# Visualizing and Understanding Artificial Neural Networks

**Anh Nguyen**

Assistant Professor

AUBURN UNIVERSITY

# Neuroscience

to understand the brains of _____


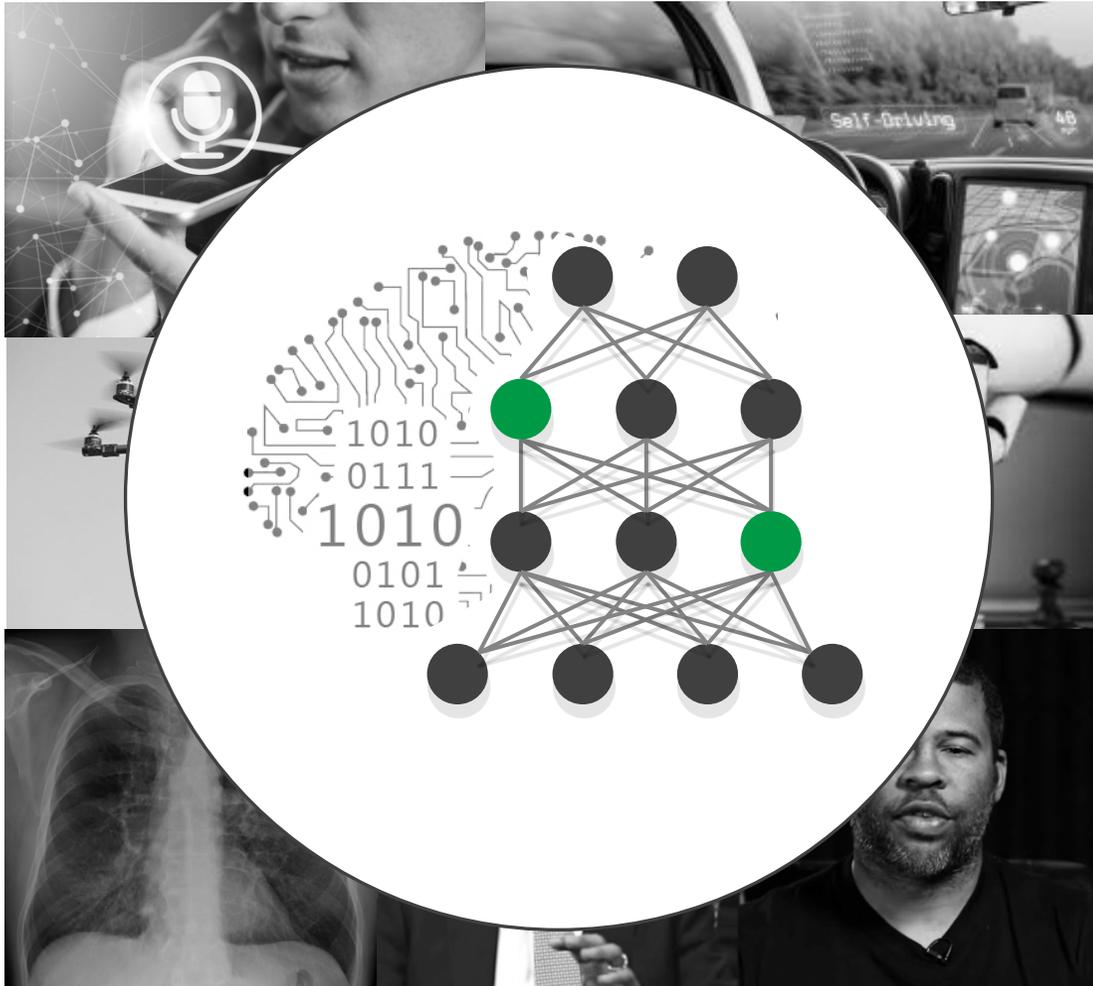
humans ☑

cats ☑

mice ☑

macaques ☑

...
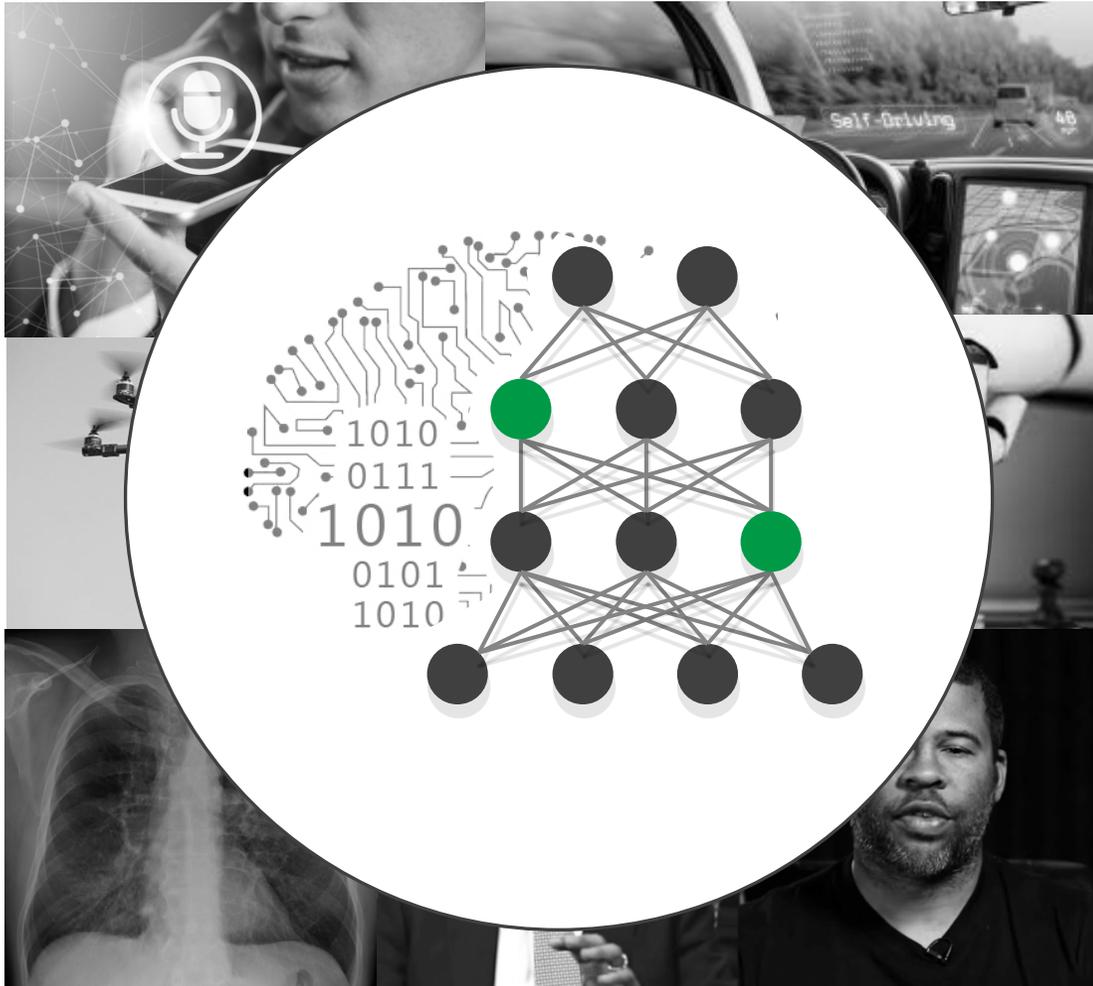
**AIs** / machines ☑

# Neuroscience

Using **AIs** to understand the brains of _____

humans ☑

cats ☑

mice ☑

macaques ☑
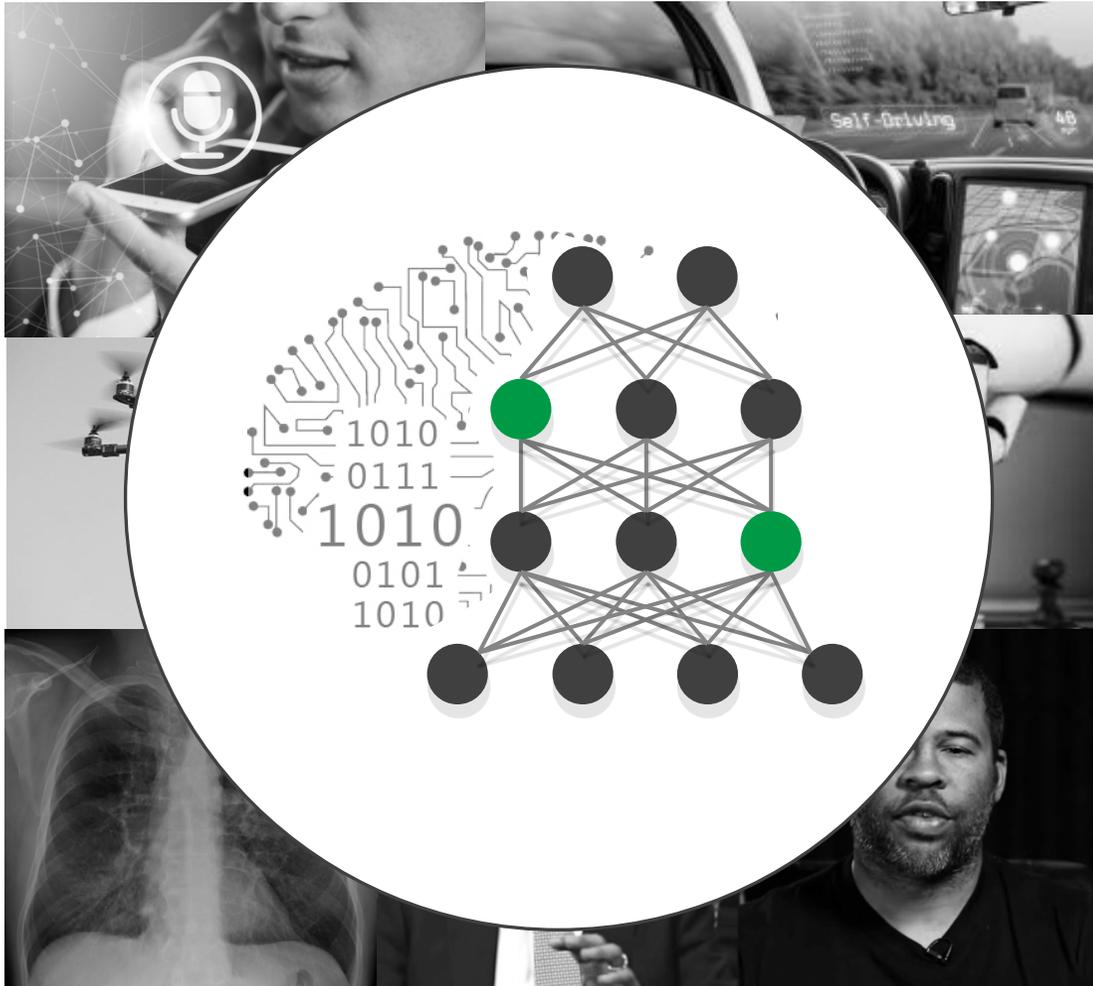
...

