# Image Synthesis
# (Plus some bonuses)

EECS 443 – David Fouhey

Winter 2022, University of Michigan

http://web.eecs.umich.edu/~fouhey/teaching/EECS442_W23/

StyleGAN2

StyleGAN3 (Ours)

https://github.com/NVlabs/stylegan3      [Karras et al., "Alias-Free Generative Adversarial Networks", 2021]
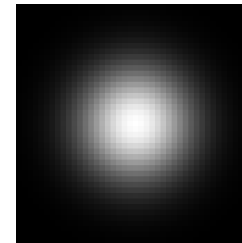
# How Many Images Are There?

- Set height and width to 1024
- Assume 256^3 (aka $2^{24}$) values per pixel
- **How many images can I create?**
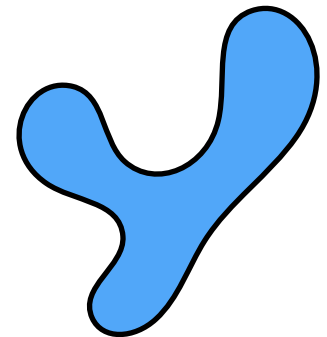- $(256^3)^{\sim 1M}$
- **Why might it be quite a bit less?**

# Learning a Mapping

- Want to learn a mapping
- **From:** a "latent" space z, often assume to be the result of sampling N-D Gaussian noise
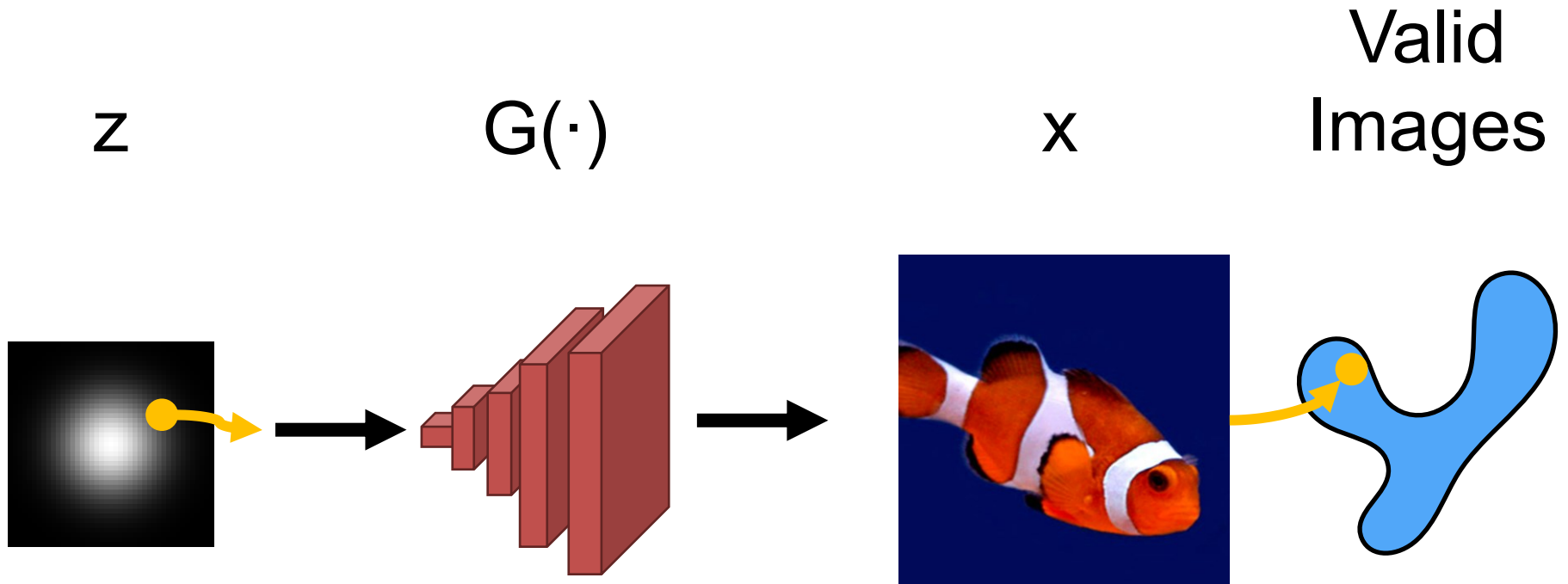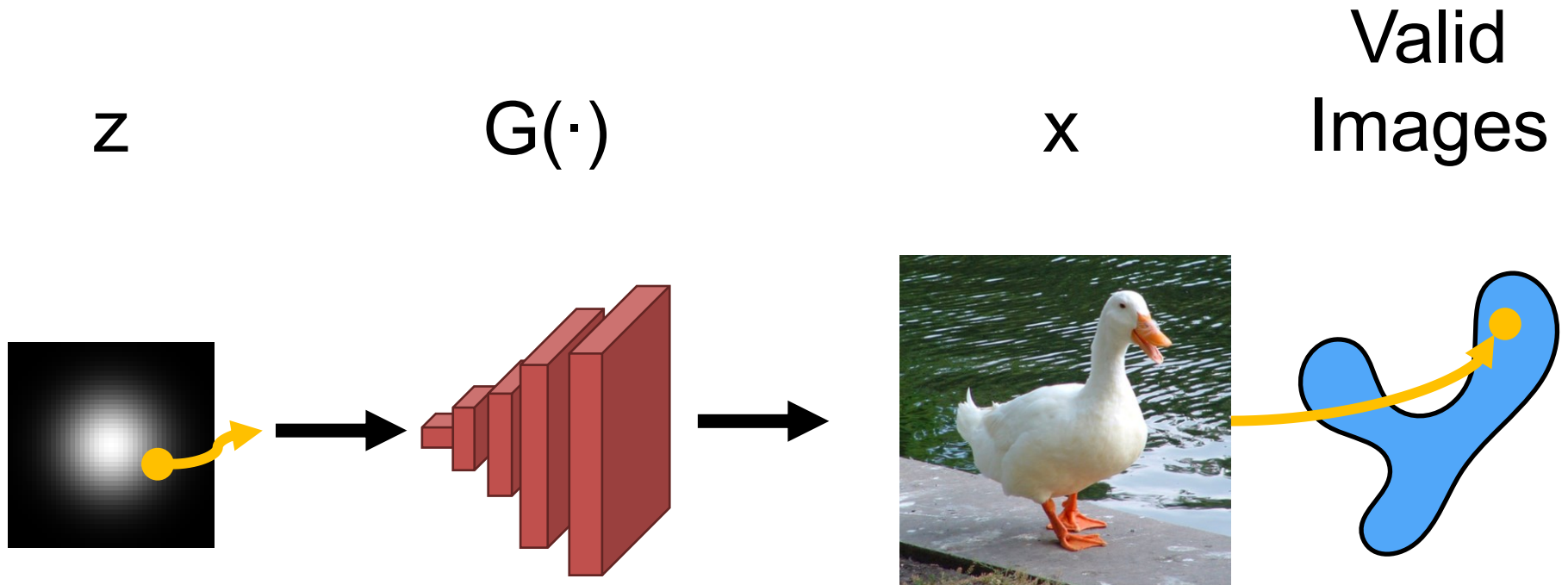- **To:** the space of valid images

