CS 4644/7643: Lecture 18 Danfei Xu

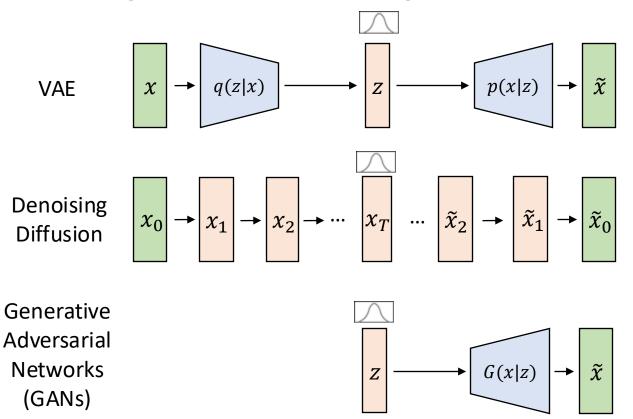
Topics:

- Generative Adversarial Networks
- Self-supervised Learning
 - Pretext task from image transformation
 - Contrastive learning

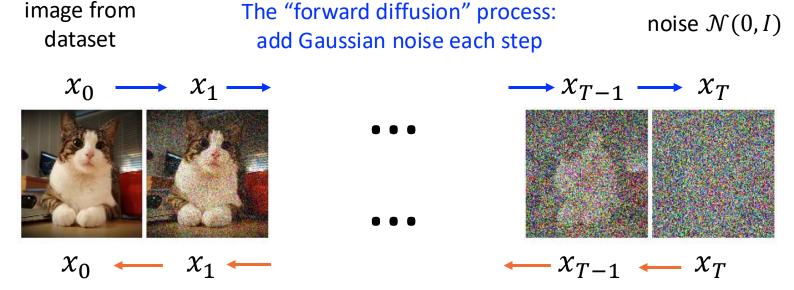
Administrative

- Start the HW4 coding part NOW --- it takes a long time to run GAN / diffusion model training on Colab GPUs.
- Milestone Report due Nov 3th. **NO GRACE PERIOD**

Denoising Diffusion: Image to Noise and Back

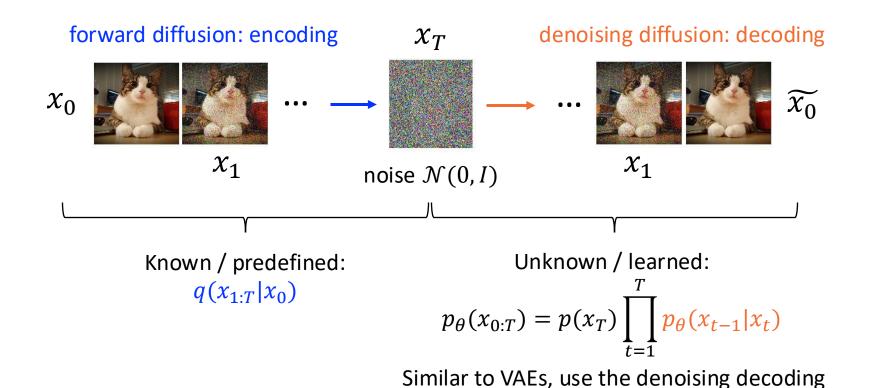


The Denoising Diffusion Process



The "denoising diffusion" process: generate an image from noise by denoising the gaussian noises

Connection to VAEs



process to generate new images.

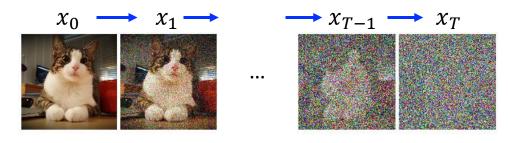
The Diffusion (Encoding) Process

The **known** forward process $x_0 \longrightarrow x_1 \longrightarrow \cdots \longrightarrow x_T$ $q(x_{1:T}|x_0) = \prod^T q(x_t|x_{t-1})$ Probability Chain Rule (Markov Chain)

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; (1-\beta_t)x_{t-1}, \beta_t I)$$
 Conditional Gaussian

 β_t is the *variance schedule* at the diffusion step t

 $0 < \beta_1 < \beta_2 < \dots < \beta_T < 1$, typical value range [0.0001, 0.02], with T = 1000



The Denoising (Decoding) Process

The **learned** denoising process
$$x_0 \leftarrow x_1 \leftarrow \cdots \leftarrow x_T$$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t) \quad \text{Probability Chain Rule (Markov Chain)}$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_q(t)) \quad \text{Conditional Gaussian}$$

$$\text{Want to learn time-dependent mean} \quad \text{Assume fixed / known variance}$$

$$\text{(simplification)}$$

How do we form a learning objective?

The Denoising (Decoding) Process

The **learned** denoising process $x_0 \leftarrow x_1 \leftarrow \cdots \leftarrow x_T$ $p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\sigma}(t))$

High-level intuition: derive a ground truth denoising distribution $q(x_{t-1}|x_t, x_0)$ and train a neural net $p_{\theta}(x_{t-1}|x_t)$ to match the distribution.

The learning objective: $\operatorname{argmin}_{\theta} D_{KL}(q(x_{t-1}|x_t,x_0)||p_{\theta}(x_{t-1}|x_t))$

What does it look like? $q(x_{t-1}|x_t,x_0) = \mathcal{N}\left(x_{t-1};\mu_q(t),\Sigma_q(t)\right)$

$$\mu_q(t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{(1 - \bar{\alpha}_t)}} \epsilon \right), \qquad \epsilon \sim \mathcal{N}(0, I) \longleftarrow \text{Recall: Gaussian reparameterization trick}$$

The "ground truth" noise that brought x_{t-1} to x_t

The Denoising (Decoding) Process

The **learned** denoising process $x_0 \longleftarrow x_1 \longleftarrow x_T$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_q(t))$$

High-level intuition: derive a ground truth denoising distribution $q(x_{t-1}|x_t, x_0)$ and train a neural net $p_{\theta}(x_{t-1}|x_t)$ to match the distribution.

The learning objective: $\operatorname{argmin}_{\theta} D_{KL}(q(x_{t-1}|x_t,x_0)||p_{\theta}(x_{t-1}|x_t))$

What does it look like?
$$q(x_{t-1}|x_t,x_0) = \mathcal{N}\left(x_{t-1};\mu_q(t),\Sigma_q(t)\right)$$

Assuming identical variance $\Sigma_q(t)$, we have:

$$\operatorname{argmin}_{\theta} D_{KL}(q(x_{t-1}|x_t,x_0)| \left| p_{\theta}(x_{t-1}|x_t) \right) = \operatorname{argmin}_{\theta} w ||\mu_q(t) - \mu_{\theta}(x_t,t)||$$

Should be variance-dependent, but constant works better in practice

$$p(x) = \int p(x|z)p(z)dz \qquad \text{Intractable to estimate!}$$

$$\log p(x) = \operatorname{E}_q \left[\log \frac{p(x|z)p(z)}{q(z|x)} \right] + D_{KL}(q(z|x)||p(z|x))$$

$$\geq \operatorname{E}_q \left[\log \frac{p(x|z)p(z)}{q(z|x)} \right] \qquad \text{Evidence Lower Bound (ELBO)}$$

$$\log p(x_0) \geq \operatorname{E}_q \left[\log \frac{p(x_0|x_{1:T})p(x_{1:T})}{q(x_{1:T}|x_0)} \right] \qquad x = x_0, \ z = x_{1:T}$$

$$= \operatorname{E}_q \left[\log \frac{p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t)}{\prod_{t=1}^T q(x_t|x_{t-1})} \right]$$

... (derivation omitted, see Sohl-Dickstein et al., 2015 Appendix B)

$$= -\mathbf{E}_{q}[D_{KL}(q(x_{T}|x_{0})||p(x_{T}))] - \sum_{t=2}^{T} D_{KL}(q(x_{t-1}|x_{t},x_{0})||p_{\theta}(x_{t-1}|x_{t})) + \log p_{\theta}(x_{0}|x_{1})$$

Maximize the agreement between the predicted reverse diffusion distribution p_{θ} and the "ground truth" reverse diffusion distribution q

Learning the Denoising Process

The **learned** denoising process
$$x_0 \leftarrow x_1 \leftarrow \cdots \leftarrow x_T$$

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$$

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma(t))$$
 Conditional Gaussian

Learning objective: $\operatorname{argmin}_{\theta} || \mu_q(t) - \mu_{\theta}(x_t, t) ||$

$$\mu_q(t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{(1 - \bar{\alpha}_t)}} \epsilon \right), \qquad \epsilon \sim \mathcal{N}(0, I)$$
 known during inference in

Idea: just learn ϵ with $\epsilon_{\theta}(x_t, t)$!