**AIs** / machines ☑

# Neuroscience

Using **AI**s to understand the brains of _____



humans ☑

cats ☑

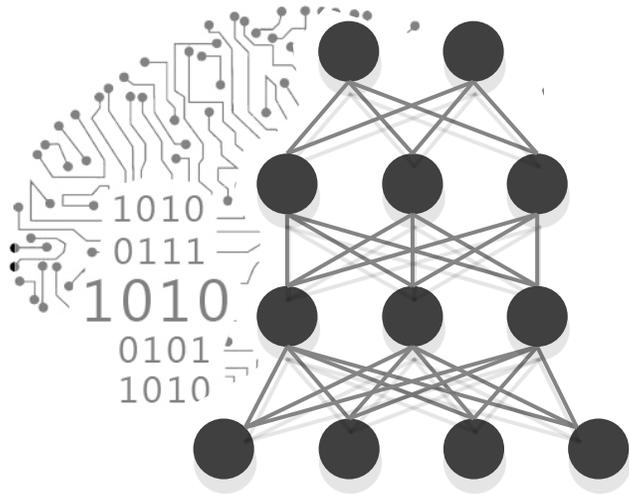mice ☑

macaques ☑

...

**AI**s / machines ☑

4

# ~~Neuroscience~~ "AI neuroscience"

Using **AIs** to understand the brains of _____



| | |
|---|---|
| humans | ☑ |
| cats | ☑ |
| mice | ☑ |
| macaques | ☑ |
| ... | |
| **AIs** / machines | ✓ |

# Subject: Image classifier



cats
mice
macaques
...
school bus
daisy

**AlexNet** (Krizhevsky et al. 2012)

# Subject: Image classifier

60M connections
500K neurons
5 conv + 5 dense layers

$f$ :

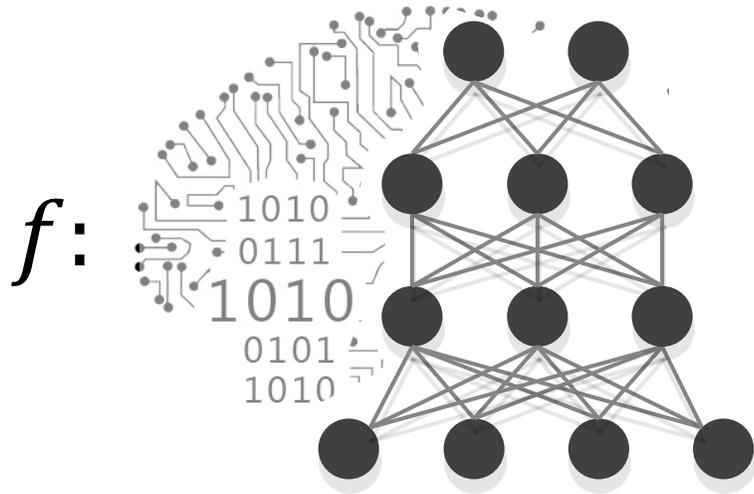**AlexNet** (Krizhevsky et al. 2012)

IM·GENET

cats         3%
mice         2%
macaques     0%
...          ...
**school bus**  92%
daisy        1%

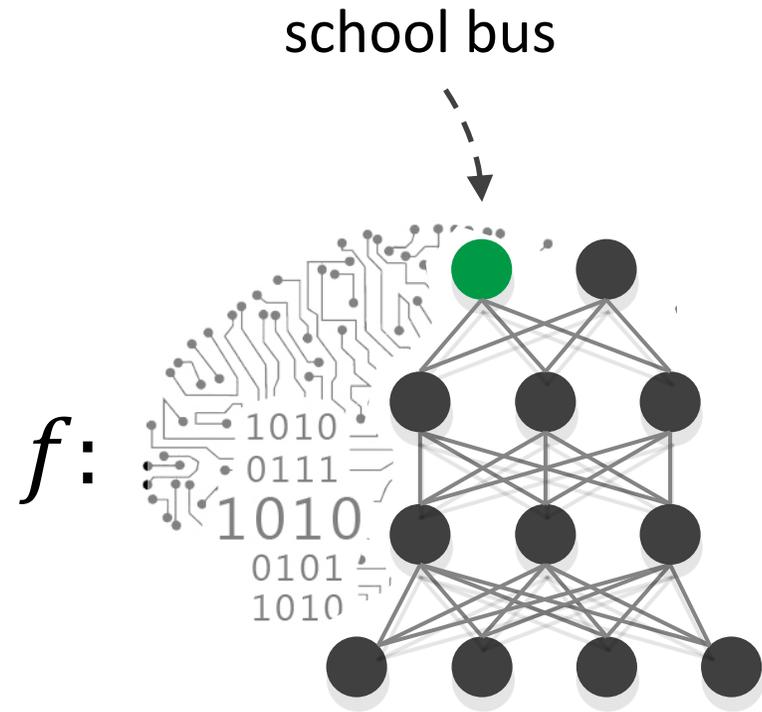# Subject: Image classifier



$f$:

**AlexNet** (Krizhevsky et al. 2012)

IMAGENET

| | |
|---|---|
| cats | 1% |
| mice | 0% |
| macaques | 2% |
| ... | ... |
| **school bus** | 95% |
| daisy | 0% |

$f$

# What does the school bus neuron want to see?



school bus

$f$:

**AlexNet** (Krizhevsky et al. 2012)

IM*A*GENET

| | |
|---|---|
| cats | 1% |
| mice | 0% |
| macaques | 2% |
| ... | ... |
| **school bus** | 95% |
| daisy | 1% |

$f$

# Finding what biological neurons want to see



Electrical signal from brain

Recording electrode

Visual area of brain

Stimulus

*Hubel & Wiesel 1954*

# Finding what <u>artificial</u> neurons want to see

## 2. Image generator
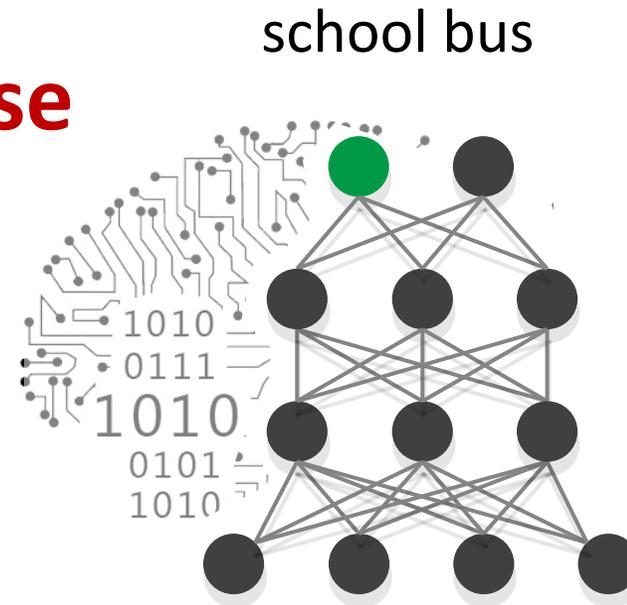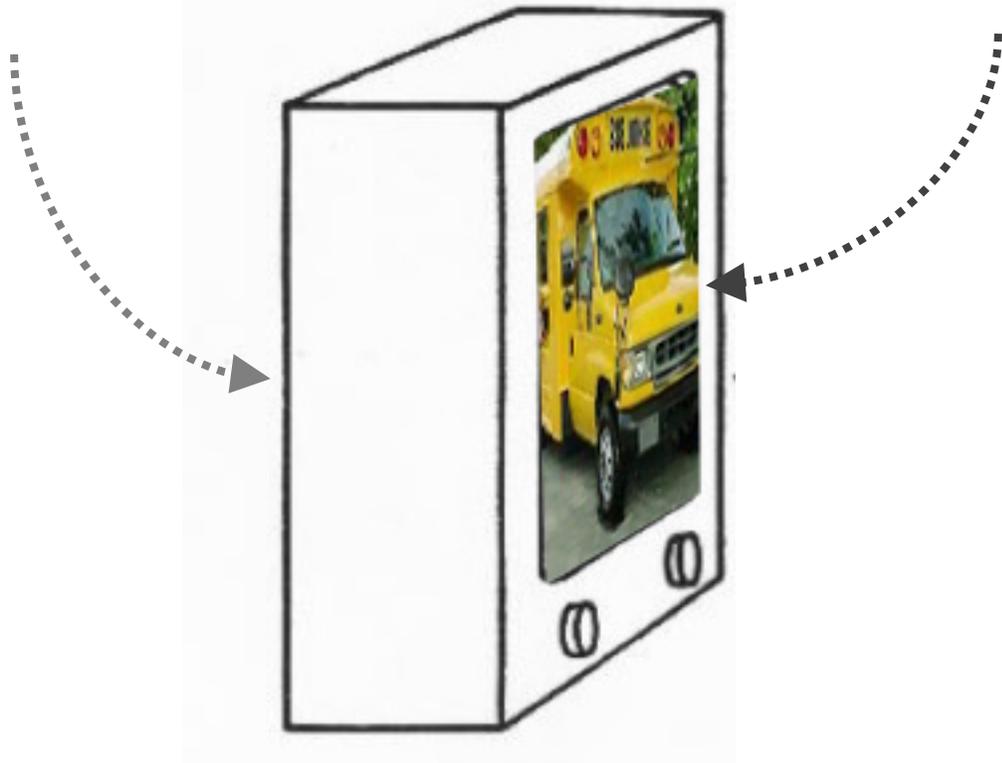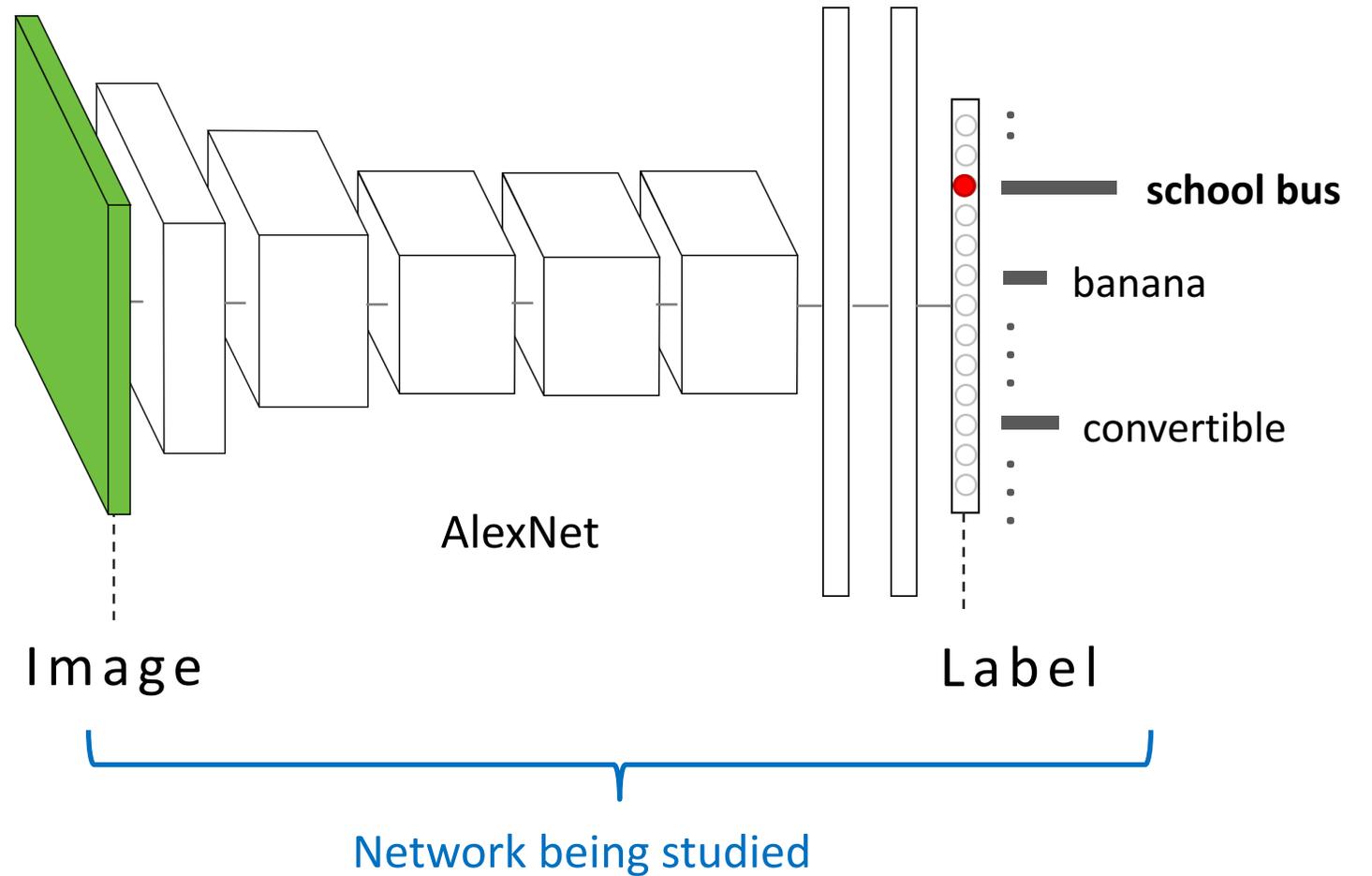## 3. 3D renderer