Valid

z        Images

# Generating Data

z          G(·)          x          Valid
                                     Images



Remake of Slides by Isola, Freeman, Torralba, taken from Owens; diagram, duck + clownfish pictures from IFT.

# Generating Data

z                    G(·)                    x                    Valid
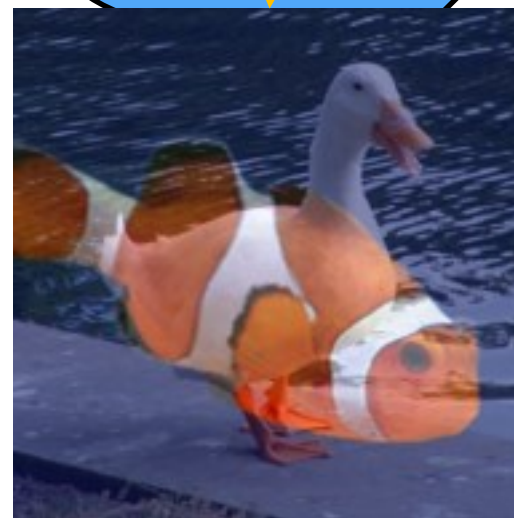                                                                  Images
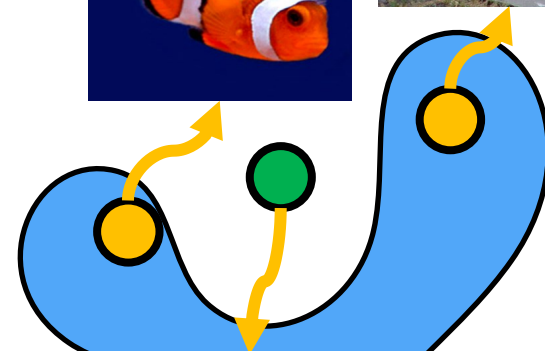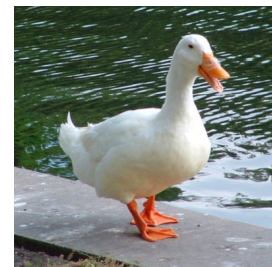
# Why The Funny Shape?