Learning the Denoising Process

The **learned** denoising process $x_0 \leftarrow x_1 \leftarrow x_T$ $p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$ $p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma(t))$ Conditional Gaussian

Simplified learning objective:
$$\operatorname{argmin}_{\theta} || \epsilon - \epsilon_{\theta} (\sqrt{\overline{\alpha}_t} x_0 + \sqrt{1 - \overline{\alpha}_t} \epsilon, t) ||$$

Inference time:
$$\mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{\beta_t}{\sqrt{(1 - \overline{\alpha}_t)}} \epsilon_{\theta}(x_t, t) \right)$$

Predicted "denoising noise"

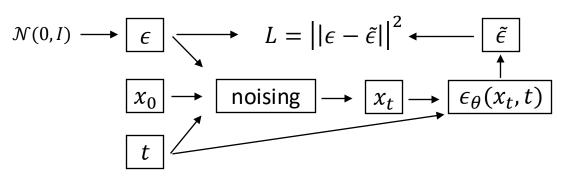
The Denoising Diffusion Algorithm

Algorithm 1 Training

- 1: repeat
- 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on

$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} (\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$

6: until converged



Compute regression loss

The Denoising Diffusion Algorithm

Algorithm 1 Training

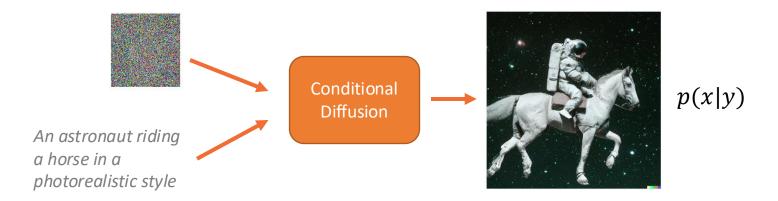
- 1: repeat
- 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3: $t \sim \text{Uniform}(\{1,\ldots,T\})$
- 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on

$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta} (\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$

6: until converged

Algorithm 2 Sampling 1: $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 2: **for** t = T, ..., 1 **do** 3: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ if t > 1, else $\mathbf{z} = \mathbf{0}$ 4: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 5: end for 6: return x_0 $p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu(t), \Sigma(t))$ $\epsilon_{\theta}(x_t, t)$

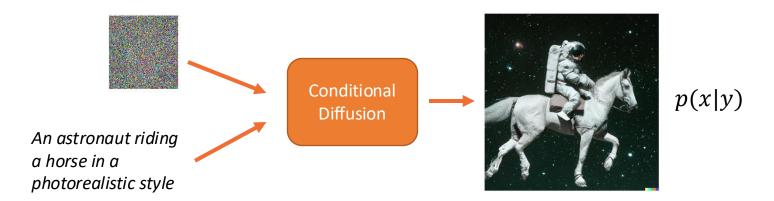
Conditional Diffusion Models



Simple idea: just condition the model on some text labels y! $\epsilon_{\theta}(x_t, y, t)$

Problem: Very blurry generation

Classifier-free Guided Diffusion



Classifier-free Guided Diffusion: estimate the gradient of the classifier model with conditional diffusion models!

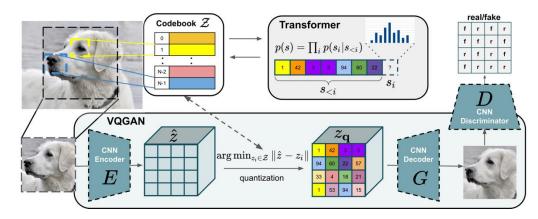
$$\nabla_{x_t} \log f_{\varphi}(y|x_t) = -\frac{1}{\sqrt{1-\bar{\alpha}_t}} (\epsilon_{\theta}(x_t, t, y) - \epsilon_{\theta}(x_t, t))$$
$$\bar{\epsilon}_{\theta}(x_t, t, y) = (w+1)\epsilon_{\theta}(x_t, t, y) - w\epsilon_{\theta}(x_t, t)$$

Linearly combine denoisers from an unconditional distribution and a conditional distribution

Latent-space Diffusion

Problem: Hard to learn diffusion process on high-resolution images

Solution: learn a low-dimensional latent space using a ViT-based autoencoder and do diffusion on the latent space!



The latent space autoencoder

Hot off arXiv!

https://www.arxiv.org/pdf/2510.21890

The Principles of Diffusion Models

From Origins to Advances

Chieh-Hsin Lai

Sony Al

Yang Song OpenAl

Dongjun Kim Stanford University

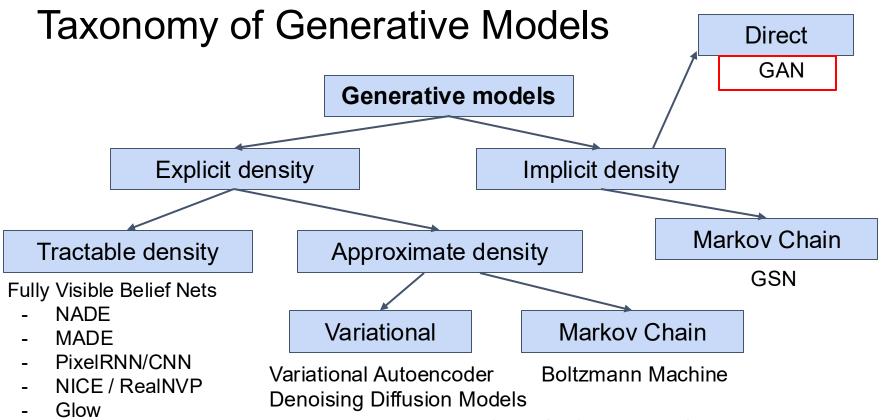
Yuki Mitsufuji

Sony Corporation, Sony Al

Stefano Ermon Stanford University

Summary

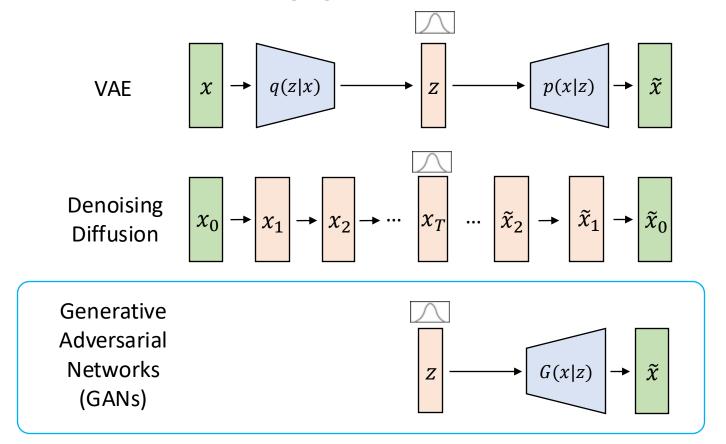
- Denoising Diffusion model is a type of generative model that learns the process of "denoising" a known noise source (Gaussian).
- We can construct a learning problem by deriving the evidence lower bound (ELBO) of the denoising process.
- The learning objective is to minimize the KL divergence between the "ground truth" and the learned denoising distribution.
- A simplified learning objective is to estimate the noise of the forward diffusion process.
- The diffusion process can be guided to generate targeted samples.
- Can be applied to many different domains. Same underlying principle.
- Very hot topic!



Ffjord

Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

GANs: Learning generate samples directly

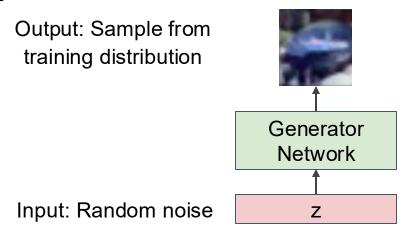


Problem: Want to sample from complex, high-dimensional training distribution. No direct way to do this!

Solution: Sample from a simple distribution we can easily sample from, e.g. random noise. Learn transformation to training distribution.

Problem: Want to sample from complex, high-dimensional training distribution. No direct way to do this!

Solution: Sample from a simple distribution we can easily sample from, e.g. random noise. Learn transformation to training distribution.



Problem: Want to sample from complex, high-dimensional training distribution. No direct way to do this!

Solution: Sample from a simple distribution we can easily sample from, e.g. random noise. Learn transformation to training distribution.

But we don't know which sample z maps to which training image -> can't learn by reconstructing training images

Output: Sample from training distribution

Generator Network

Input: Random noise

Problem: Want to sample from complex, high-dimensional training distribution. No direct way to do this!

Solution: Sample from a simple distribution we can easily sample from, e.g. random noise. Learn transformation to training distribution.

But we don't know which sample z maps to which training image -> can't learn by reconstructing training images

Output: Sample from training distribution

Objective: generated images should look "real"

Generator Network

Input: Random noise

Problem: Want to sample from complex, high-dimensional training distribution. No direct way to do this!

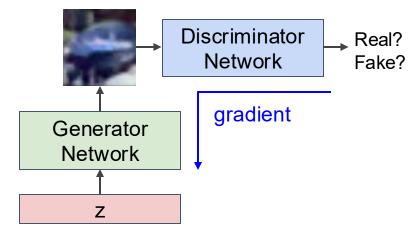
Solution: Sample from a simple distribution we can easily sample from, e.g. random noise. Learn transformation to training distribution.