1. Pixel-wise

school bus

$$x^* = \arg\max_{x}(\phi_{layer,idx}(x))$$

"Activation maximization"    Erhan et al. 2009

# Finding what <u>artificial</u> neurons want to see

2. Image generator

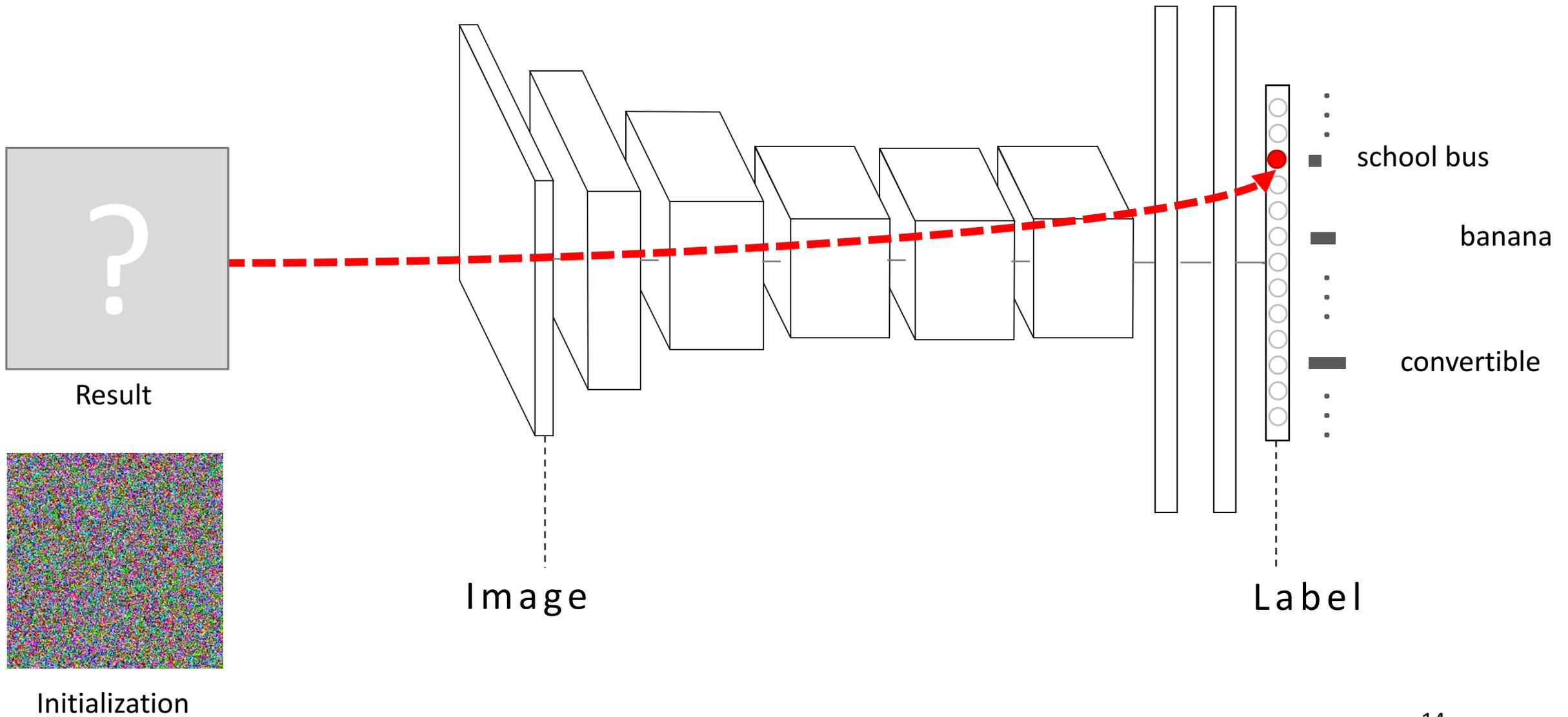3. 3D renderer

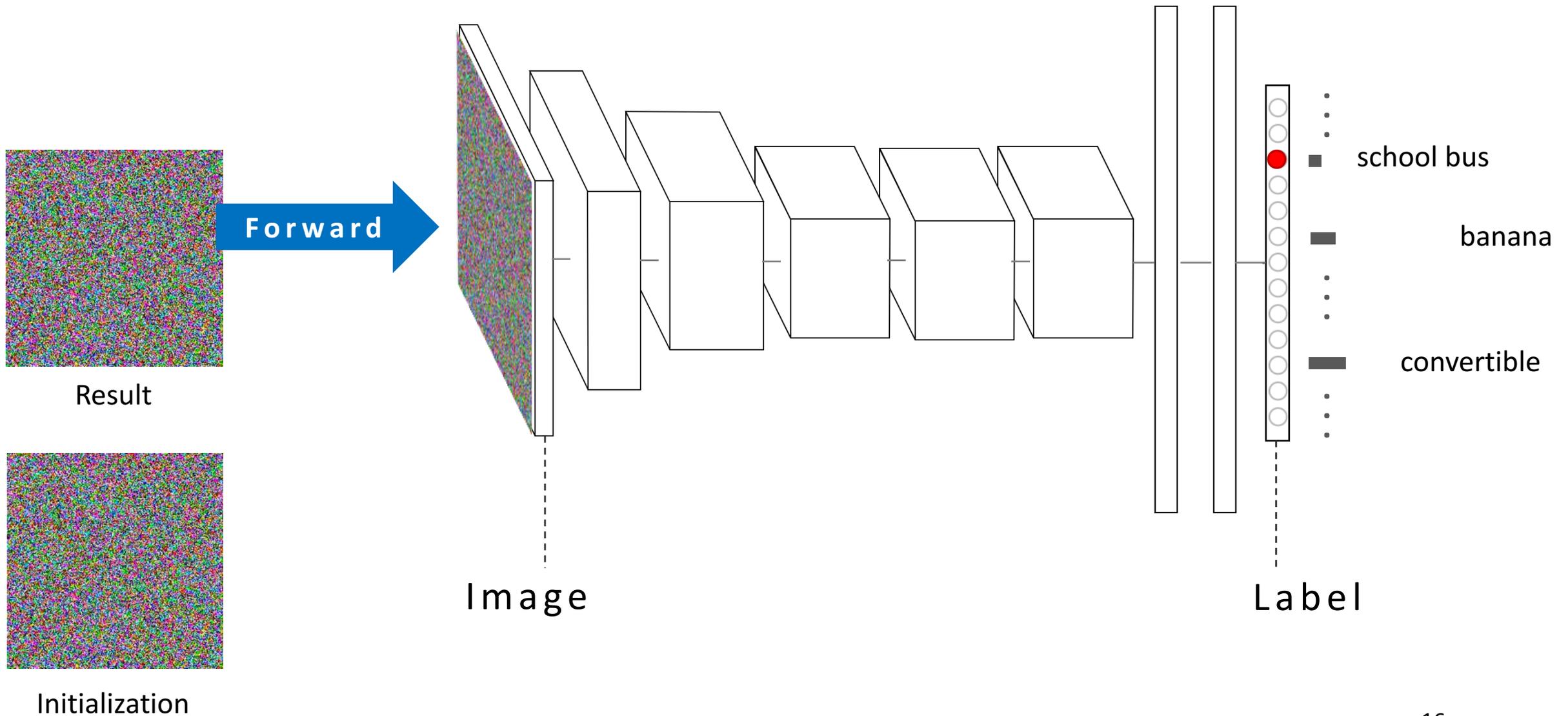**1. Pixel-wise**

school bus

# Pixel: Evolutionary algorithms



Cross-over

Mutation    **Evolution**    Evaluation

Selection

AlexNet

school bus

banana

convertible

Image

Label

Network being studied
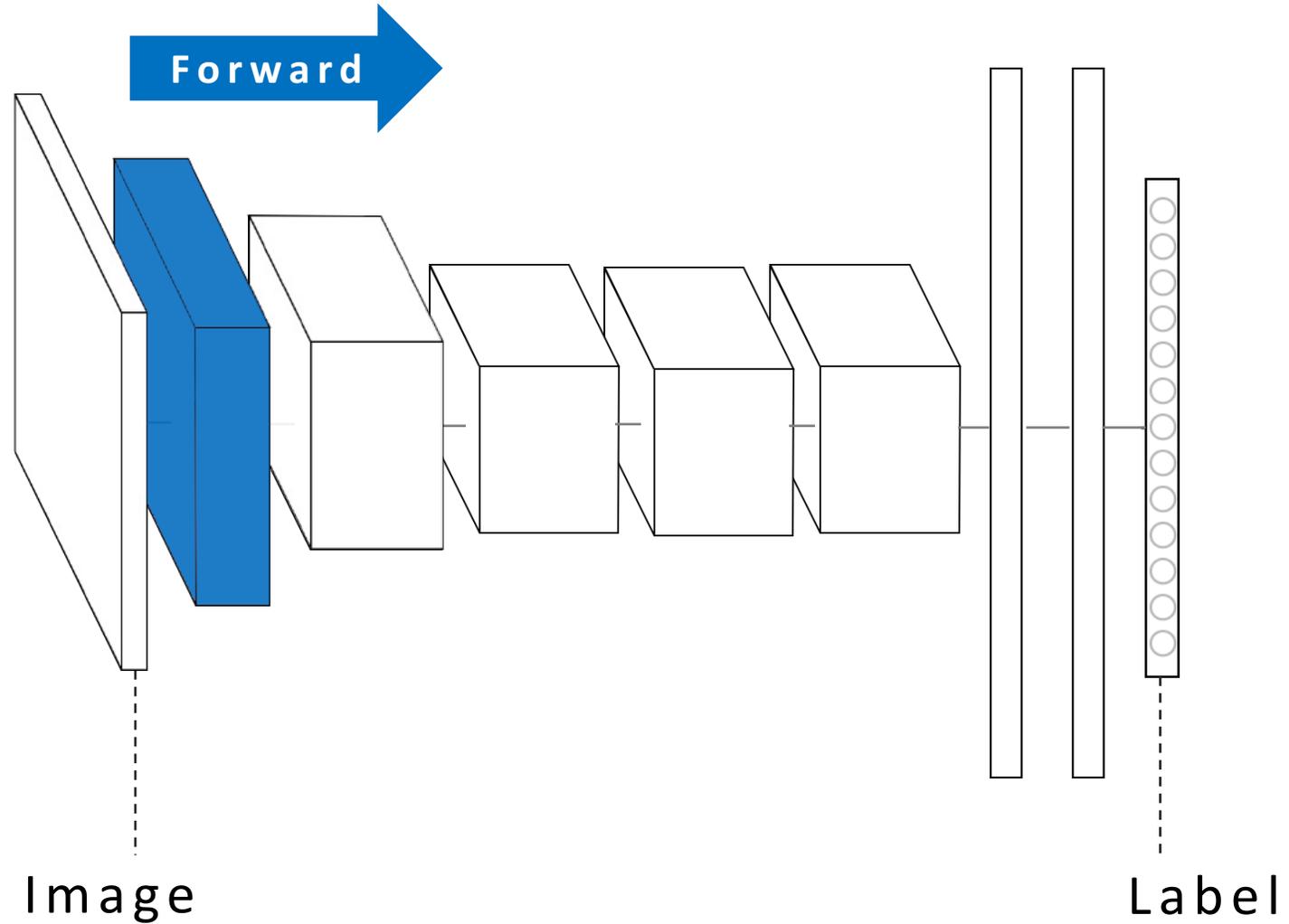
# Pixel: Gradient descent



Result

Initialization

Image

Label

school bus

banana

convertible

# Pixel: Gradient descent



Result

