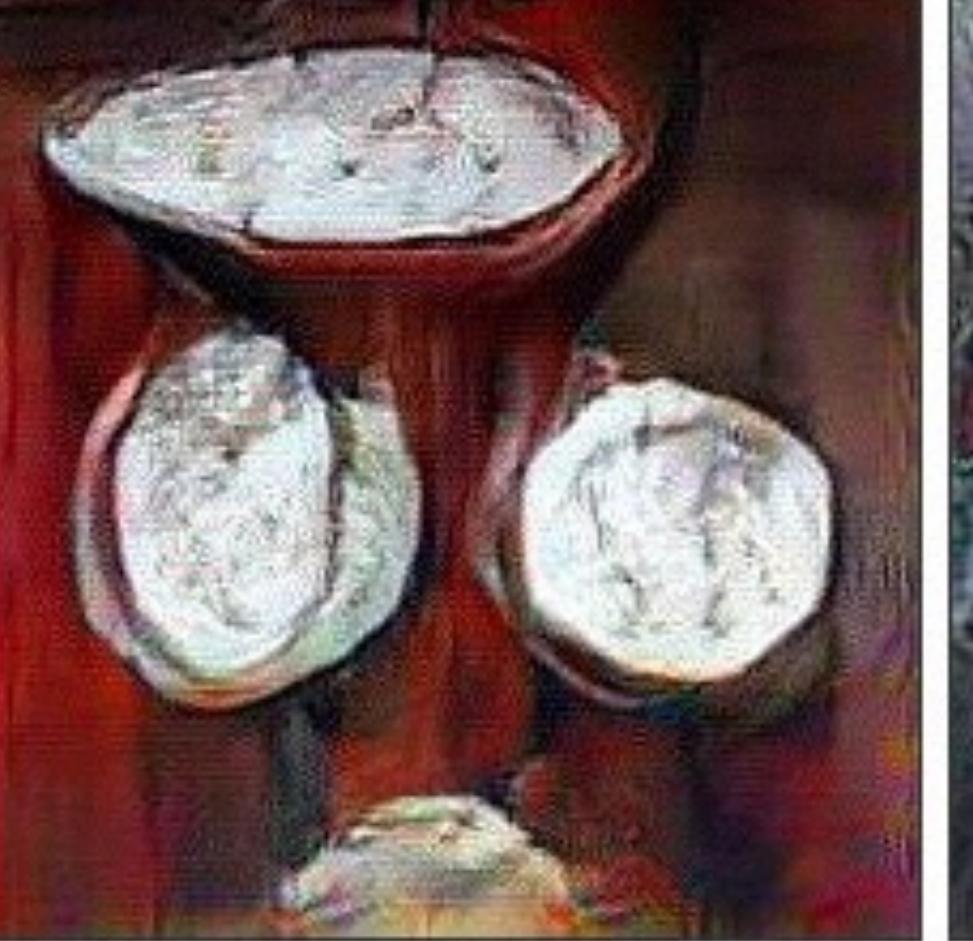
# Layer 5



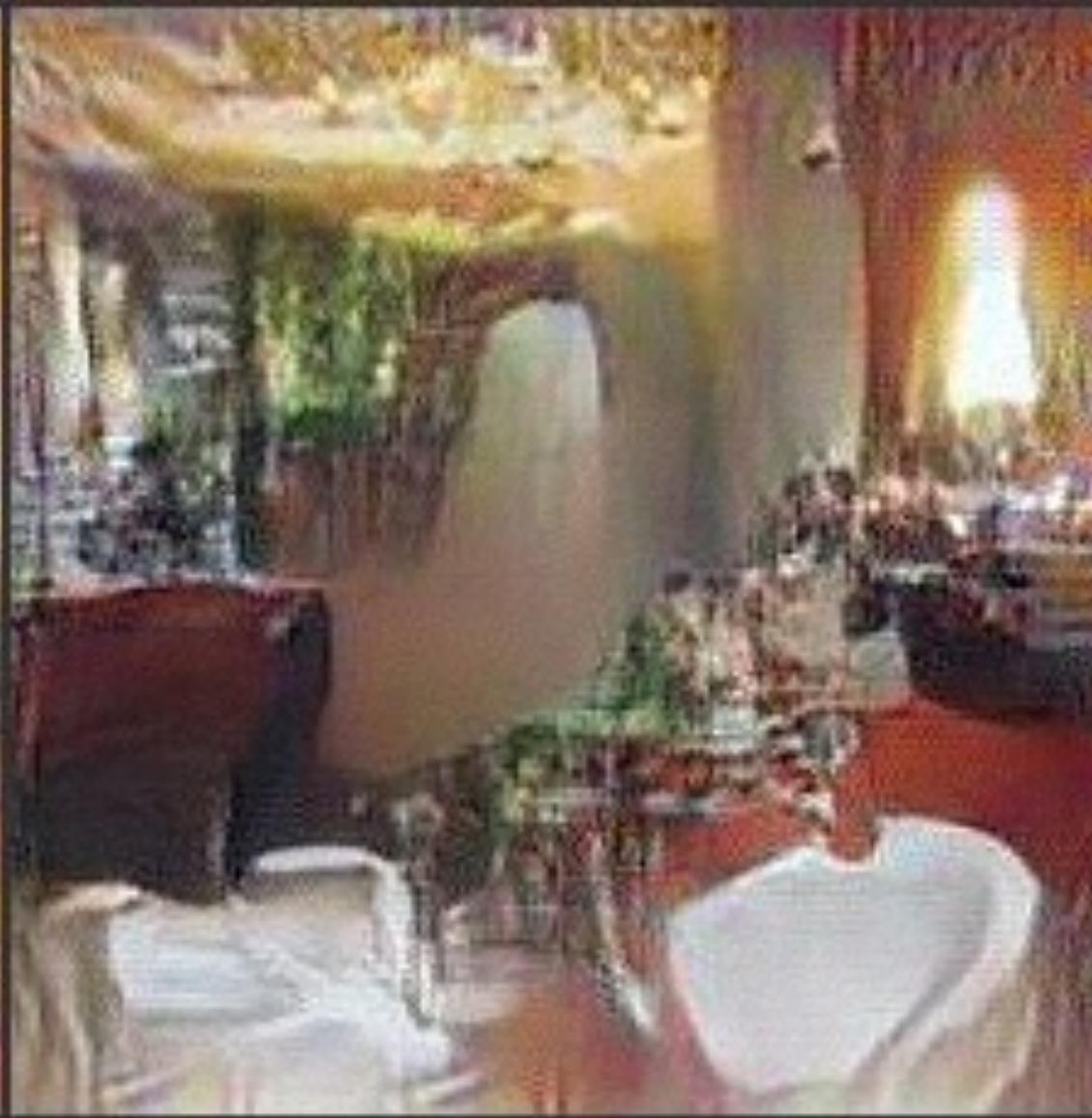
# Layer 4



# 3



Layer 8

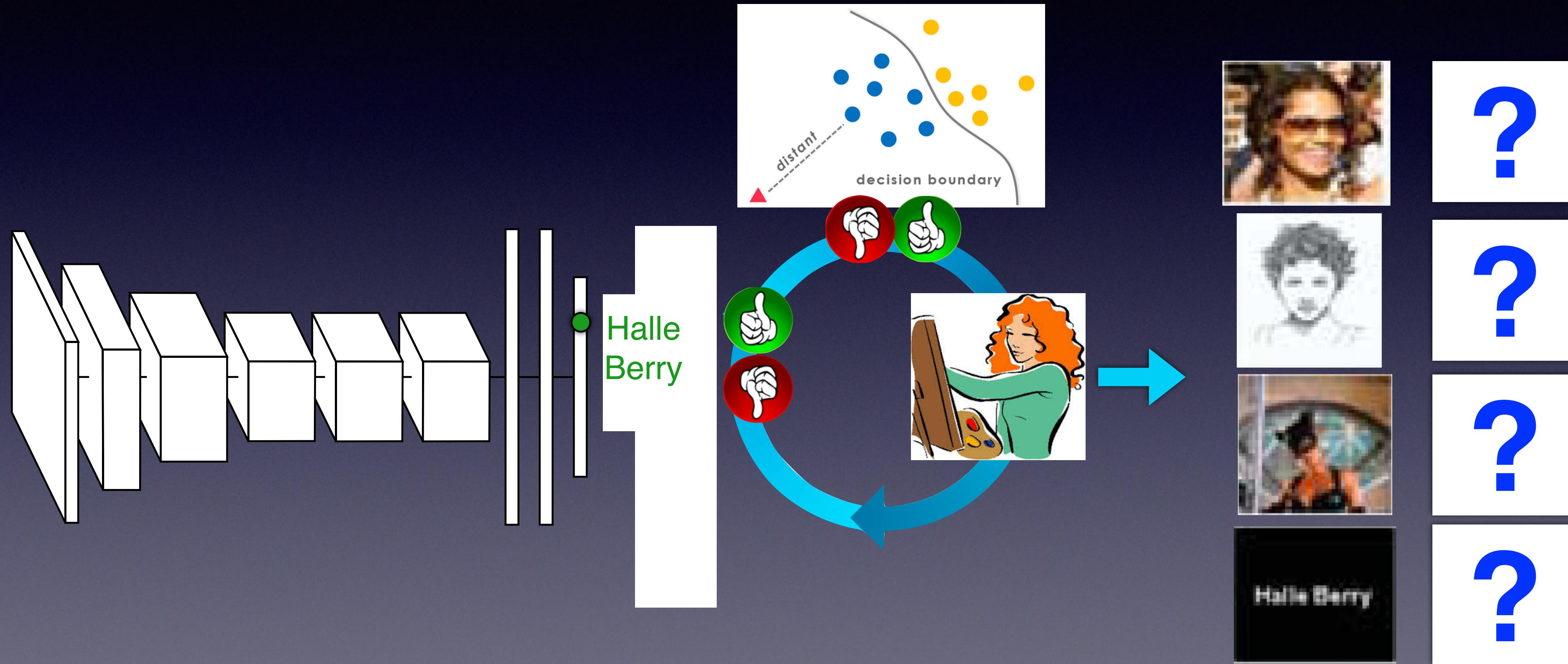


Restaurant