But we don't know which sample z maps to which training image -> can't learn by reconstructing training images

Solution: Use a discriminator network to tell whether the generate image is within data distribution ("real") or not

Output: Sample from training distribution

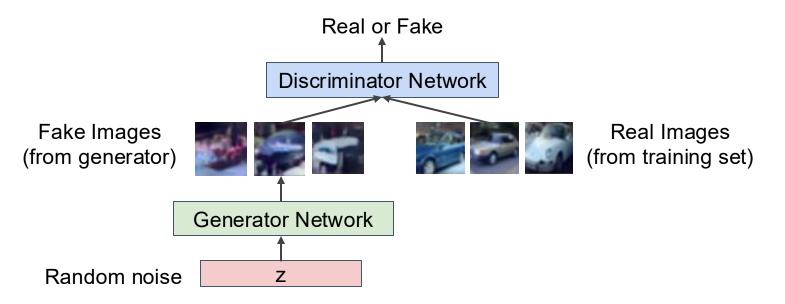
Input: Random noise



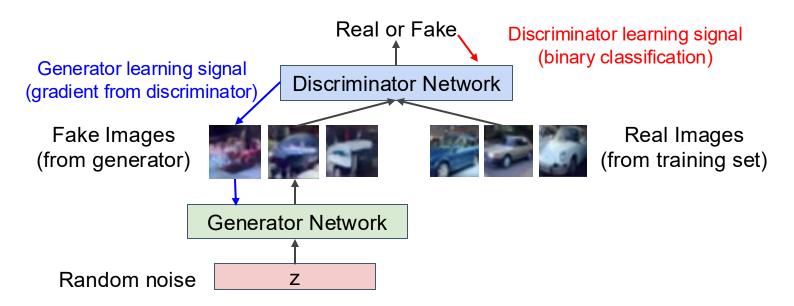
Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Discriminator network: try to distinguish between real and fake images **Generator network**: try to fool the discriminator by generating real-looking images

Discriminator network: try to distinguish between real and fake images **Generator network**: try to fool the discriminator by generating real-looking images



Discriminator network: try to distinguish between real and fake images **Generator network**: try to fool the discriminator by generating real-looking images



Fake and real images copyright Emily Denton et al. 2015. Reproduced with permission.

Discriminator network: try to distinguish between real and fake images **Generator network**: try to fool the discriminator by generating real-looking images

Train jointly in minimax game

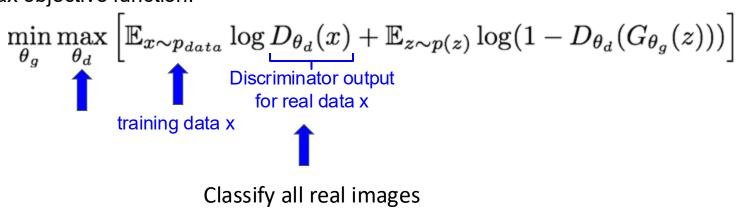
$$\min_{\substack{\theta_g \text{ Generator} \\ \text{objective}}} \max_{\substack{\theta_d \text{ Discriminator} \\ \text{objective}}} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Discriminator network: try to distinguish between real and fake images **Generator network**: try to fool the discriminator by generating real-looking images

as real

Train jointly in minimax game

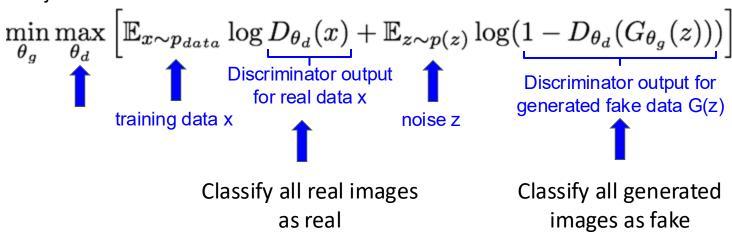
Discriminator outputs likelihood in (0,1) of real image



Discriminator network: try to distinguish between real and fake images **Generator network**: try to fool the discriminator by generating real-looking images

Train jointly in minimax game

Discriminator outputs likelihood in (0,1) of real image



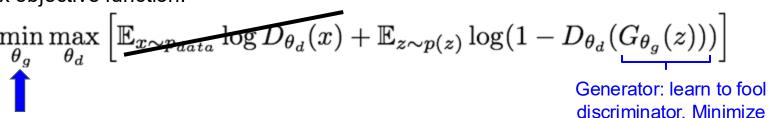
 $\log(1 - p_{\theta_d}(x_{qen}))$

Training GANs: Two-player game

Discriminator network: try to distinguish between real and fake images **Generator network**: try to fool the discriminator by generating real-looking images

Train jointly in minimax game

Discriminator outputs likelihood in (0,1) of real image



Discriminator network: try to distinguish between real and fake images **Generator network**: try to fool the discriminator by generating real-looking images

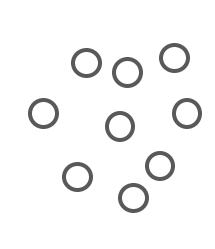
Train jointly in minimax game

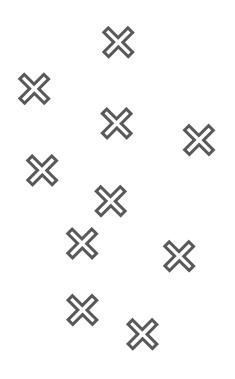
Discriminator outputs likelihood in (0,1) of real image

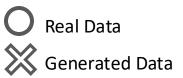
$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

- Discriminator (θ_d) wants to **maximize objective** such that D(x) is close to 1 (real) and D(G(z)) is close to 0 (fake)
- Generator (θ_g) wants to **minimize objective** such that D(G(z)) is close to 1 (discriminator is fooled into thinking generated G(z) is real)

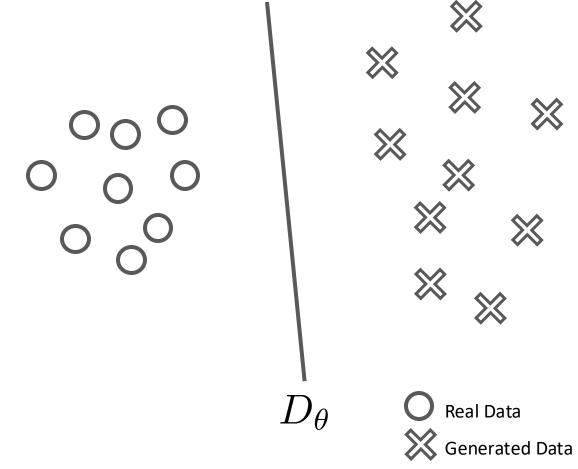
GAN Learning Process



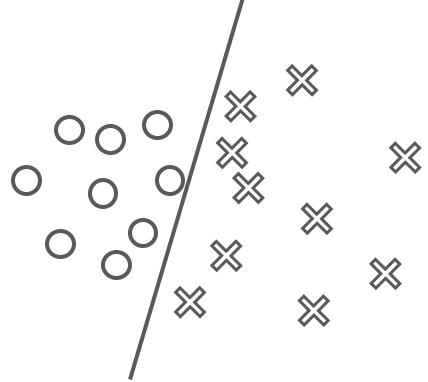




GAN Learning Process

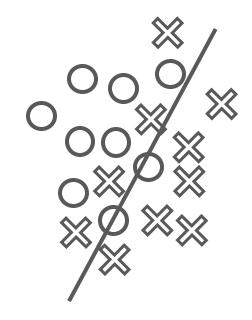


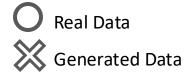
GAN Learning Process



Real Data
Generated Data

GAN Learning Process





Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Gradient descent on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient ascent on discriminator

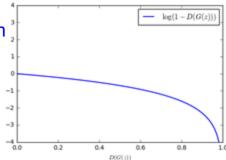
$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Gradient descent on generator

$$\min_{ heta_g} \mathbb{E}_{z \sim p(z)} \log (1 - D_{ heta_d}(G_{ heta_g}(z)))$$
 fake, want to learn from

In practice, optimizing this generator objective does not work well!

When sample is likely fake, want to learn from it to improve generator (move to the right on X axis).



Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Gradient descent on generator

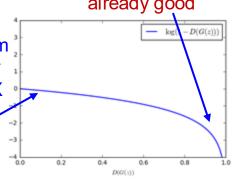
$$\min_{ heta_g} \mathbb{E}_{z \sim p(z)} \log (1 - D_{ heta_d}(G_{ heta_g}(z)))$$
 fake, want to learn from

In practice, optimizing this generator objective does not work well!

When sample is likely fake, want to learn from it to improve generator (move to the right on X axis).

But gradient in this region is relatively flat!

Gradient signal dominated by region where sample is already good



Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

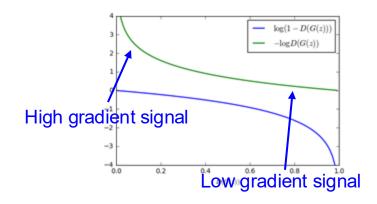
1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log (1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Instead: Gradient ascent on generator, different objective $\max_{\theta_{c}} \mathbb{E}_{z \sim p(z)} \log(D_{\theta_{d}}(G_{\theta_{g}}(z)))$

Instead of minimizing likelihood of discriminator being correct, now maximize likelihood of discriminator being wrong.