Initialization

Image

Label

school bus

banana

convertible

# Pixel: Gradient descent



Result

Forward

Initialization

Image

Label

school bus

banana

convertible

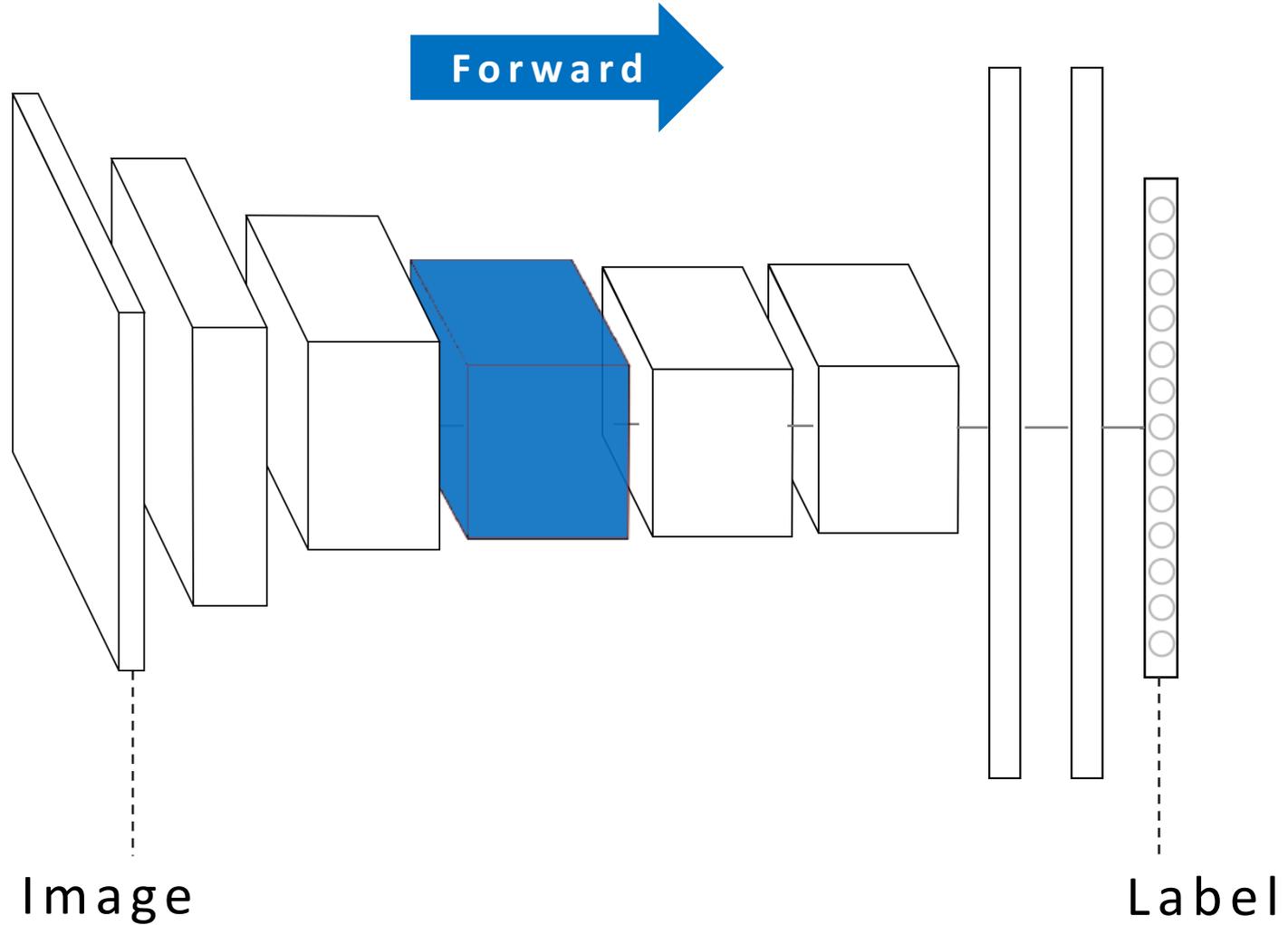# Pixel: Gradient descent



Forward

Result

Initialization

Image

Label

# Pixel: Gradient descent



Result

Initialization

Forward

Image

Label

# Pixel: Gradient descent



Forward

Result

Initialization

Image

Label

# Pixel: Gradient descent

Result

Initialization

Image

Label

# Pixel: Gradient descent



Result

Initialization

Image
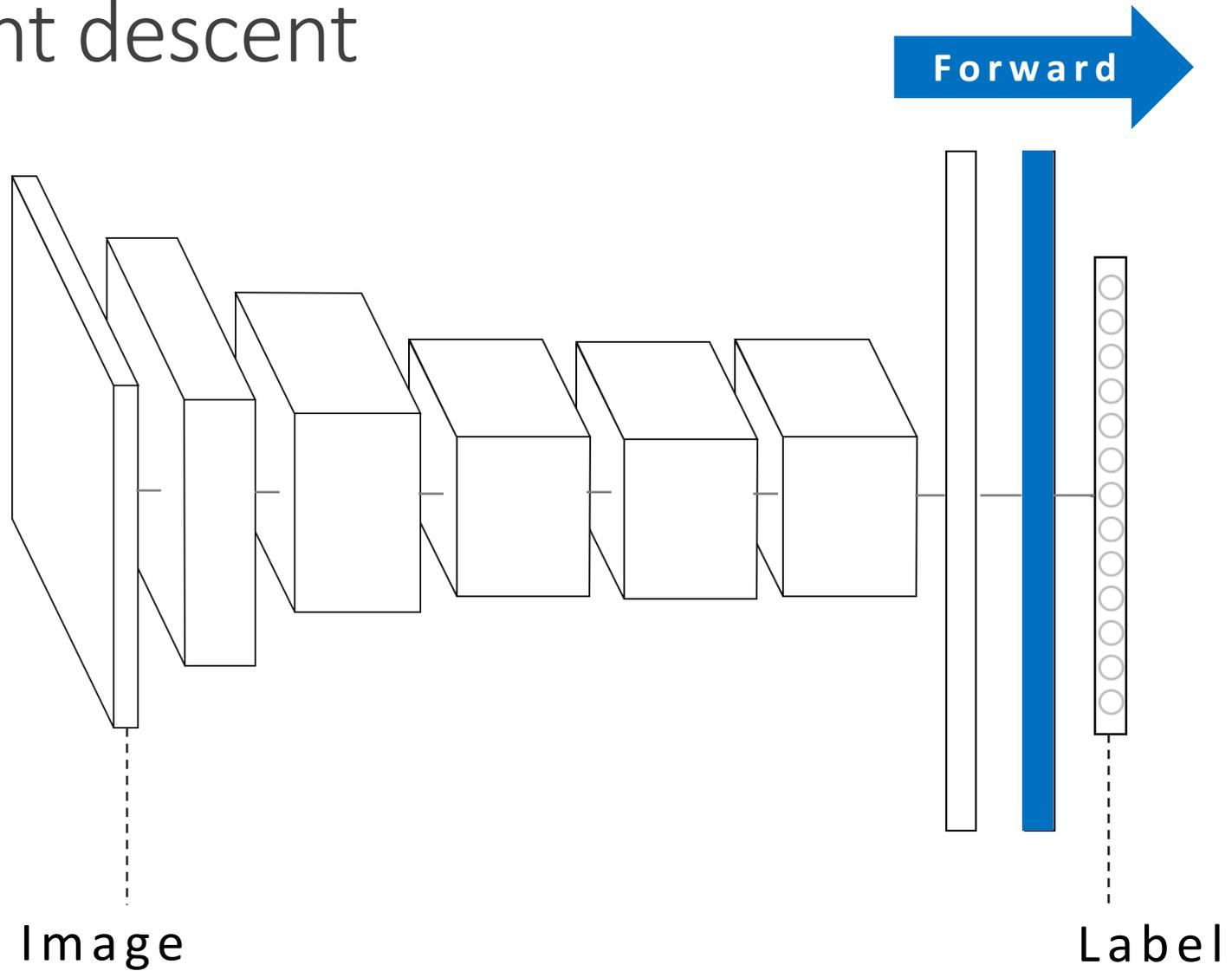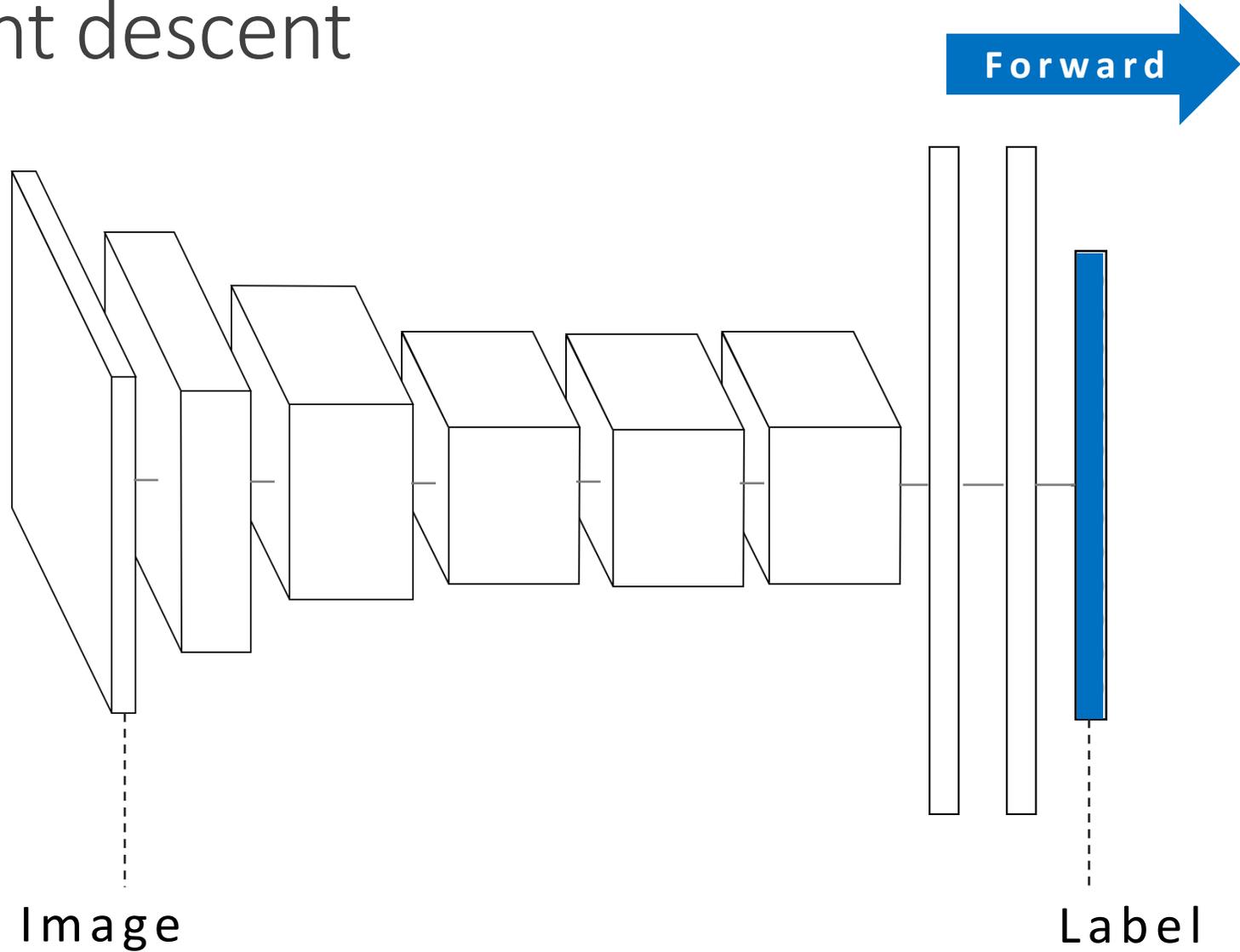
Label