Given two valid images, what about their average?
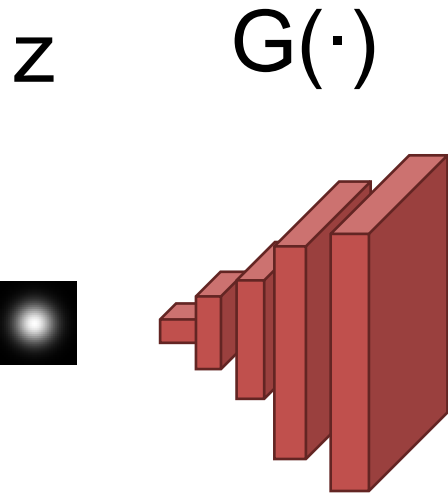
Key things to remember
- Linear combinations of images aren't images.
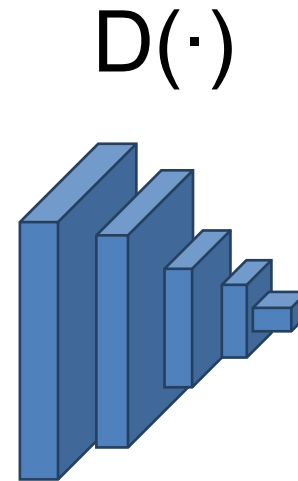- Explains funny shape ("manifold") and need for a deep network

# Generative Adversarial Networks

- Generator tries to make fake images – accepts noise and makes an image
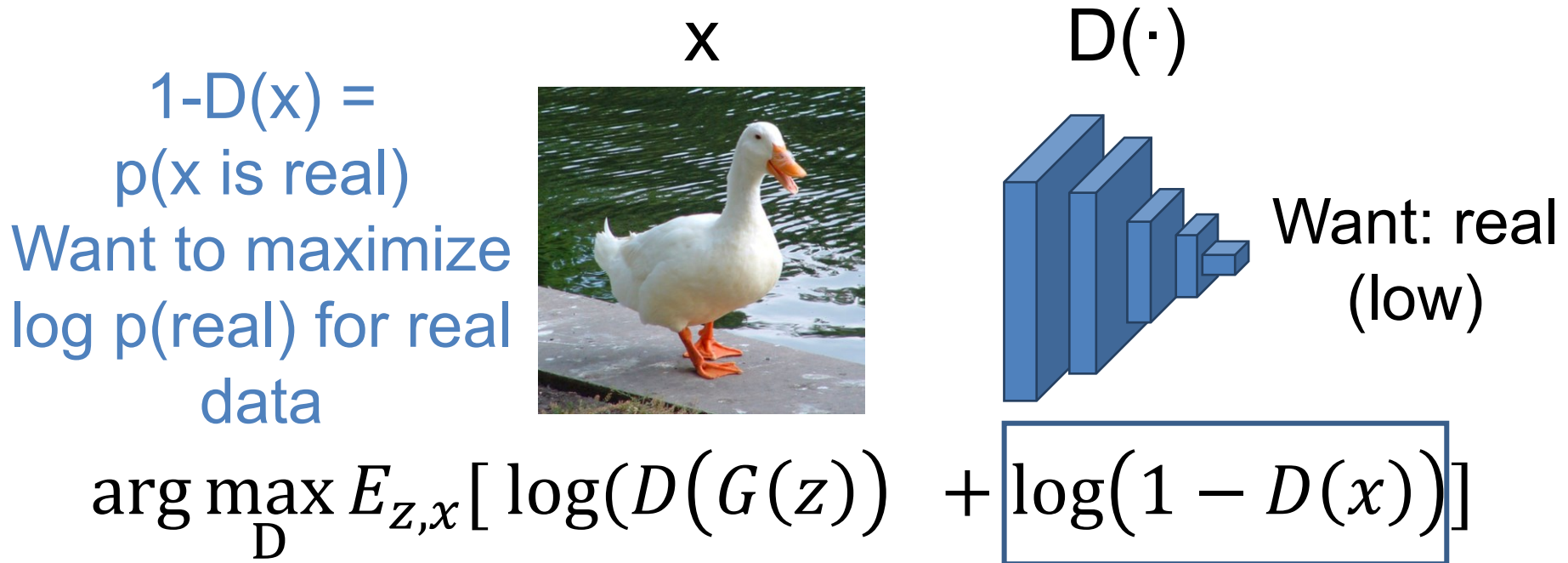- Discriminator tries to identify fakes – outputs p(fake)