Same objective of fooling discriminator, but now higher gradient signal for bad samples => works much better! Standard in practice.



Putting it together: GAN training algorithm

for number of training iterations do

for k steps do

- Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- · Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D_{\theta_d}(x^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(z^{(i)}))) \right]$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_q(z)$.
- Update the generator by ascending its stochastic gradient (improved objective):

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d}(G_{\theta_g}(z^{(i)})))$$

end for

Putting it together: GAN training algorithm

for number of training iterations do

for k steps do

- Sample minibatch of m noise samples $\{z^{(1)}, \ldots, z^{(m)}\}$ from noise prior $p_q(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \ldots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- · Update the discriminator by ascending its stochastic gradient:

 $\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D_{\theta_d}(x^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(z^{(i)}))) \right]$

Update scriminator

Followup work (e.g. Wasserstein GAN, BEGAN) alleviates this problem, better stability!

Some find k=1

others use k > 1,

more stable,

no best rule.

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_q(z)$.
- Update the generator by ascending its stochastic gradient (improved objective):

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d}(G_{\theta_g}(z^{(i)})))$$

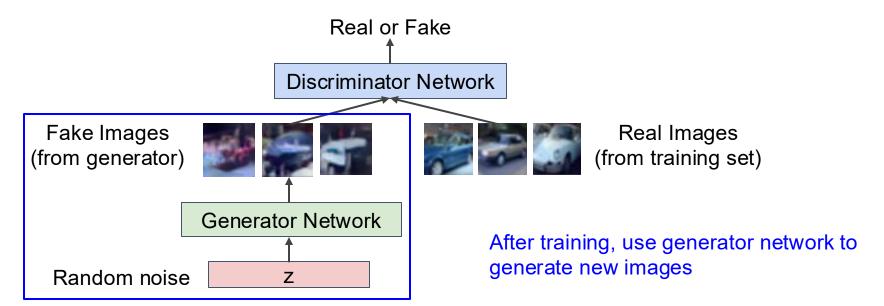
Update generator

end for

Arjovsky et al. "Wasserstein gan." arXiv preprint arXiv:1701.07875 (2017)

Berthelot, et al. "Began: Boundary equilibrium generative adversarial networks." arXiv preprint arXiv:1703.10717 (2017)

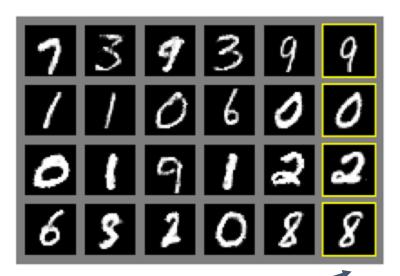
Generator network: try to fool the discriminator by generating real-looking images **Discriminator network**: try to distinguish between real and fake images



Fake and real images copyright Emily Denton et al. 2015. Reproduced with permission.

Generative Adversarial Nets

Generated samples





Nearest neighbor from training set

Generative Adversarial Nets

Generated samples (CIFAR-10)





Nearest neighbor from training set

Generative Adversarial Nets: Convolutional Architectures

Generator is an upsampling network with fractionally-strided convolutions Discriminator is a convolutional network

Architecture guidelines for stable Deep Convolutional GANs

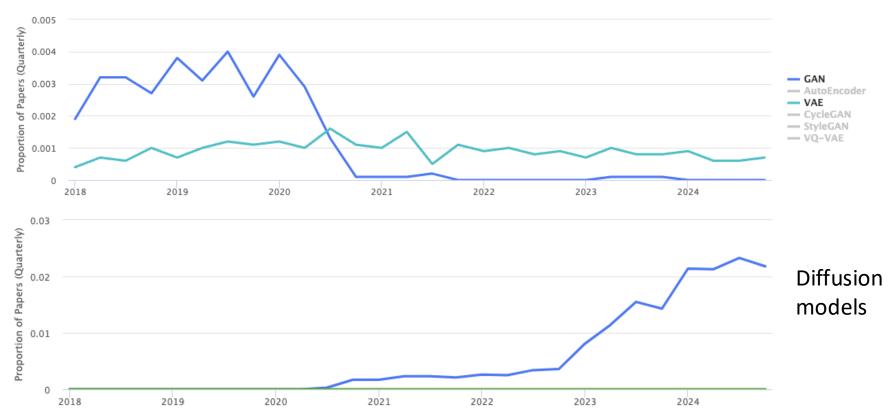
- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except for the output, which uses Tanh.
- Use LeakyReLU activation in the discriminator for all layers.

Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

2019: BigGAN

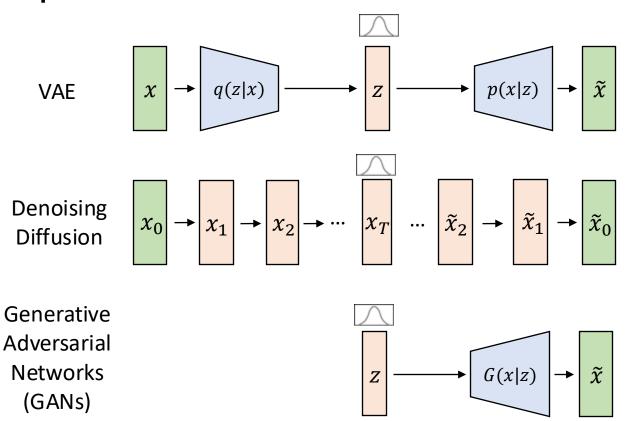


GANs were popular ...



Source: https://paperswithcode.com

Deep Generative Models



Generative Models: Closing Thoughts

- Learn without supervision = ability to leverage large, raw dataset
- Realism: Generate plausible samples given dataset
- Diversity: Generate diverse samples (avoid mode collapse)
- Controllability: Generate based on instruction / conditioning
- Healthy combination of theory and deep learning magic
- Generative Modeling is extremely hot in 2025. More will come ...

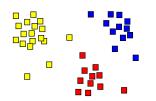
Supervised Learning

- Train Input: {*X*, *Y*}
- Learning output: $f: X \to Y, P(y|x)$
- e.g. classification

Sheep Dog Cat Lion Giraffe

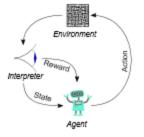
Unsupervised Learning

- Input: {X}
- Learning output: P(x)
- Example: Clustering, density estimation, generative modeling



Reinforcement Learning

- Evaluative feedback in the form of reward
- No supervision on the right action



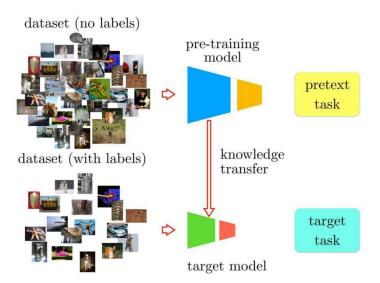
Self-Supervised Learning:

Create your own supervision

Self-supervised Learning

In short: still supervised learning, with two important distinctions:

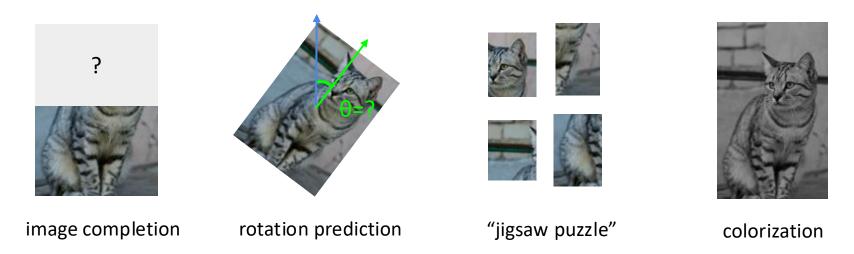
- 1. Learn from labels generated *autonomously* instead of human annotations.
- 2. The goal is to learn good representations for other target tasks.



Source: Noroozi et al., 2018

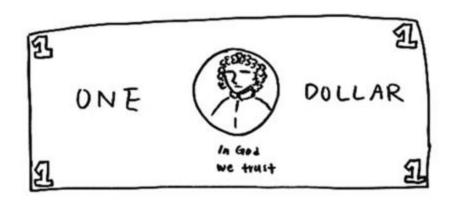
Self-supervised pretext tasks

Example: learn to predict image transformations / complete corrupted images



- 1. Solving the pretext tasks allow the model to learn good features.
- 2. We can automatically generate labels for the pretext tasks.

Generative vs. Self-supervised Learning





Left: Drawing of a dollar bill from memory. Right: Drawing subsequently made with a dollar bill present. Image source: Epstein, 2016

Learning to generate pixel-level details is often unnecessary; learn abstract features with pretext tasks instead

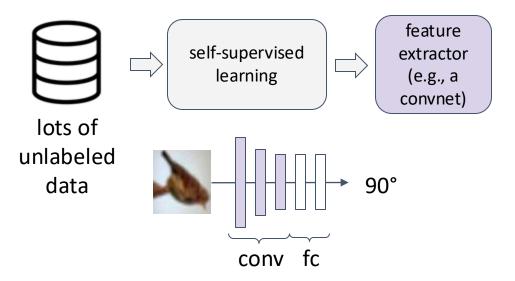
Source: Anand, 2020

How to evaluate a self-supervised learning method?