school bus

banana

convertible

# Pixel: Gradient descent



Result

Initialization

Image

Label

Gradient

school bus

banana

convertible

# Pixel: Gradient descent



Backprop

<u>Gradient</u>

Result

Initialization

Image

Label

school bus

banana

convertible

# Pixel: Gradient descent



Backprop

Gradient

school bus

banana

convertible

Result

Initialization

Image

Label

# Pixel: Gradient descent



Backprop  Gradient

Result

Initialization

Image

Label

school bus

banana

convertible

# Pixel: Gradient descent



Backprop

Gradient

Result

Initialization

Image

Label

school bus

banana

convertible

# Pixel: Gradient descent



Result

Update image

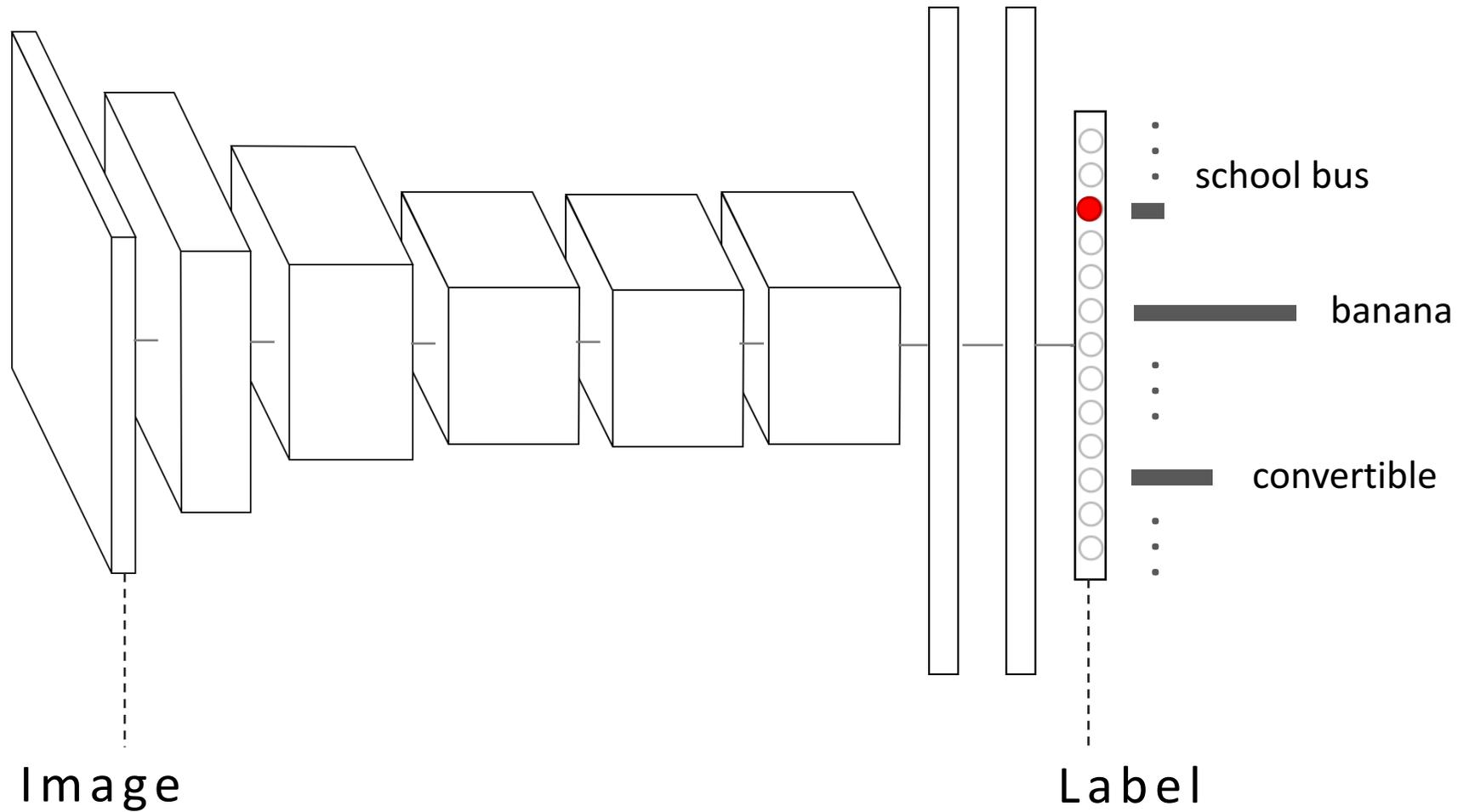Initialization

Image

Label

school bus

banana

convertible
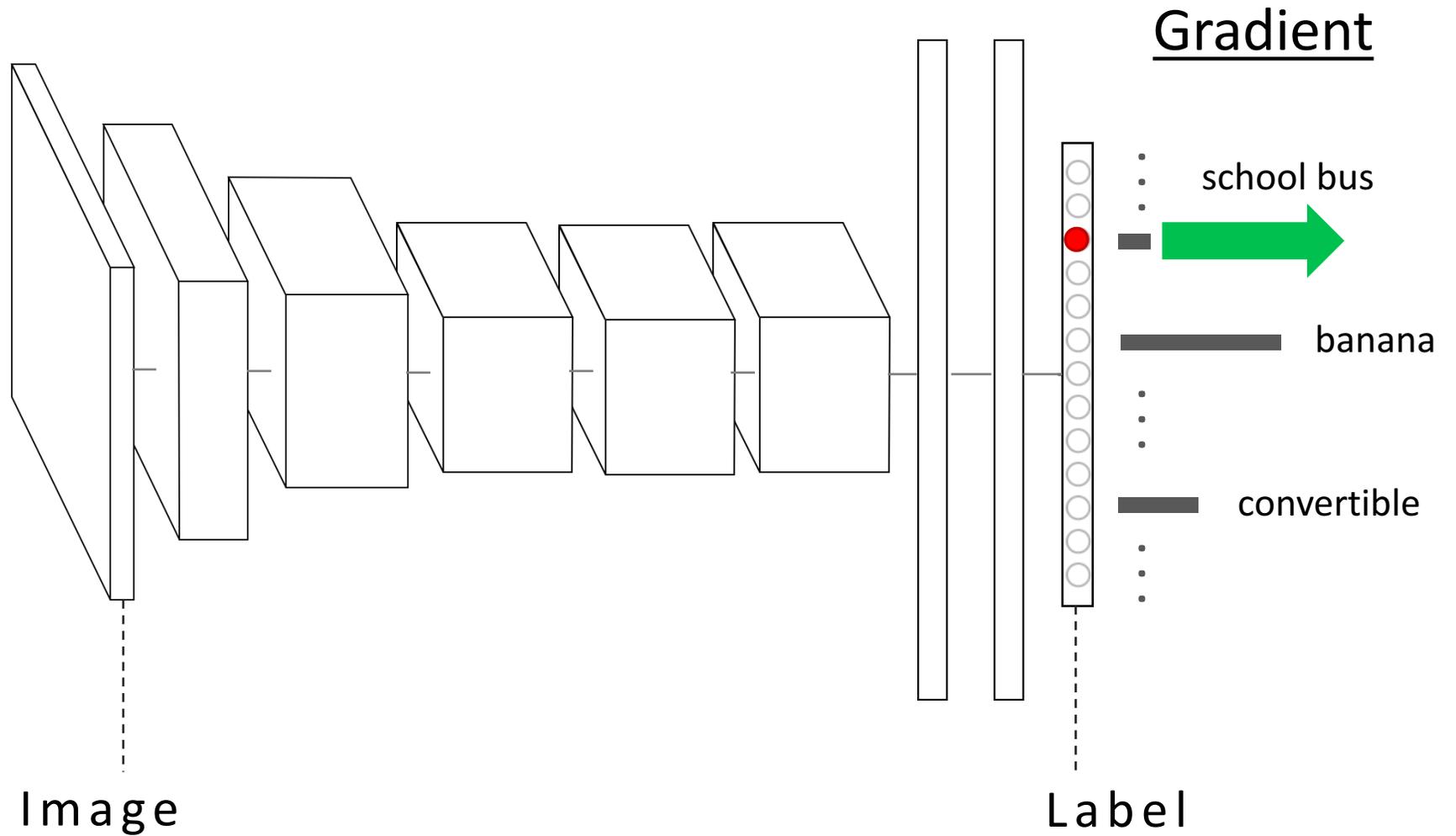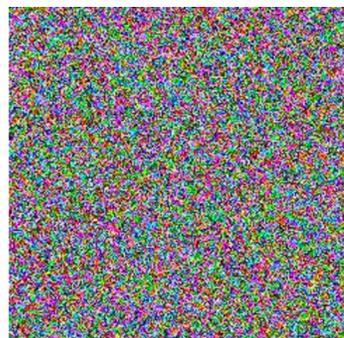
# Pixel: Gradient descent