Ostrich

# Multi-Faceted Feature Visualization

Multifaceted Feature Visualization. Nguyen, Yosinski, Clune 2016, ICML Workshop







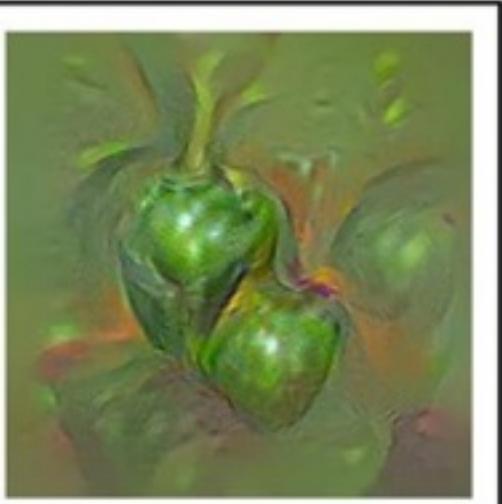
1



2



3



4



5



6



7



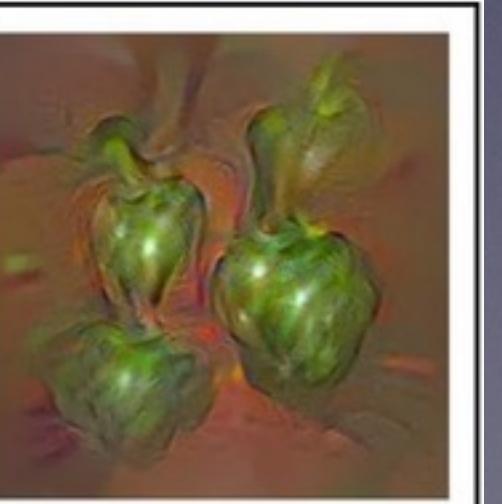
8



9



10







# Deep Visualization Toolbox

[yosinski.com/deepvis](http://yosinski.com/deepvis)