We usually don't care about the performance of the self-supervised learning task, e.g., we don't care if the model learns to predict image rotation perfectly.

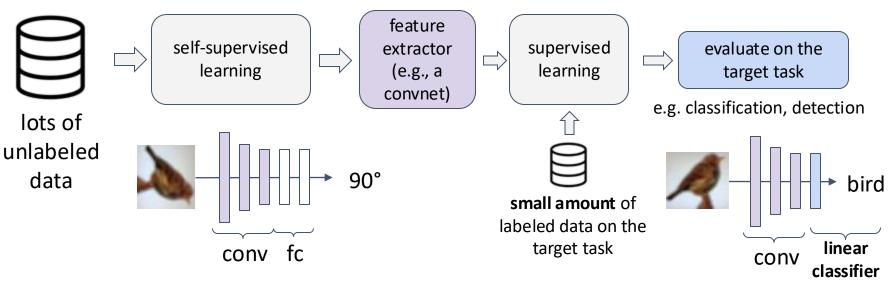
Evaluate the learned feature encoders on downstream target tasks

How to evaluate a self-supervised learning method?



1. Learn good feature extractors from self-supervised pretext tasks, e.g., predicting image rotations

How to evaluate a self-supervised learning method?



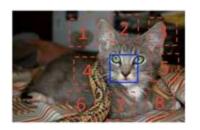
1. Learn good feature extractors from self-supervised pretext tasks, e.g., predicting image rotations

2. Attach a shallow network on the feature extractor; train the shallow network on the target task with small amount of labeled data

Broader picture

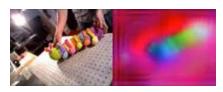
Today's lecture

computer vision



Doersch et al., 2015

robot / reinforcement learning



Dense Object Net (Florence and Manuelli et al., 2018)

language modeling

Language Models are Few-Shot Learners

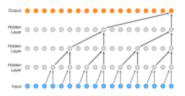
Tom B. Be	rows"	Benjamin	Mann"	Nick I	Ryder*	Melanie Subbiah*	
Jared Kaplan [†]	Profeile	Dharival	Arvind Noth	Austus	Pranar Sh	ram Girish Sastry	
Amundu Askell	Sandhini	Agarwal	Ariel Herbert	Vons	Greichen Kru	eger Tom Henighan	
Revon Child	Aditya	Ramesh	Daniel M. Zi	ingler	Jeffrey Wu	Clemens Winter	
Christopher I	Book	Mark Chen	Eric Sig	for	Mateusz Litw	in Souti Gray	
Benjamin Chess			Jack Clark		Christopher Berner		
Sum McCandlish		Aloc Ra	althred	Bya Sutskever		Durio Amodei	
			OpenAl				

Abstract

Recent work has demonstrated substantial gains on many NLP tasks and benchmarks by pre-train on a large corpus of text followed by fine-tuning on a specific task. While typically task-agnostic in architecture, this method still requires task-specific fine-tuning datasets of thousands or tens of thousands of examples. By contrast, humans can generally perform a new language task from only a few examples or from simple instructions - something which current NLP systems still largely struggle to do. Here we show that scaling up language models greatly improves task-agnostic few-shot performance, sometimes even maching competitiveness with prior state-of-the-art fine tuning approaches. Specifically, we train GPT-3, an autoregressive language model with 175 billion parameters, 10s more than any previous non-sparse language model, and test its performance in the few-shot setting. For all tasks, GPT-3 is applied without any gradient updates or fine-tuning with tasks and few-shot demonstrations specified purely via test interaction with the model. GPTachieves strong performance on many NLF datasets, including translation, question-unswering, and close tasks, as well as several tasks that require on-the-fly reasoning or domain adaptation, such as unscrambling words, using a novel word in a sentence, or performing 3-digit arithmetic. At the same time, we also identify some datasets where GPT-3's few shot learning still struggles, as well as some datasets where GPT-3 faces methodological issues related to training on large web-corpora. Finally we find that GPT-5 can generate samples of news articles which human evaluators have difficulty distinguishing from articles written by humans. We discuss broader societal impacts of this finding

GPT3 (Brown, Mann, Ryder, Subbiah et al., 2020)

speech synthesis



Wavenet (van den Oord et al., 2016)

• • •

Today's Agenda

Pretext tasks from image transformations

- Rotation, inpainting, rearrangement, coloring

Contrastive representation learning

- Intuition and formulation
- Instance contrastive learning: SimCLR and MOCO
- Sequence contrastive learning: CPC

Today's Agenda

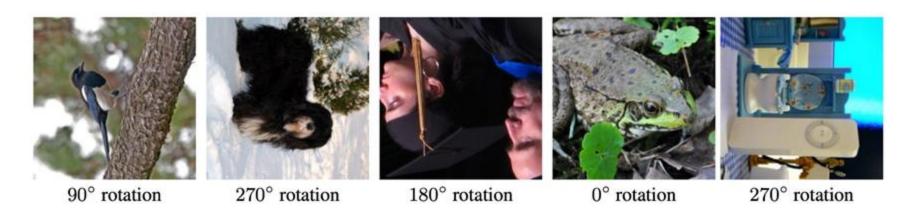
Pretext tasks from image transformations

- Rotation, inpainting, rearrangement, coloring

Contrastive representation learning

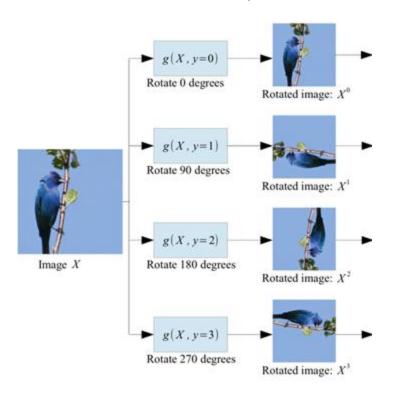
- Intuition and formulation
- Instance contrastive learning: SimCLR and MOCO
- Sequence contrastive learning: CPC

Pretext task: predict rotations



Hypothesis: a model could recognize the correct rotation of an object only if it has the "visual commonsense" of what the object should look like unperturbed.

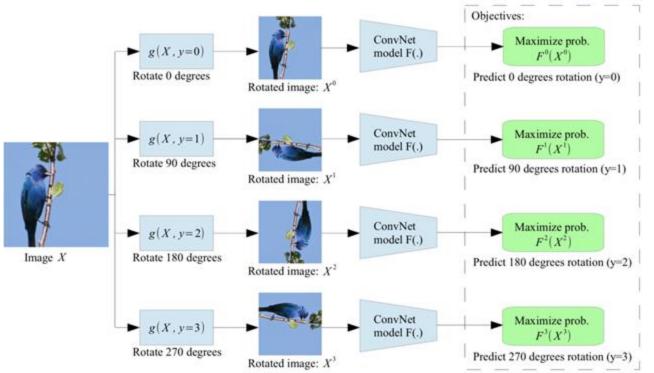
Pretext task: predict rotations



Self-supervised learning by rotating the entire input images.

The model learns to predict which rotation is applied (4-way classification)

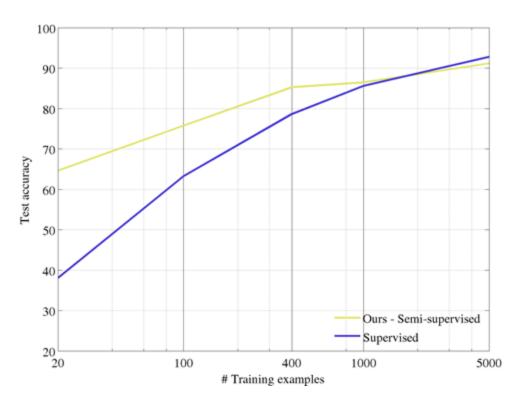
Pretext task: predict rotations



Self-supervised learning by rotating the entire input images.

The model learns to predict which rotation is applied (4-way classification)

Evaluation on semi-supervised learning



Self-supervised learning on **CIFAR10** (entire training set).

Freeze conv1 + conv2 Learn conv3 + linear layers with subset of labeled CIFAR10 data (classification).

Transfer learned features to supervised learning

		fication nAP)	Detection (%mAP)	Segmentation (%mIoU)
Trained layers	fc6-8	all	all	all
ImageNet labels	78.9	79.9	56.8	48.0
Random		53.3	43.4	19.8
Random rescaled Krähenbühl et al. (2015)	39.2	56.6	45.6	32.6
Egomotion (Agrawal et al., 2015)	31.0	54.2	43.9	
Context Encoders (Pathak et al., 2016b)	34.6	56.5	44.5	29.7
Tracking (Wang & Gupta, 2015)	55.6	63.1	47.4	
Context (Doersch et al., 2015)	55.1	65.3	51.1	
Colorization (Zhang et al., 2016a)	61.5	65.6	46.9	35.6
BIGAN (Donahue et al., 2016)	52.3	60.1	46.9	34.9
Jigsaw Puzzles (Noroozi & Favaro, 2016)	-	67.6	53.2	37.6
NAT (Bojanowski & Joulin, 2017)	56.7	65.3	49.4	
Split-Brain (Zhang et al., 2016b)	63.0	67.1	46.7	36.0
ColorProxy (Larsson et al., 2017)		65.9		38.4
Counting (Noroozi et al., 2017)	-	67.7	51.4	36.6
(Ours) RotNet	70.87	72.97	54.4	39.1

Pretrained with fullImageNet supervisionNo pretraining

Self-supervised learning on ImageNet (entire training set) with AlexNet.