Generator

z       G(·)                    G(z)
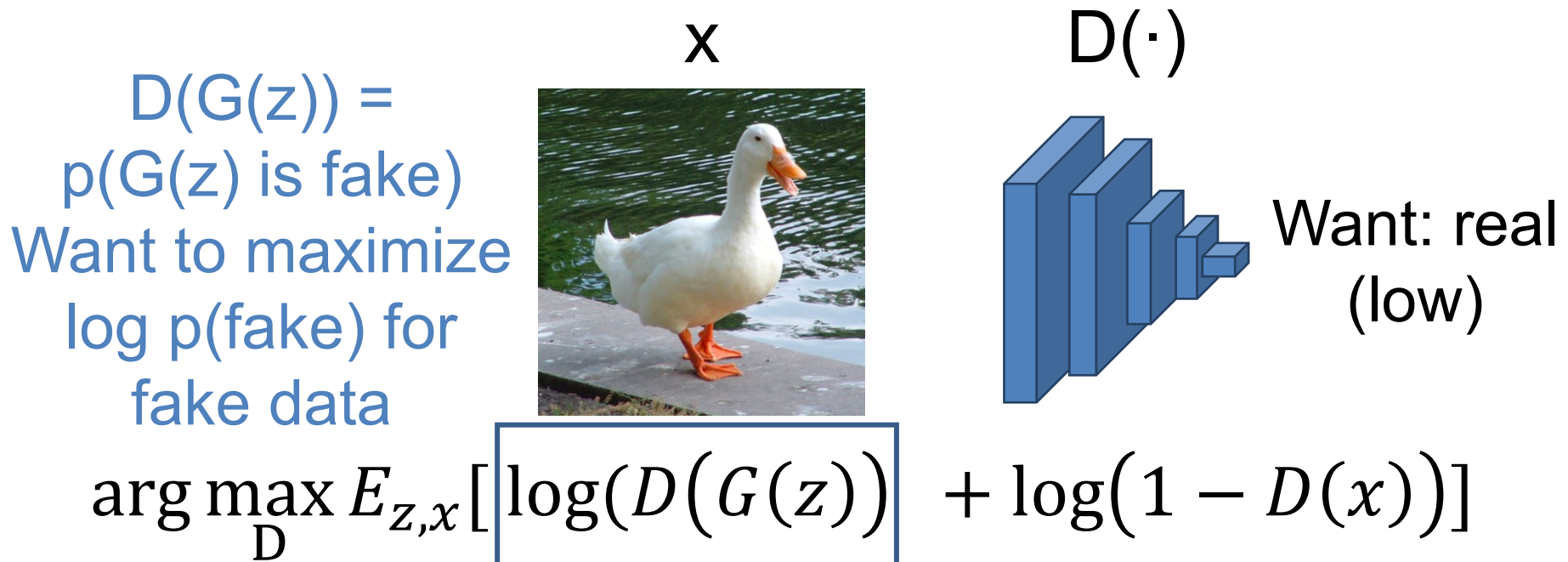
Discriminator

D(·)



Real vs fake

# Generative Adversarial Networks

z     G(·)      G(z)      D(·)



Pretend this is fixed

Want: fake (high)

D(·)

x



1-D(x) =
p(x is real)
Want to maximize
log p(real) for real
data

Want: real (low)

$$\arg \max_{\mathrm{D}} E_{z,x} [\ \log(D(G(z))\ + \boxed{\log(1 - D(x))}]$$

Remake of Slides by Isola, Freeman, Torralba, taken from Owens; diagram, duck + clownfish pictures from IFT.

# Generative Adversarial Networks

z     G(·)     G(z)     D(·)

Pretend this is fixed

**FAKE!**

Want: fake (high)

---

x     D(·)

D(G(z)) = p(G(z) is fake) Want to maximize log p(fake) for fake data

Want: real (low)

$$\arg\max_{D} E_{z,x}[\,\boxed{\log(D(G(z))}\, + \log(1 - D(x))]$$

# Generative Adversarial Networks

z     G(·)        G(z)        D(·)



Pretend this is fixed

Want: fake (low)

---

Goal of generator G: fool the discriminator D while getting to use gradients from D
*Analogy*: art forger and art detective

**How good are you at spotting forgeries?**

$$\arg \min_G E_{z,x}[\, \log(D\big(G(z)\big) \; + \log\big(1 - D(x)\big)]$$

# Generative Adversarial Networks

z     G(·)        G(z)        D(·)



Want: fake (high)

---

Final goal: find the generator that fools the best D that you could find.

In practice, important not to let the discriminator get too good. **Why?**

Theory: optimum when G produces distribution

$$\arg \min_G \max_D E_{z,x}[\, \log(D(G(z)) \; + \log(1 - D(x)))]$$

# Revisiting Averages

Can use z to walk latent space
$G(\alpha z_1 + (1-\alpha)z_2)$ for $\alpha$ in $[0,1]$



$x_1$

$x_2$

$z_1$

$z_2$

StyleGAN2

StyleGAN3 (Ours)

[Karras et al., "Alias-Free Generative Adversarial Networks", 2021]

StyleGAN2

StyleGAN3 (Ours)

[Karras et al., "Alias-Free Generative Adversarial Networks", 2021]

# Conditioning on Things

- Turning noise into pictures of things is all fun and good, but what if we want control over our synthetic images?

# Conditional GANs (Pix2Pix)

- Generator tries to make fake images – accepts **image** and makes an image
- Discriminator tries to identify fakes – outputs p(fake), potentially at each pixel

G(·)        G(x)        D(·)



Real
v
Fake

# Conditional GANs – Discriminator

G(·)　　　G(x)　　　D(·)



Pretend this is fixed

Want: Fake

---

y　　　D(·)

Want real outputs y to be low, fake output G(x) high



Want: Real

$$\arg\max_{D} E_{z,x}[\ \log(D\big(G(x)\big)\ + \log\big(1 - D(y)\big)]$$

Remake of Slides by Isola, Freeman, Torralba, taken from Owens.