Result

Initialization

Update image

$$x_{t+1} = x_t + \frac{\partial \phi(x_t)}{\partial x_t}$$

school bus

banana

convertible

Label

Result

Initialization

Forward

Backprop

*Many iterations*

school bus

banana

convertible

Image

Label

29

Deep Neural Networks are Easily Fooled
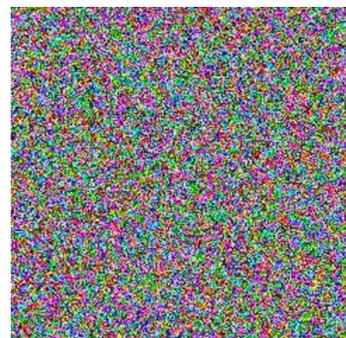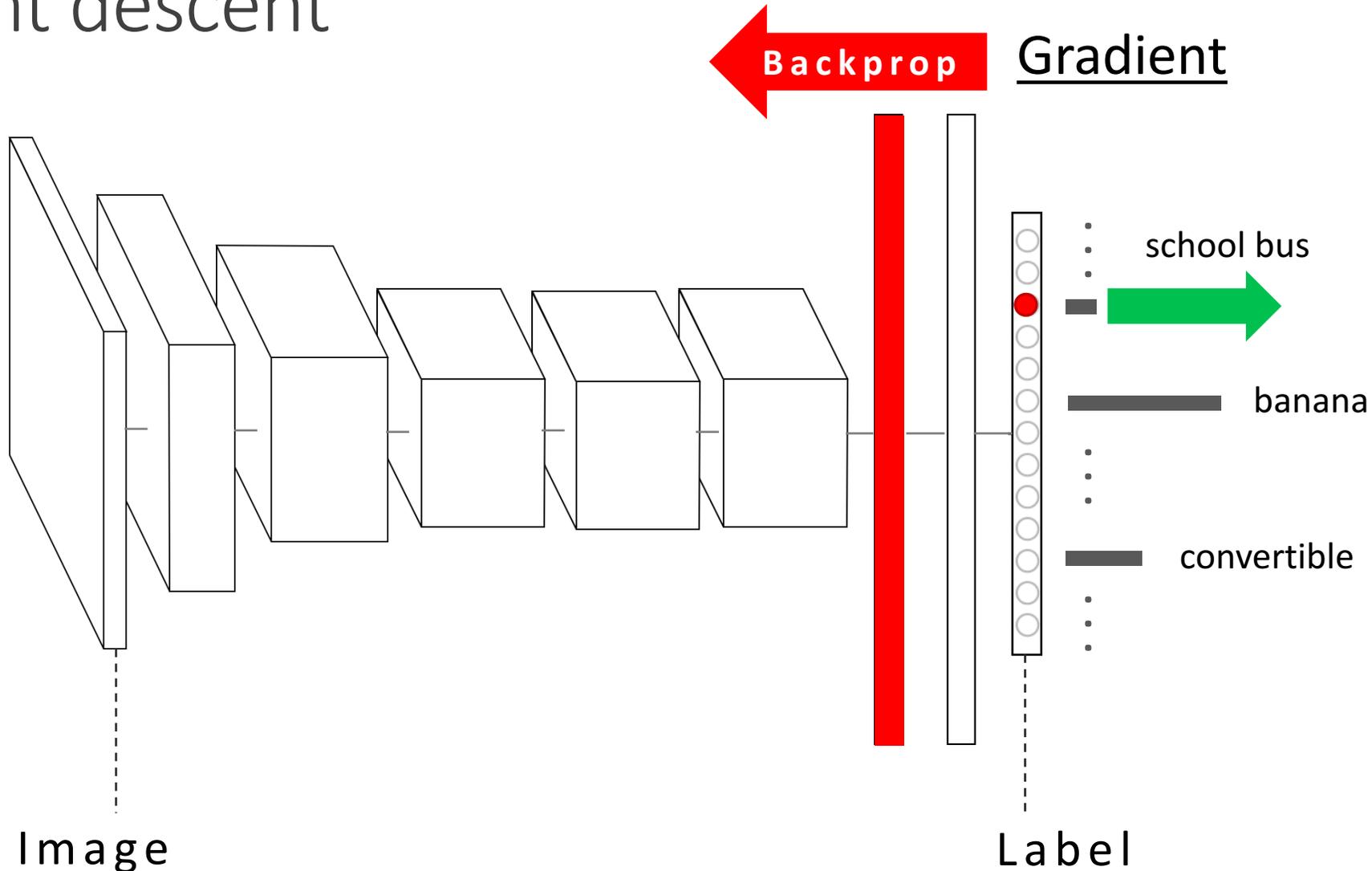


**Rubbish**

Yosinski    Clune

Result

Initialization

Tibetan terrier    golden retriever    Brittany spaniel    gorilla

chimpanzee    eel    backpack    cliff dwelling

confectionery    greenhouse    mask    parking meter

# Probabilistic interpretation

*image*          *class*

Classifier:   $x \longrightarrow y$

$$\underline{p(y|x)}$$

class

# **Problems:**

1. Poor interpretability

2. Mode collapse

Classifier: $x \longrightarrow y$

*image*  *class*

$$p(x, y) = p(x)p(y|x)$$

prior    class

Random initializations → Visualizations → Real images

# Solutions:

- total variation (Mahendran & Vedaldi, 2015)
- Gaussian blur (Yosinski et al, 2015)
- $\alpha$-norms (Simonyan et al, 2014)
- jitter (Mordvintsev et al, 2015)
- center bias (Nguyen et al, 2016)
- GMM (Mordvintsev et al, 2016)

...

*image*   *class*

Classifier:   $x \longrightarrow y$

$$p(x, y) = p(x)p(y|x)$$

prior   class

Update:   $$x_{t+1} = x_t + \frac{\partial \phi(x_t)}{\partial x_t} + \frac{\partial R(x_t)}{\partial x_t}$$

higher activation

more realistic and diverse

33

| Lemon | Keyboard | Dumbbell | Kit fox | Bell pepper | Cup | Beacon | Volcano |

**Priors**

2013

$L_2$ norm — Simonyan et al 2014

Gaussian blur — Yosinski et al 2015

Patch statistics — Wei et al 2015

Total variation — Mahendran & Vedaldi 2015

Center bias — Nguyen et al 2016

Mean image — Nguyen et al 2016

Deep generator — Nguyen et al 2016 / Nguyen et al 2017

2017

| | Lemon | Keyboard | Dumbbell | Kit fox | Bell pepper | Cup | Beacon | Volcano | **Priors** | |
|---|---|---|---|---|---|---|---|---|---|---|
| **2013** | | | | | | | | | $L_2$ norm | Simonyan et al 2014 |
| | | | | | | | | | Gaussian blur | Yosinski et al 2015 |
| | | | | | | | | | Patch statistics | Wei et al 2015 |
| | | | | | | | | | Total variation | Mahendran & Vedaldi 2015 |
| | | | | | | | | | Center bias | Nguyen et al 2016 |
| | | | | | | | | | Mean image | Nguyen et al 2016 |
| **2017** | | | | | | | | | Deep generator | Nguyen et al 2016 Nguyen et al 2017 |

# Finding what <u>artificial</u> neurons want to see

## 2. Image generator
## 3. 3D renderer

1. Pixel-wise

school bus

# Finding what <u>artificial</u> neurons want to see

**2. Image generator**

3. 3D renderer

1. Pixel-wise

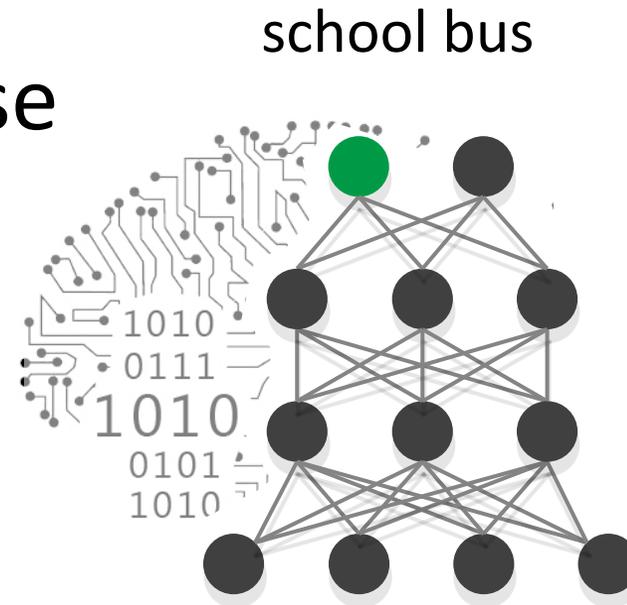school bus

Image

Label

Bell pepper

banana

convertible

DNN being visualized

Image

fc6

**Encoder:**



Image                    fc6

**Encoder**:

Image ........................ fc6

Generator / **Decoder**:

**Training losses:**

1. Reconstruction
2. Adversarial (GAN)
3. Feature matching

fc6 ........................ Image

Code

Image

Label

school bus

banana

convertible

Generator network (prior)

DNN being visualized

42

$\mathbb{R}^{4096}$

Code

Image

Label

school bus

banana

convertible

Generator network (prior)

DNN being visualized

43

Code

Image

Label

junco

banana

convertible

Generator network (prior)

DNN being visualized

44

junco

Code

Image

Label

junco

banana

convertible

Generator network (prior)

DNN being visualized

45

Real image        Pixel-wise (Nguyen et al. 2015)

**school bus**

banana

convertible

Code       Image       Label

Generator network (prior)

DNN being visualized

46

Nguyen et al. 2016    Synthesizing the preferred inputs for neurons in neural networks...



mosque | lipstick | brambling | leaf beetle | badger | toaster

library | cheeseburger | swimming trunks | barn | candle | table lamp

chest | running shoe | water jug | pool table | broom | cellphone

Dosovitskiy    Yosinski

Brox    Clune

Image

Label

**Output layer**

lipstick  brambling  leaf beetle

cheeseburger  swimming trunks  barn

running shoe  water jug  pool table

48

# Hidden neurons



Edges and colors

Layer 1

Image

Missile

banana

convertible

Label

# Hidden neurons



Layer 2

corners and textures

Missile

banana

convertible

Label

# Hidden neurons



Layer 3

dog faces, mountains, trees...