Finetune on labeled data from **Pascal VOC 2007**.

Self-supervised learning with rotation prediction

source: Gidaris et al. 2018

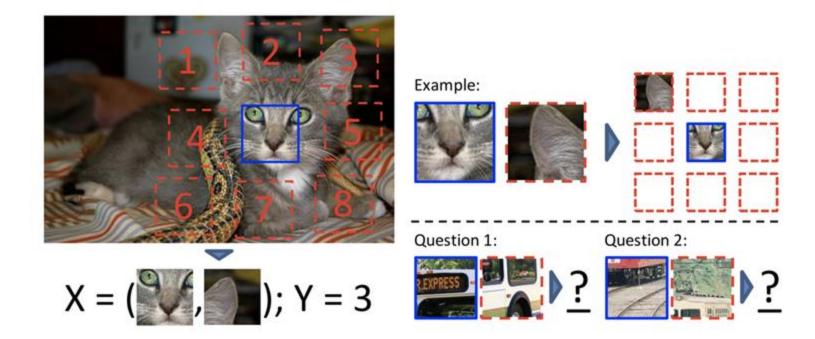
Visualize learned visual attentions



(a) Attention maps of supervised model

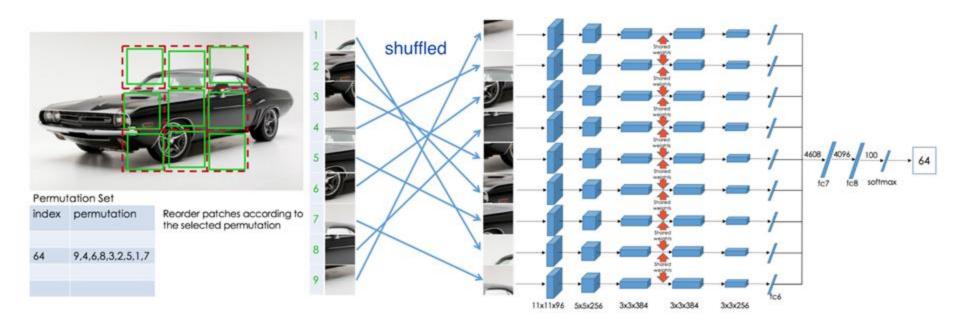
(b) Attention maps of our self-supervised model

Pretext task: predict relative patch locations



(Image source: <u>Doersch et al., 2015</u>)

Pretext task: solving "jigsaw puzzles"



(Image source: Noroozi & Favaro, 2016)

Transfer learned features to supervised learning

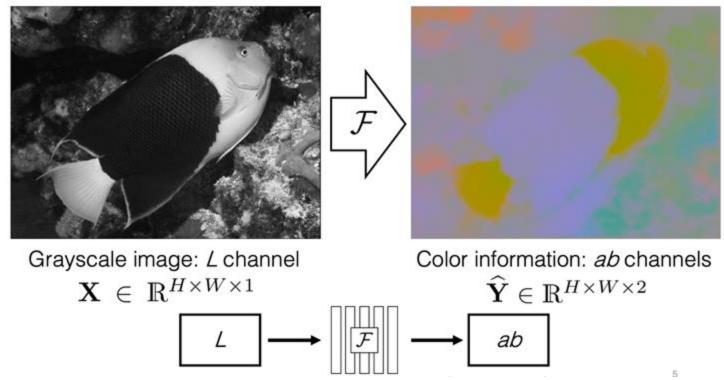
Table 1: Results on PASCAL VOC 2007 Detection and Classification. The results of the other methods are taken from Pathak et al. [30].

Method	Pretraining time	Supervision	Classification	Detection	Segmentation
Krizhevsky <i>et al.</i> [25]	3 days	1000 class labels	78.2%	56.8%	48.0%
Wang and Gupta[39]	1 week	motion	58.4%	44.0%	-
Doersch et al. [10]	4 weeks	context	55.3%	46.6%	-
Pathak et al. [30]	14 hours	context	56.5%	44.5%	29.7%
Ours	$2.5 \mathrm{days}$	context	$\boldsymbol{67.6\%}$	$\boldsymbol{53.2\%}$	37.6 %

"Ours" is feature learned from solving image Jigsaw puzzles (Noroozi & Favaro, 2016). Doersch et al. is the method with relative patch location

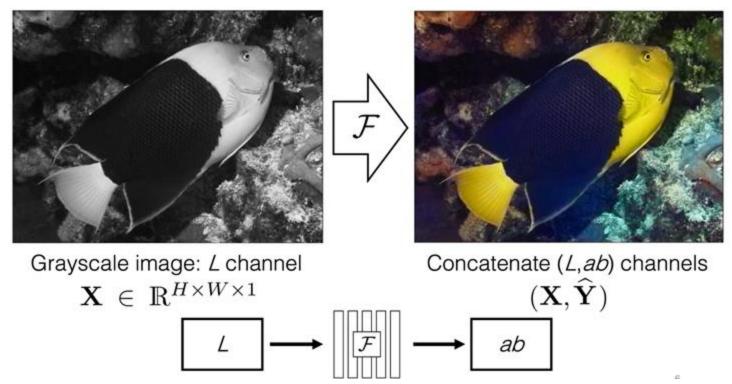
(source: Noroozi & Favaro, 2016)

Pretext task: image coloring



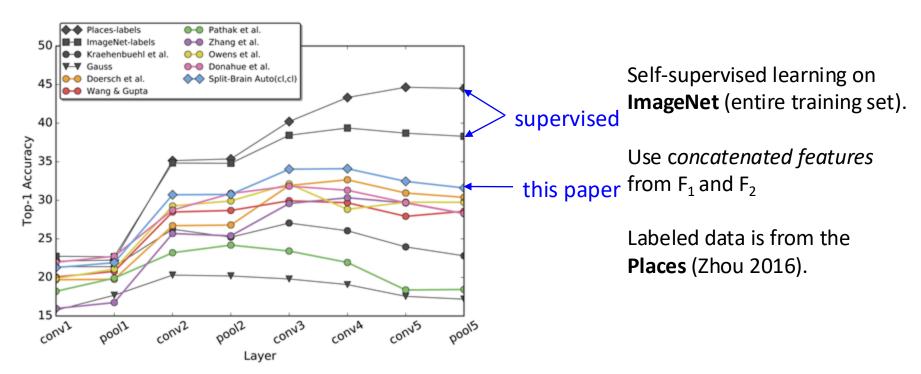
Source: Richard Zhang / Phillip Isola

Pretext task: image coloring



Source: Richard Zhang / Phillip Isola

Transfer learned features to supervised learning



Source: Zhang et al., 2017

Pretext task: video coloring

Idea: model the *temporal coherence* of colors in videos

reference frame



t = 0

how should I color these frames?







t = 2



t = 3

Source: Vondrick et al.,

<u> 2018</u>

Pretext task: video coloring

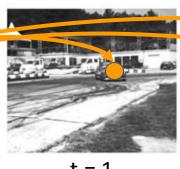
Idea: model the *temporal coherence* of colors in videos

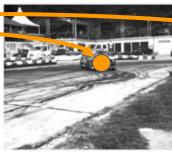
reference frame

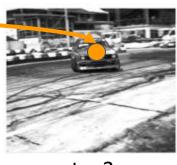
how should I color these frames?

Should be the same color!









• • •

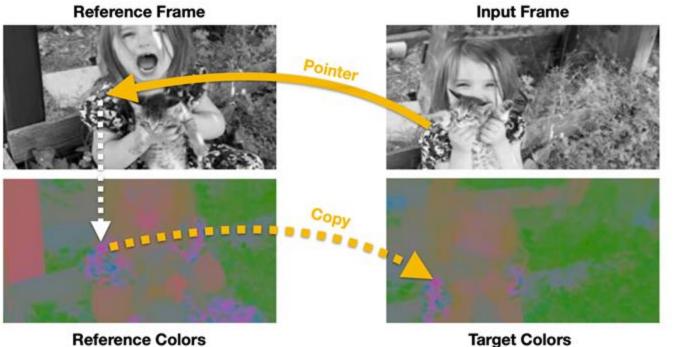
t = 0

t=1 t=2 t=3

Hypothesis: learning to color video frames should allow model to learn to track regions or objects without labels!