# Conditional GANs – Generator

G(·)    G(x)    D(·)



Pretend this is fixed

Want: Real

If you're the generator, want to make fakes that fool the discriminator into think they're real

$$\arg \min_{G} E_{z,x}[\, \log(D(G(x)) \ + \log(1 - D(y))]$$

Same min/max game as before

$$\arg \min_{G} \max_{D} E_{z,x}[\, \log(D(G(x)) \ + \log(1 - D(y))]$$

# One Catch

- G can just output random good images.
- Solution – make the D look at the input too

G(·)        G(x)        D(·)



Real
v
Fake

# One Catch

G(·)          G(x)          D(·)



Real
v
Fake

# More Broadly

- Neural networks are lazy and will do precisely what you ask and no more

- You *have to* be careful what you ask them to do



3. Find x.

3 cm

4 cm

Here it is ✓ -1

| Input | Output | Groundtruth |
|---|---|---|



Data from
[maps.google.com]

Input                    Output                    Groundtruth

Slide by Isola, Freeman, Torralba, taken from Owens.

# Input

# L1 loss



# Why is it blurry junk?  Hold that thought!

## Input

## L1 loss + discriminator

# Let's Talk About Blurry Pictures

# What Color Is This Bird?



To make things more concrete: what color is this the pixel under this gold circle?

# What Color Is This Bird?



Many options. **What minimizes mean-squared error?**
**What minimizes the L1 distance?**

Example credit and diagrams: Richard Zhang

# What Color Is This Bird?

Option 1: Discretize / quantize

*Before learning*: assign pixel nearest color index

*After*: convert index to value

Works because network can express its uncertainty. Loads of details if you want a good version.

# Discretized Values – Angles

Imagine predicting an angle from 0° to 180°. Having bins enables:

- Expressing bimodal distributions (e.g., either 30° or 150°)
- Getting a confidence from the prediction



0°                    180°

# Discretized Values – Angles

## *SDO*/HMI Pipeline

## Neural Net Emulation



2016 May 10
06:48:00 TAI

0        degrees        180

# Discretized Values – Angles



SDO/HMI
Pipeline

Lower
90% CI

Upper
90% CI

2016 May 10
06:48:00 TAI

0          degrees          180

# Option 2 – GAN



GAN – the discriminator will prevent us from making birds grey or brown. **Why?**

# What to Take Away

- Be careful what you ask a deep net to solve.
- The objective you're asking it to solve bakes in assumptions
- Most solutions broken in one way or another
- Deep learning is not magic

# Aside: Perceptual Losses

G(·)          G(x)



y

Conventionally, minimize distance in pixel space:
$$\|G(\boldsymbol{x}) - \boldsymbol{y}\|$$

See Johnson et al. Perceptual Losses for Real-Time Style Transfer and Super-Resolution, ECCV 2016.

# Aside: Perceptual Losses

G(·)  G(x)  F(·)  F(G(x))



y



F(y)

Instead measure distance after passing through network
$$\|F(G(\boldsymbol{x})) - F(\boldsymbol{y})\|$$

See Johnson et al. Perceptual Losses for Real-Time Style Transfer and Super-Resolution, ECCV 2016.

# And Now For Something Completely Different

# ImageNet + Deep Learning



Beagle

- Image Retrieval
- Detection
- Segmentation
- Depth Estimation
- …

Slide Credit: C. Doersch

# ImageNet + Deep Learning



Beagle

Materials?
Pose?
Parts?
Geometry?
Boundaries?

*Do we even need semantic labels?*
*Do we need this task?*

# To Make It Super Clear

- w = **weights_from_somewhere_else**
- for batch in batches:
  - inputs, labels = batch
  - calculate gradient of loss function with respect to w applied to samples in inputs
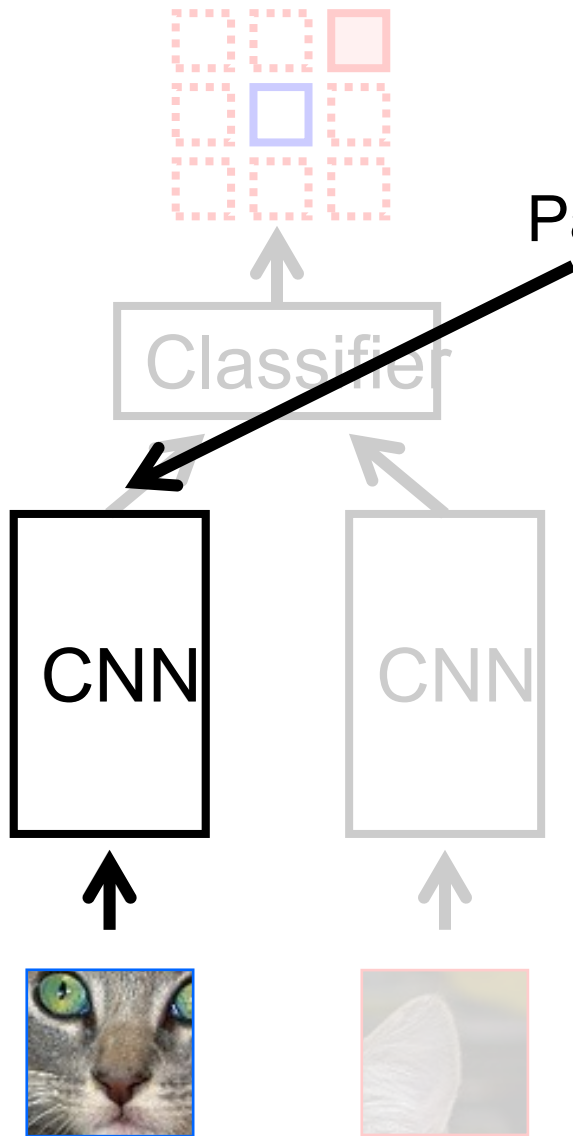  - w += gradient