#deepvis



Jason Yosinski



Jeff Clune



Anh Nguyen



Thomas Fuchs



Hod Lipson



Cornell University



UNIVERSITY  
OF WYOMING



**Jet Propulsion Laboratory**  
California Institute of Technology

# Are DNN representations local or distributed?

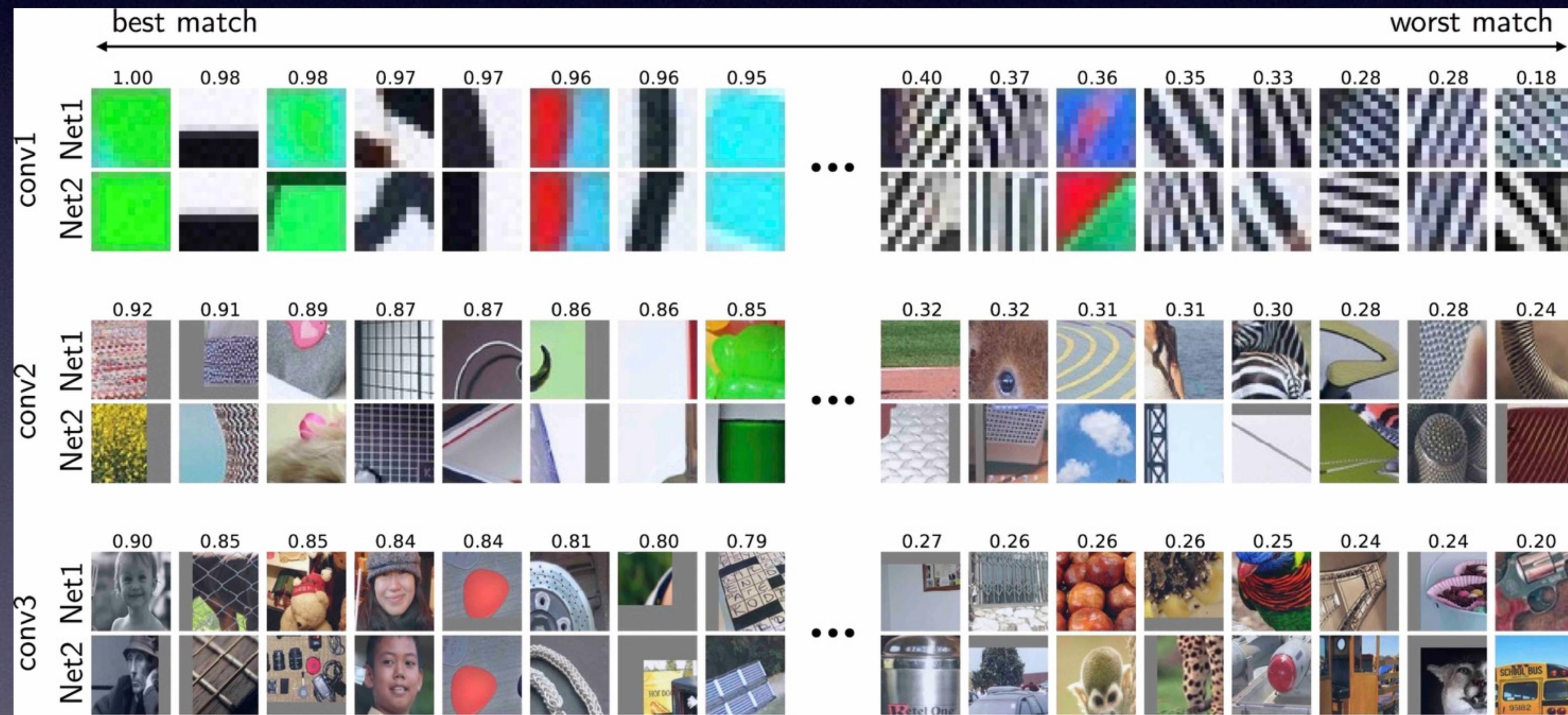
Li Y, Yosinski J, Clune J, Lipson H, Hopcroft J (2016) Convergent Learning: Do different neural networks learn the same representations? International Conference on Learning Representations (ICLR '16).