Missile

banana

convertible

Label

# Hidden neurons



lighthouse  (9)

screen  (106)

**Layer 5**

Missile

banana

convertible

Label

# Fully-connected layer 6 and 7



Image

Missile

banana

convertible

Label

Ima...

**Layer 6**

Missile

banana

convertible

Label

**Layer 6**
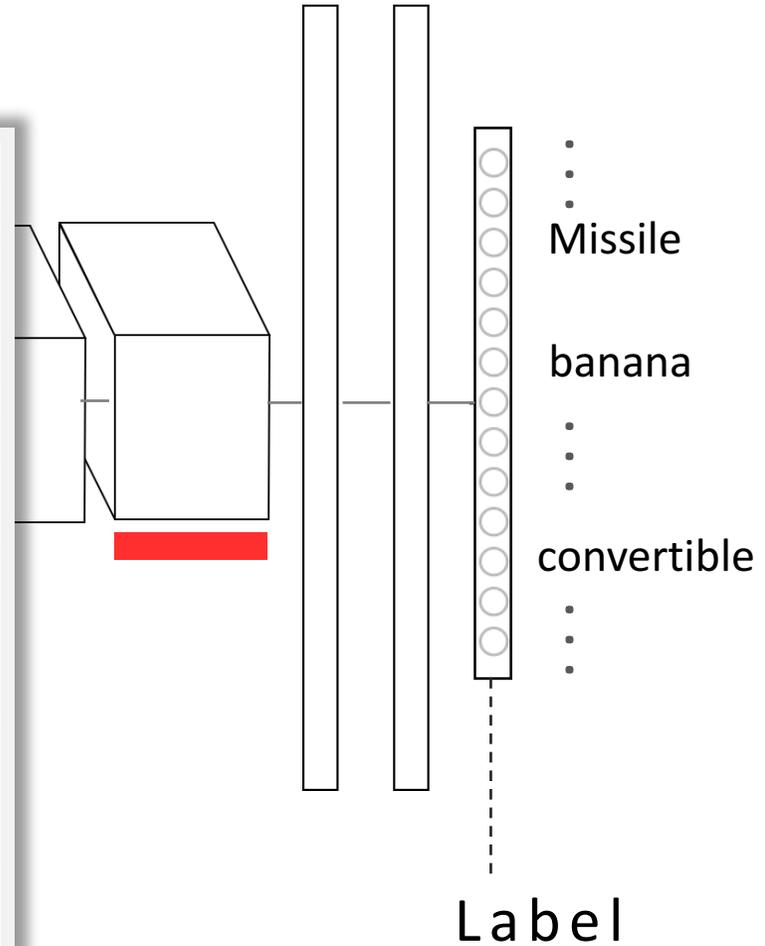
Ima...

Missile

banana

convertible

Label

# Uninterpretable stimuli – distributed coding?



Missile

banana

convertible

Label

er 6

ImageNet

Places365 ImageNet

school bus

banana

convertible

Code

Image

Label        Label

Generator network (prior)

DNN being visualized

# BigGAN-AM on Places365

**NEW**

Qi Li    Long Mai

**Real images**    **Synthesized stimuli**



plaza    plaza

hotel room    hotel room

⬚ :

**VS.** Nguyen et al. 2016, 2017
- Higher image fidelity
- Synthesizing a batch instead of a single image

⬚ :

Label

# Optimizing a set of stimuli simultaneously



ImageNet

Places365

**Latent** $\in \mathbb{R}^{140}$
- colors,
- poses,
- background

**Semantic** $\in \mathbb{R}^{128}$

256x256 BigGAN

plaza

hotel room

studio

Code

Label

Generator network (prior)

DNN being visualized

# Finding what <u>artificial</u> neurons want to see

2. Image generator

3. 3D renderer

1. Pixel-wise

school bus

# Finding what <u>artificial</u> neurons want to see

2. Image generator

**3. 3D renderer**

1. Pixel-wise

school bus

# Fine-grained control over stimuli changes



Alcorn et al. 2019

Alcorn    Qi Li    Gong    Wang    Long Mai    Jeff Ku

z_delta

Density of **school bus** predictions over object-camera distances

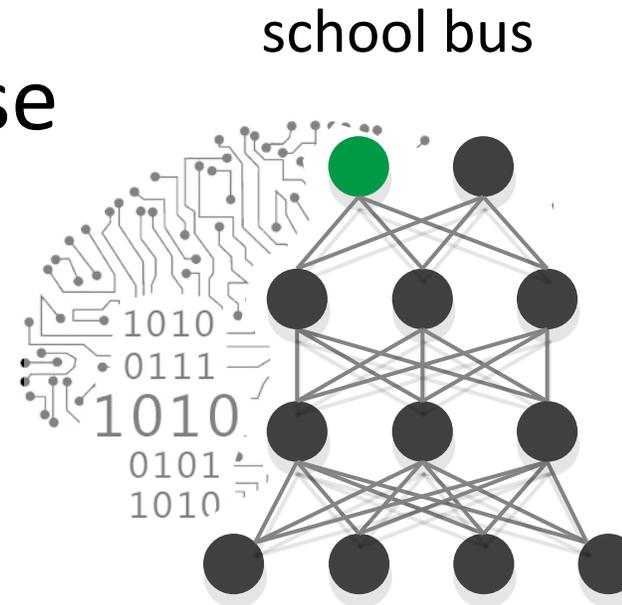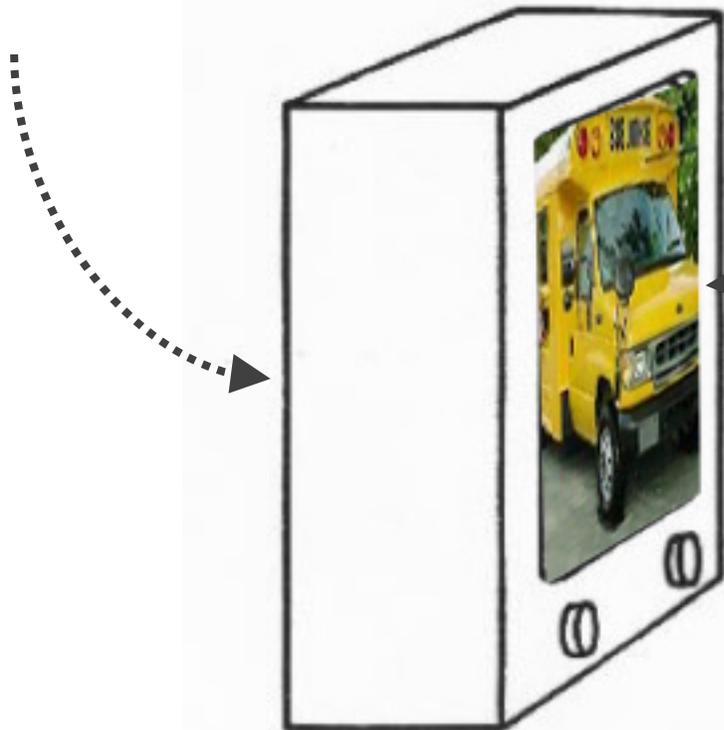**99%**

Alcorn et al. 2019
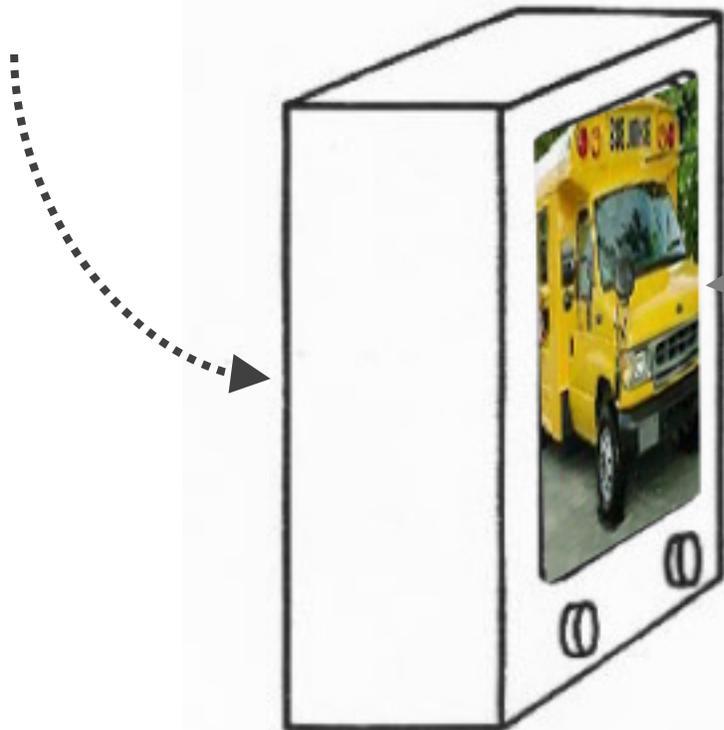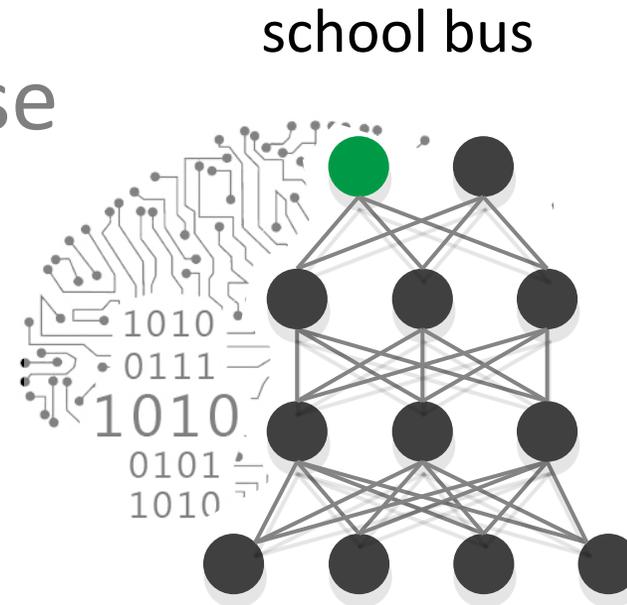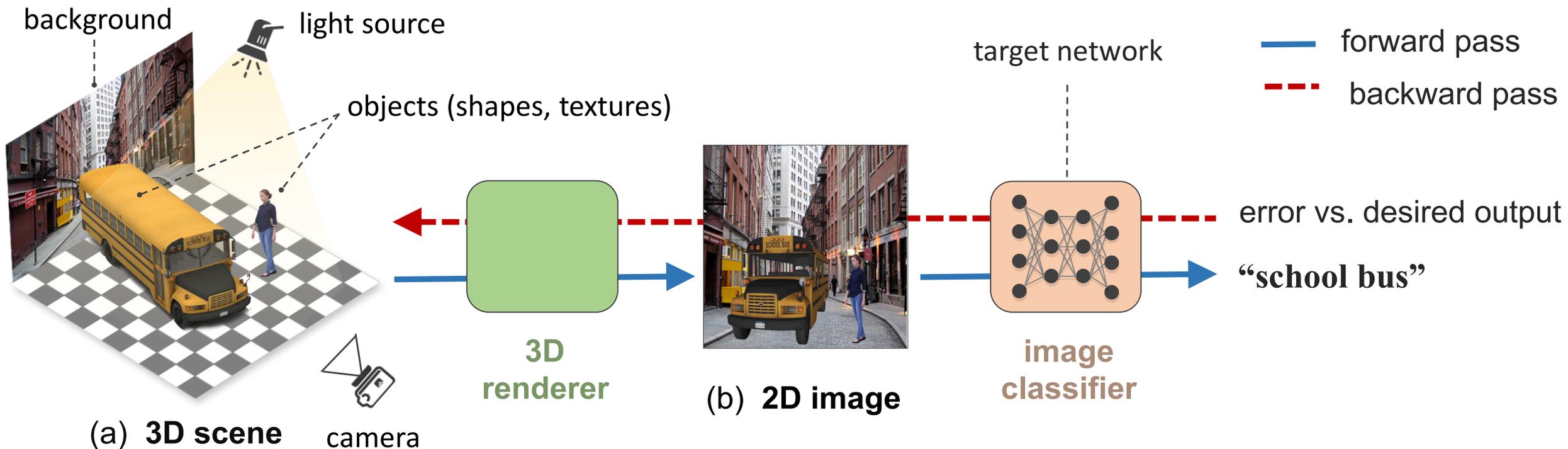
# Finding what <u>artificial</u> neurons want to see