# Context as Supervision
## [Collobert & Weston 2008; Mikolov et al. 2013]

# Context Prediction for Images



A    B
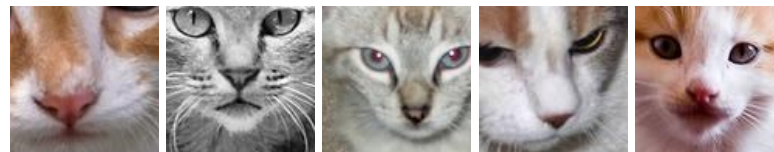
# Semantics from a non-semantic task

# Relative Position Task



8 possible locations

Classifier

CNN

CNN

Randomly Sample Patch

Sample Second Patch

Slide Credit: C. Doersch

Patch Embedding

Classifier

CNN CNN

Input          Nearest Neighbors

Note: connects *across* instances

# Avoiding Trivial Shortcuts



Include a gap

Jitter the patch locations

Slide Credit: C. Doersch

# A Not-So "Trivial" Shortcut



Position in image

# Chromatic Aberration

# Chromatic Aberration

# What is learned?

Input  Ours  Random Initialization  ImageNet AlexNet



Slide Credit: C. Doersch

# Pre-Training for R-CNN



warped region
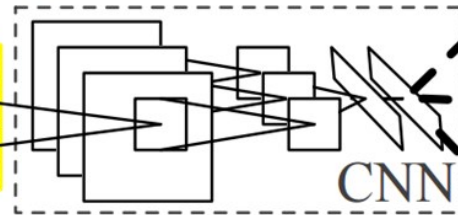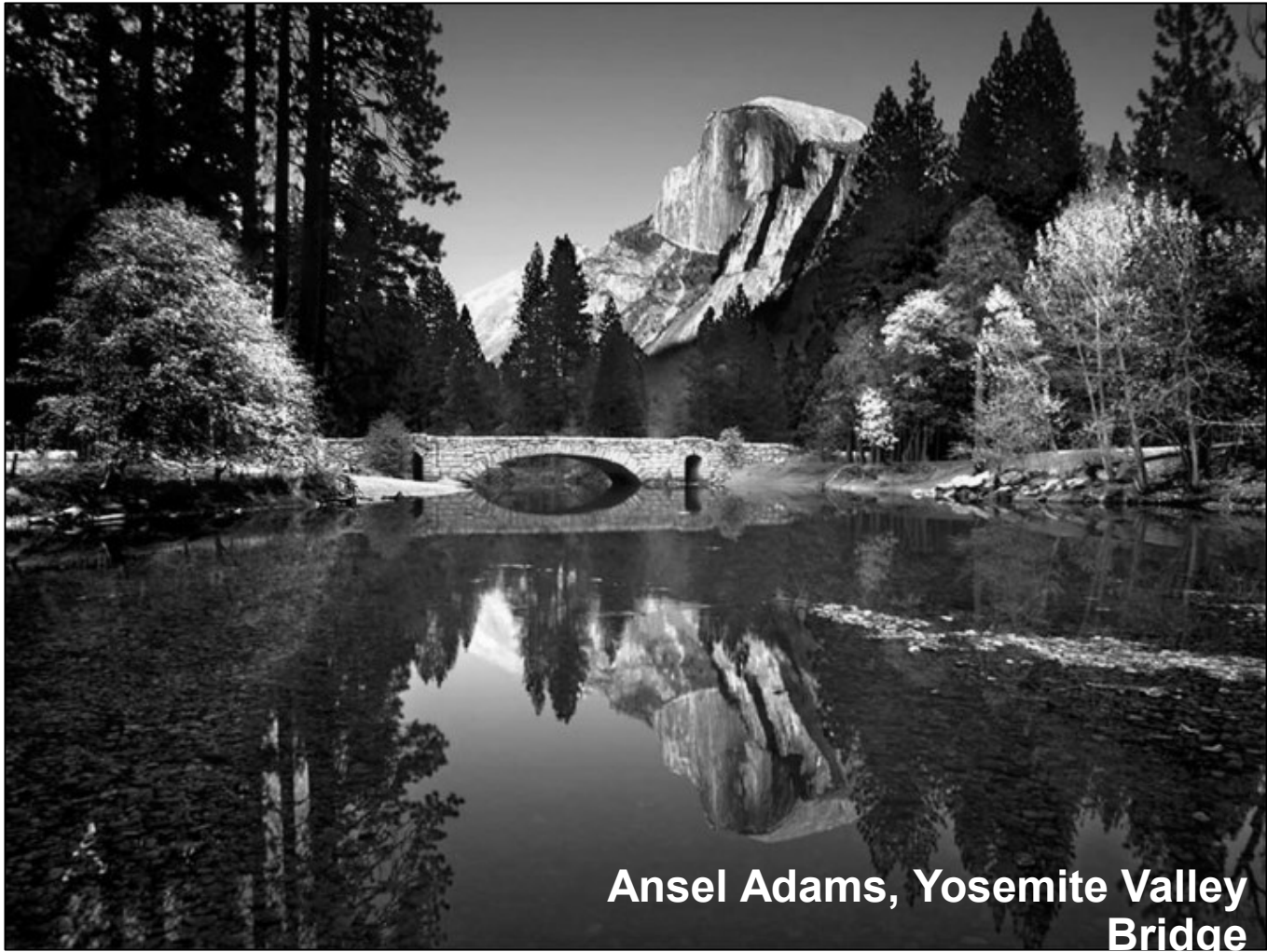
1. Input image