Source: Vondrick et al.,

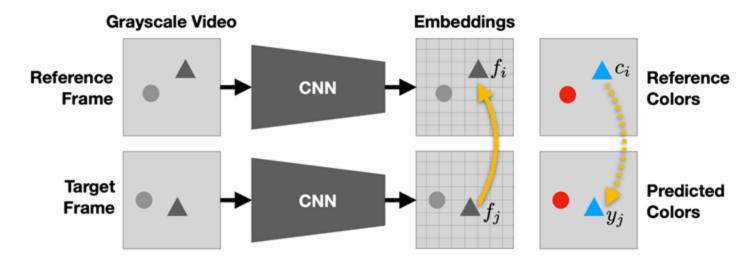
2018



Learning objective:

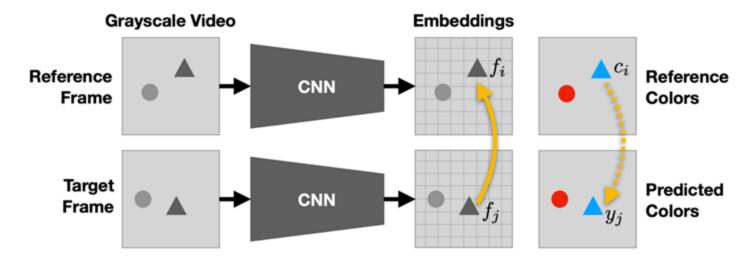
Establish mappings between reference and target frames in a learned feature space.

Use the mapping as "pointers" to copy the correct color (LAB).



attention map on the reference frame

$$A_{ij} = rac{\exp\left(f_i^T f_j\right)}{\sum_k \exp\left(f_k^T f_j\right)}$$

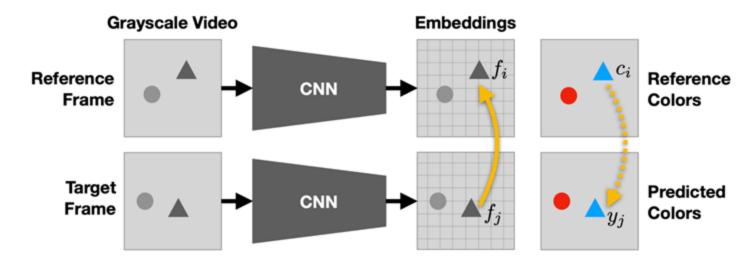


attention map on the reference frame

$$A_{ij} = \frac{\exp(f_i^T f_j)}{\sum_k \exp(f_k^T f_j)}$$

predicted color = weighted sum of the reference color

$$y_j = \sum_i A_{ij} c_i$$



attention map on the reference frame

$$A_{ij} = \frac{\exp(f_i^T f_j)}{\sum_k \exp(f_k^T f_j)}$$

predicted color = weighted sum of the reference color

$$y_j = \sum_i A_{ij} c_i$$

loss between predicted color and ground truth color

$$\min_{\theta} \sum_{j} \mathcal{L}\left(y_{j}, c_{j}\right)$$

Colorizing videos (qualitative)

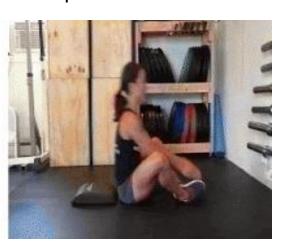
reference frame



target frames (gray)



predicted color



Source: Google Al blog

Colorizing videos (qualitative)

reference frame



target frames (gray)



predicted color



Source: Google Al blog

Tracking emerges from colorization

Propagate segmentation masks using learned attention



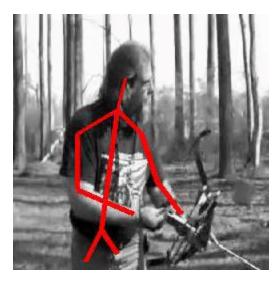




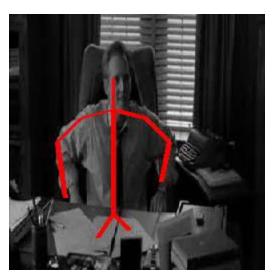
Source: Google Al blog

Tracking emerges from colorization

Propagate pose keypoints using learned attention







Source: Google AI blog

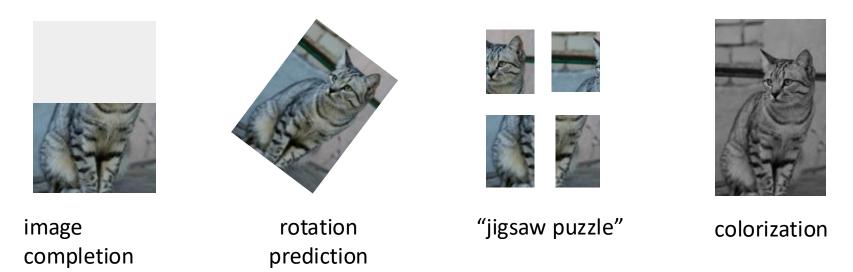
Summary: pretext tasks from image transformations

- Pretext tasks focus on "visual common sense", e.g., predict rotations, inpainting, rearrangement, and colorization.
- The models are forced learn good features about natural images, e.g., semantic representation of an object category, in order to solve the pretext tasks.
- We don't care about the performance of these pretext tasks, but rather how useful the learned features are for downstream tasks (classification, detection, segmentation).

Summary: pretext tasks from image transformations

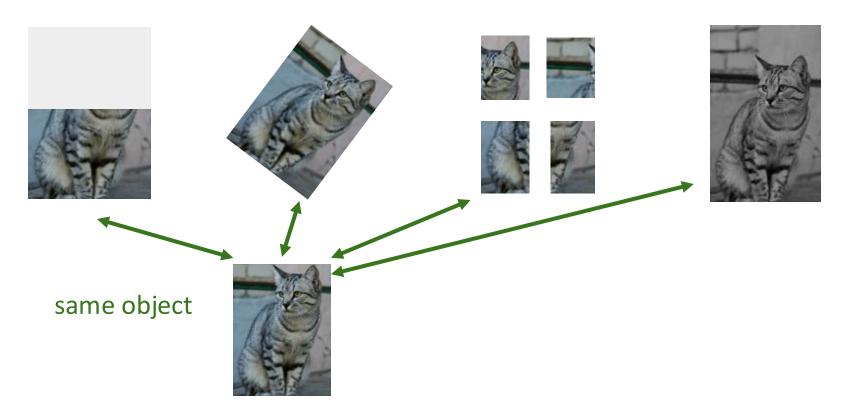
- Pretext tasks focus on "visual common sense", e.g., predict rotations, inpainting, rearrangement, and colorization.
- The models are forced learn good features about natural images, e.g., semantic representation of an object category, in order to solve the pretext tasks.
- We don't care about the performance of these pretext tasks, but rather how useful the learned features are for downstream tasks (classification, detection, segmentation).
- Problems: 1) coming up with individual pretext tasks is tedious, and 2) the learned representations may not be general.

Pretext tasks from image transformations

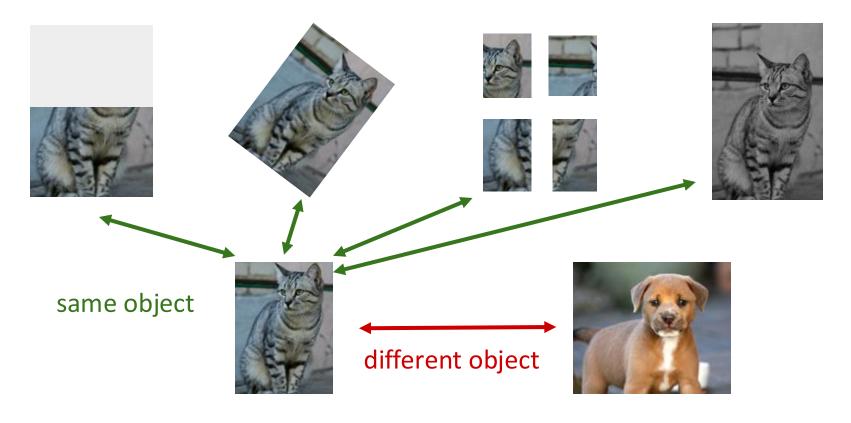


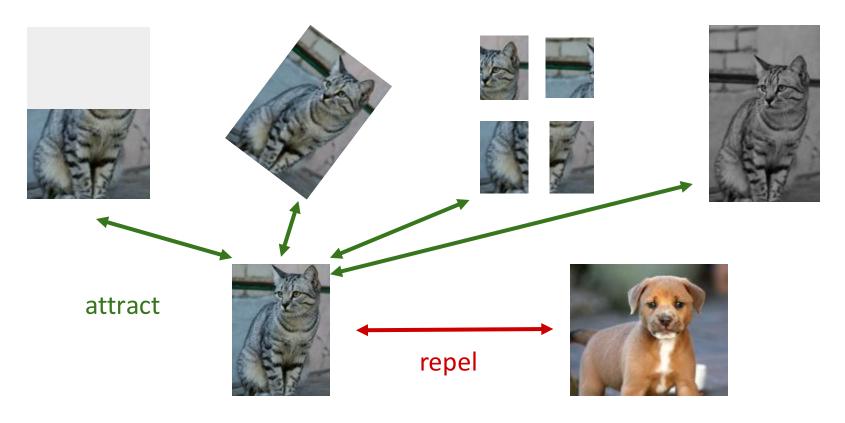
Learned representations may be tied to a specific pretext task! Can we come up with a more general pretext task?