- Answer: “It’s complicated.”

# Do different DNNs learn the same features?

Li Y, Yosinski J, Clune J, Lipson H, Hopcroft J (2016) Convergent Learning: Do different neural networks learn the same representations? International Conference on Learning Representations (ICLR '16).

Answer: Yes for most active features. No for least active ones.



# Conclusions

- Initial “fooling” work suggested Deep Learning is a Clever Hans
- Work since reveals DNNs learn a surprising amount about the concepts they classify
  - Understand an entire object & its context
  - Understand the multifaceted nature of concepts



# Conclusions

- Why fooling?
  - Asking for maximum confidence
  - Space of images is vast
  - We are not training for fool-proof images



# Future Work Potential

- Demonstrate all of this in other modes
  - Speech recognition
  - Music classification
  - Including learning natural priors to prevent fooling
- Investigate with unsupervised learning
  - and reinforcement learning
- Can we generate such “preferred stimuli” for biological brains?



# Thanks



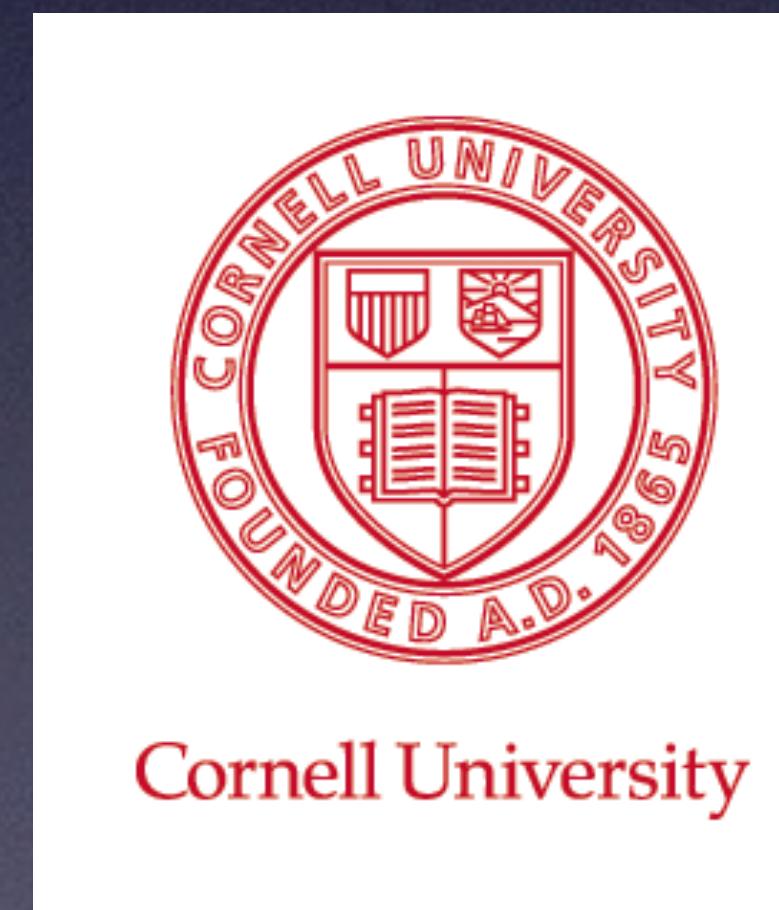
Anh Nguyen



Jason Yosinski



Alexey Dosovitskiy



Cornell University