2. Extract region proposals (~2k)

3. Compute CNN features

CNN

aeroplane? no.

person? yes.

tvmonitor? no.

4. Classify regions

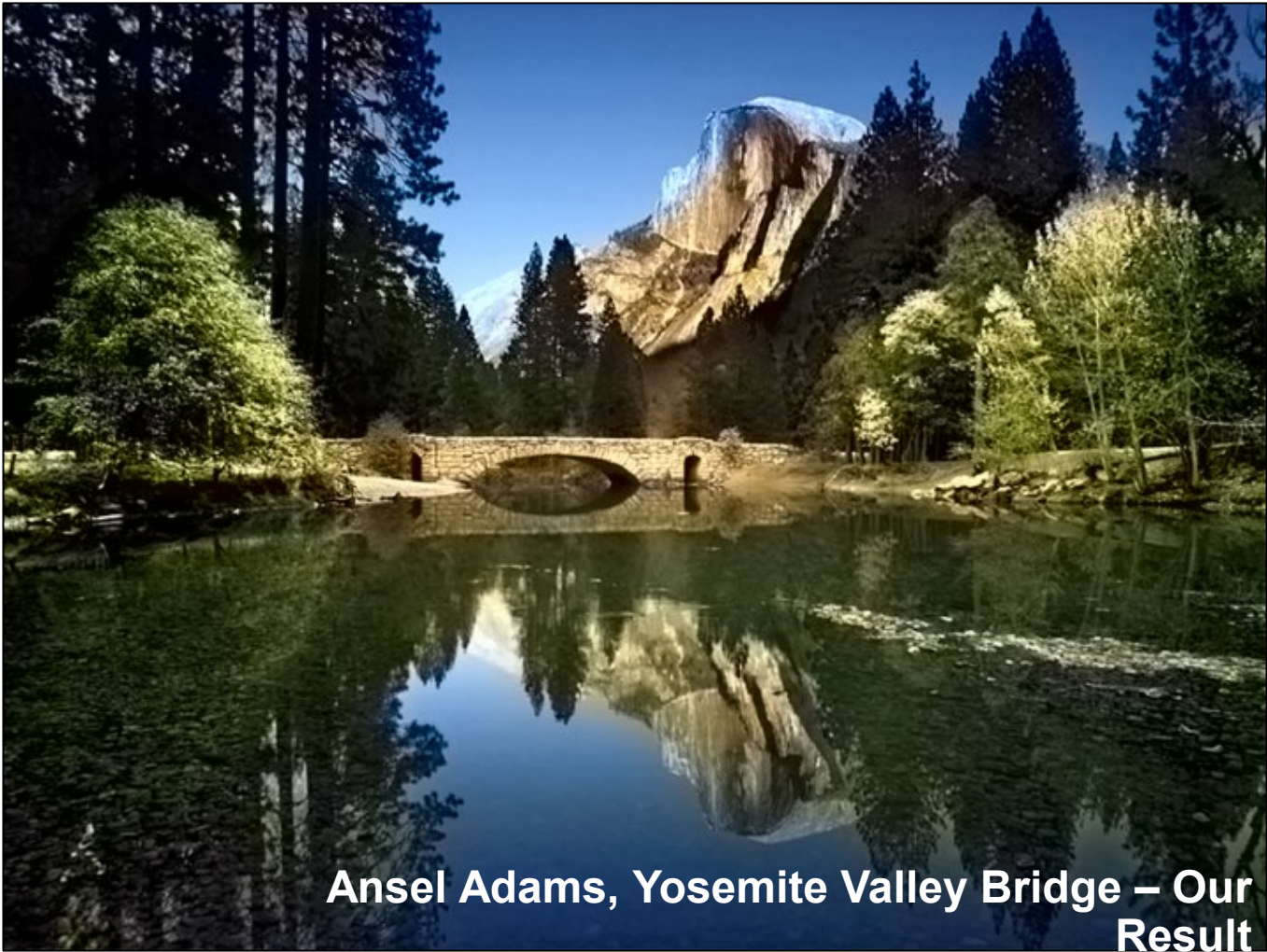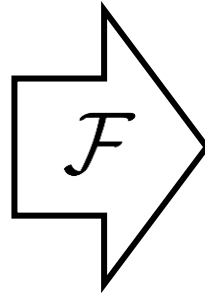Pre-train on relative-position task, w/o labels