2. Image generator

3. 3D renderer

1. Pixel-wise

school bus

# Finding what **biological** neurons want to see

Electrical signal from brain

2. Image generator

3. 3D renderer

Recording electrode

1. Pixel-wise

# Finding what biological neurons want to see

**2. Image generator**

3. 3D renderer

1. Pixel-wise



Electrical signal from brain

Recording electrode

# Using image generators to decode **real brain** signals



**fMRI**

Seen/imagined image

Feature decoder

Decoded features

Input image
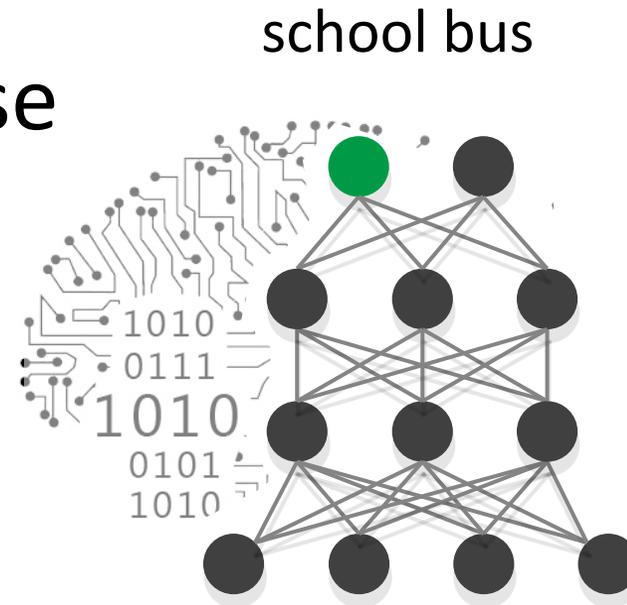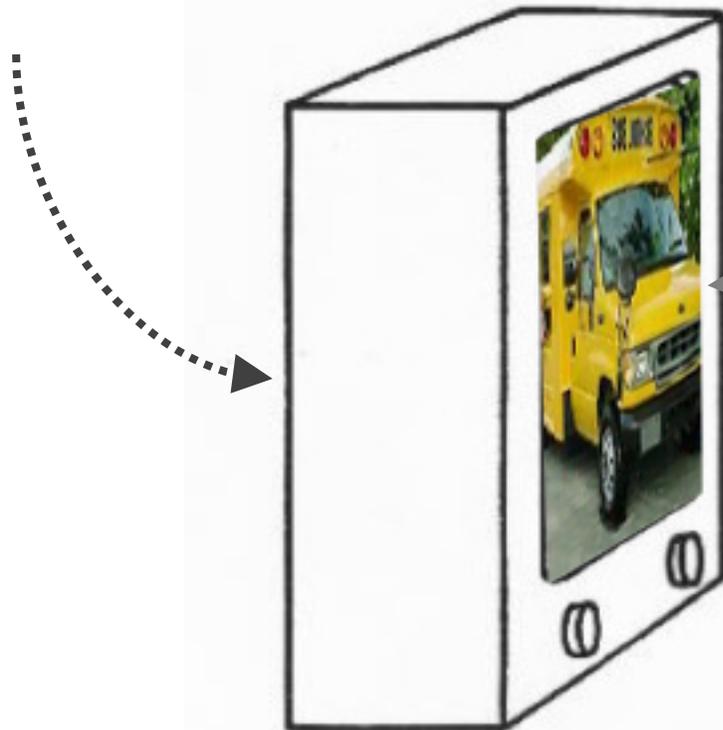
Deep generator network (DGN)

Deep neural network (DNN)

Error

Input image features

Iteratively optimize image

Reconstructed image

Shen et al. 2017

**EEG signals**

Palazzo et al. 2017

(b) Jack-o'-Lantern

**Neural activations**

Malakhova 2018

Ponce, Xiao, Schade, et al. 2019

# Thank you!



cinema 1.0    forklift 1.0

2. Image generator
3. 3D renderer

1. Pixel-wise

More info: http://AnhNguyen.me

AUBURN
UNIVERSITY

# References

○ Nguyen, A., Yosinski, J., & Clune, J. (2015). Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 427-436).

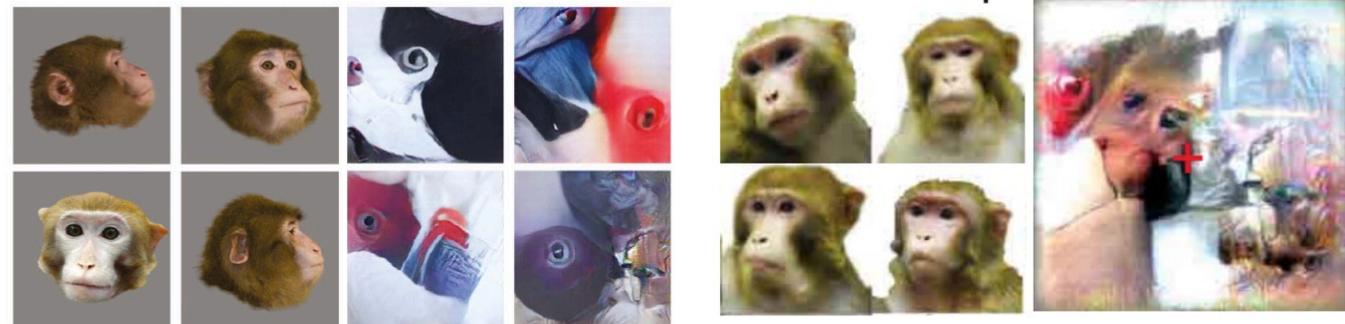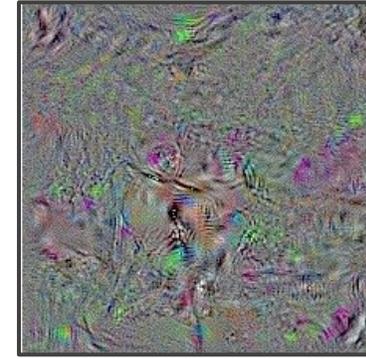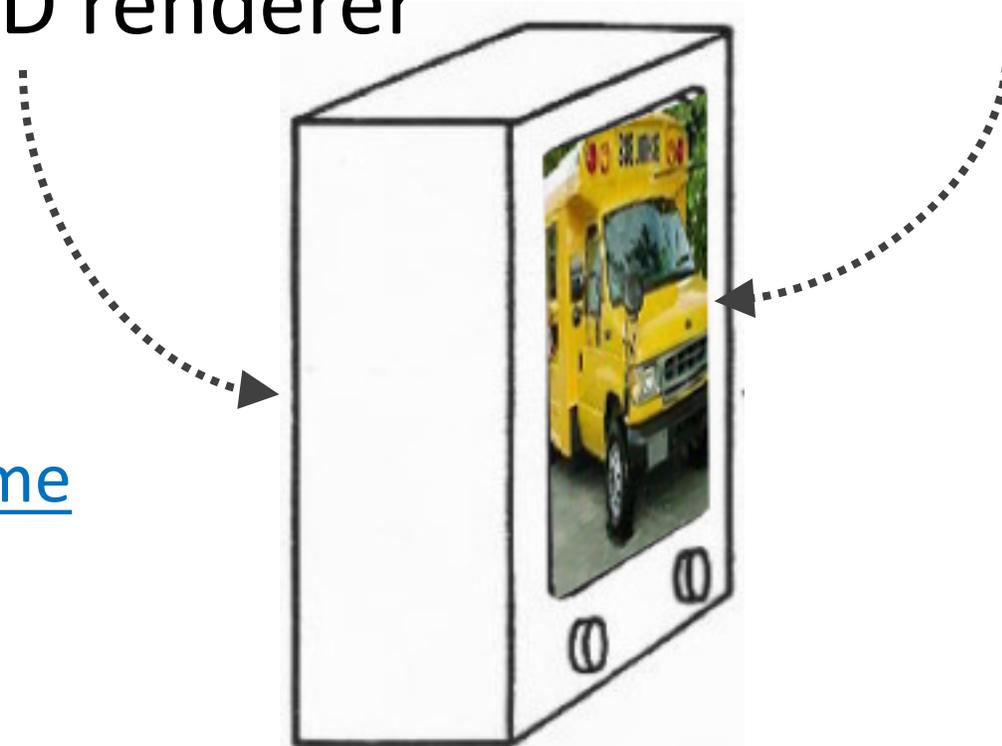○ Nguyen, A., Dosovitskiy, A., Yosinski, J., Brox, T., & Clune, J. (2016). Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In *Advances in Neural Information Processing Systems* (pp. 3387-3395).

○ Nguyen, A., Clune, J., Bengio, Y., Dosovitskiy, A., & Yosinski, J. (2017). Plug & play generative networks: Conditional iterative generation of images in latent space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 4467-4477).

○ Alcorn, M. A., Li, Q., Gong, Z., Wang, C., Mai, L., Ku, W. S., & Nguyen, A. (2019). Strike (with) a pose: Neural networks are easily fooled by strange poses of familiar objects. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (pp. 4845-4854).

○ Li, Q., Mai, L., & Nguyen, A. (2019). Improving sample diversity of a pre-trained, class-conditional GAN by changing its class embeddings. *arXiv preprint arXiv:1910.04760*.