[Girshick et al. 2014]

# Other Sources Of Signal

**Ansel Adams, Yosemite Valley Bridge**

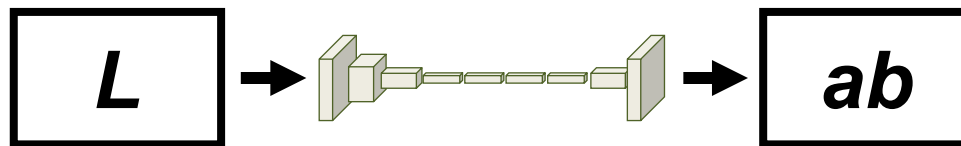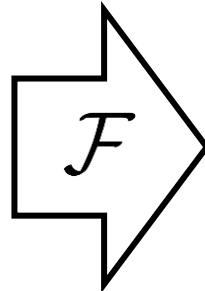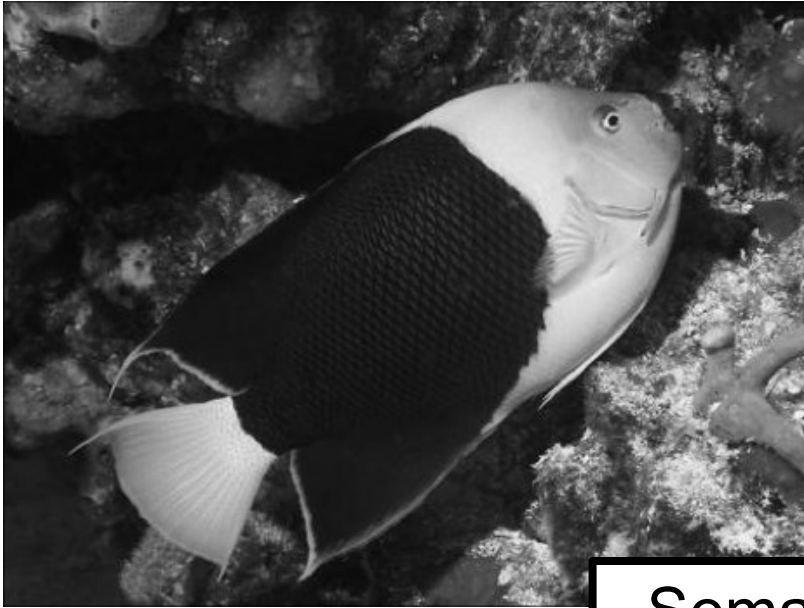Ansel Adams, Yosemite Valley Bridge – Our Result

Slide Credit: R. Zhang

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

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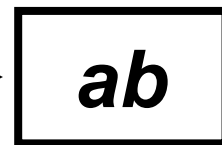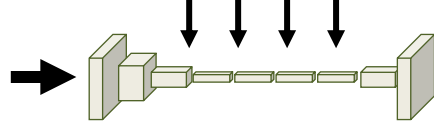$$\widehat{\mathbf{Y}} \in \mathbb{R}^{H \times W \times 2}$$

Grayscale image: L ... oncatenate (L,ab)

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

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

Semantics? Higher-level abstraction?

L

ab

"Free" supervisory signal

# Input

# Ground Truth

# Output

# Contrastive Learning



Weights shared between networks, but details important

Given sample, construct feature z; take a bunch of other images, minimize:

$$\frac{\exp(\boldsymbol{z}^T \boldsymbol{z})}{\exp(\boldsymbol{z}^T \boldsymbol{z} + \sum_i \boldsymbol{z}^T \boldsymbol{x_i})}$$

Basically, a scoring function with w = z:

$$\frac{\exp(\boldsymbol{w}^T \boldsymbol{z})}{\exp(\boldsymbol{w}^T \boldsymbol{z} + \sum_i \boldsymbol{w}^T \boldsymbol{x_i})}$$

He et al. Momentum Contrastive Learning, 2019. Wu et al. Instance discrimination 2018

# Contrastive Learning



**Best performing methods measure distance to augmented sample $z'$:**

$$\frac{\exp(z^T z')}{\exp(z^T z' + \sum_i z^T x_i)}$$

**Goal: score augmented sample higher than everything else.**

$z$

$x_i$

$z'$

He et al. Momentum Contrastive Learning, 2019. Wu et al. Instance discrimination 2018