A more general pretext task?



A more general pretext task?





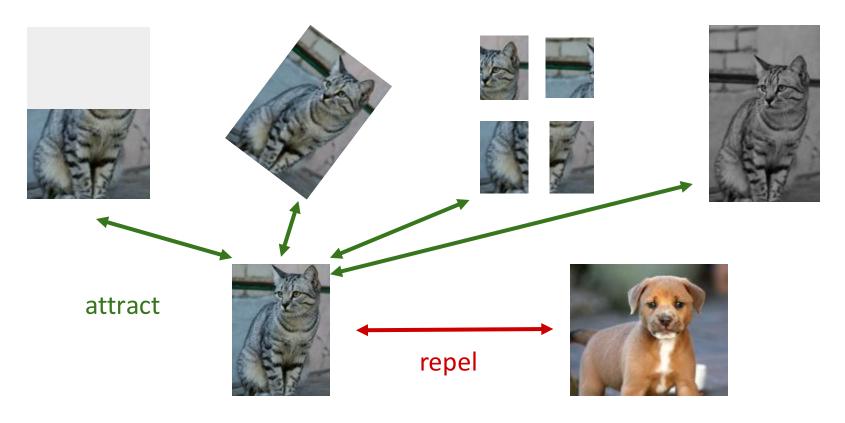
Today's Agenda

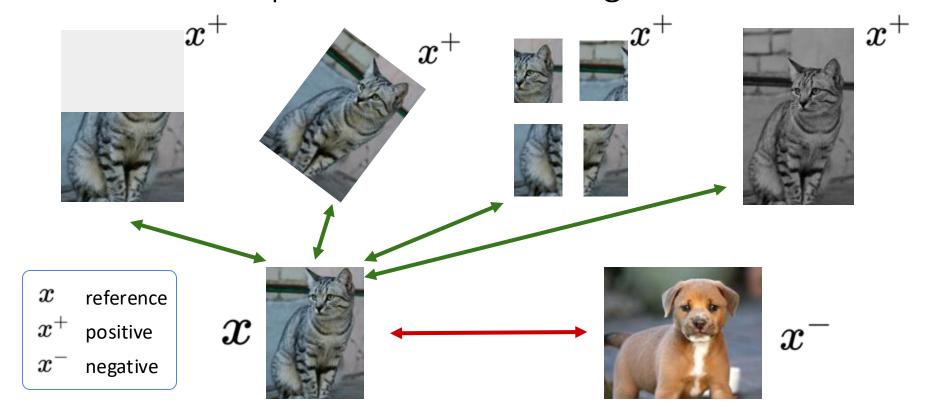
Pretext tasks from image transformations

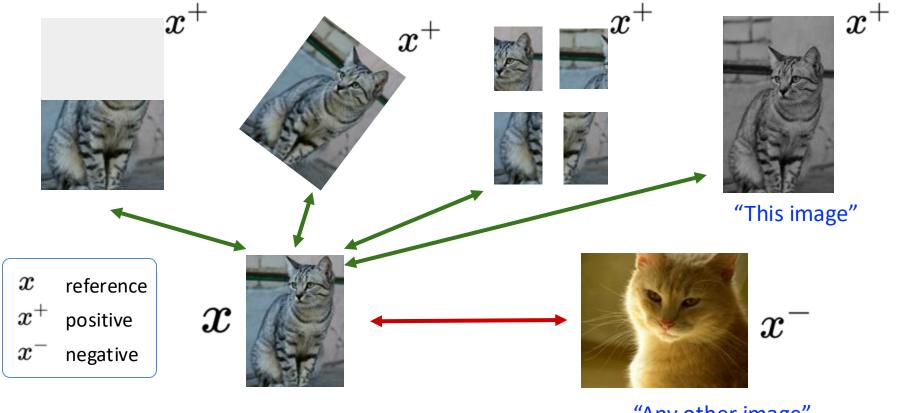
- Rotation, inpainting, rearrangement, coloring

Contrastive representation learning

- Intuition and formulation
- Instance contrastive learning: SimCLR and MOCO
- Sequence contrastive learning: CPC







"Any other image"

What we want:

$$score(f(x), f(x^+)) >> score(f(x), f(x^-))$$

x: reference sample; x⁺ positive sample; x⁻ negative sample

Given a chosen score function, we aim to learn an **encoder function** f that yields high score for positive pairs (x, x^+) and low scores for negative pairs (x, x^-) .

Loss function given 1 positive sample and N - 1 negative samples:

$$L = -\mathbb{E}_X \left[\log \frac{\exp(s(f(x), f(x^+)))}{\exp(s(f(x), f(x^+)) + \sum_{j=1}^{N-1} \exp(s(f(x), f(x_j^-))))} \right]$$

Loss function given 1 positive sample and N - 1 negative samples:

$$L = -\mathbb{E}_X \left[\log \frac{ \exp(s(f(x), f(x^+))}{ \exp(s(f(x), f(x^+)) + \sum_{j=1}^{N-1} \exp(s(f(x), f(x_j^-)))} \right]$$

$$x \quad x^+ \quad x \quad x^- \quad x_3^- \quad x_3^-$$

. . .

Loss function given 1 positive sample and N - 1 negative samples:

$$L = -\mathbb{E}_X \left[\log \frac{\exp(s(f(x), f(x^+))}{\exp(s(f(x), f(x^+)) + \sum_{j=1}^{N-1} \exp(s(f(x), f(x_j^-)))} \right]$$
 score for the positive score for the N-1 negative pair

This seems familiar ...

Loss function given 1 positive sample and N - 1 negative samples:

$$L = -\mathbb{E}_X \left[\log \frac{\exp(s(f(x), f(x^+))}{\exp(s(f(x), f(x^+)) + \sum_{j=1}^{N-1} \exp(s(f(x), f(x_j^-)))} \right]$$
 score for the positive score for the N-1 negative pair

This seems familiar ...

Cross entropy loss for a N-way softmax classifier!
I.e., learn to find the positive sample from the N samples

Loss function given 1 positive sample and N - 1 negative samples:

$$L = -\mathbb{E}_X \left[\log \frac{\exp(s(f(x), f(x^+)))}{\exp(s(f(x), f(x^+)) + \sum_{j=1}^{N-1} \exp(s(f(x), f(x_j^-))))} \right]$$

Commonly known as the InfoNCE loss (van den Oord et al., 2018)

A lower bound on the mutual information between f(x) and $f(x^{+})$

$$MI[f(x),f(x^+)] - \log(N) \geq -L$$

The larger the negative sample size (N), the tighter the bound

Detailed derivation: Poole et al., 2019

SimCLR: A Simple Framework for Contrastive Learning

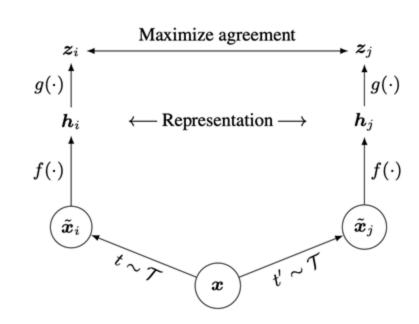
Cosine similarity as the score function:

$$s(u,v)=rac{u^Tv}{||u||||v||}$$

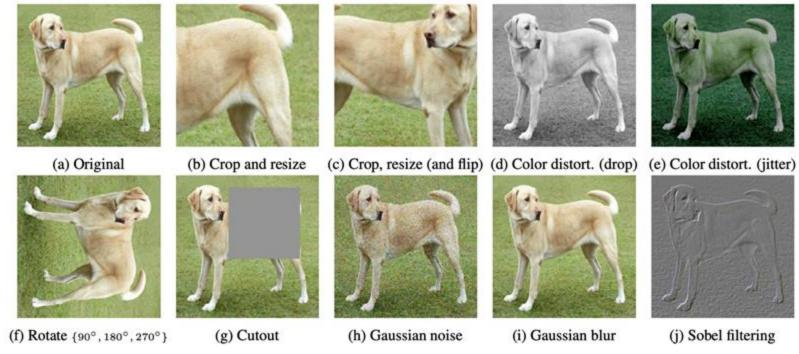
Use a projection network *h(·)* to project features to a space where contrastive learning is applied

Generate positive samples through data augmentation:

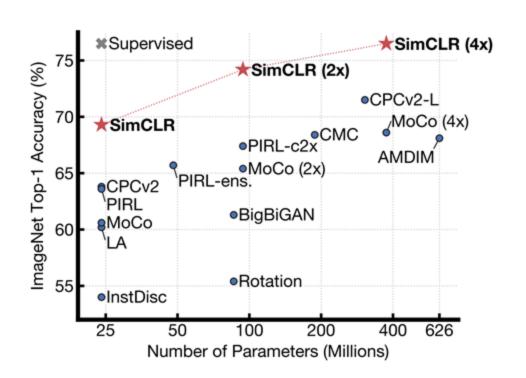
 random cropping, random color distortion, and random blur.



SimCLR: generating positive samples from data augmentation



Training linear classifier on SimCLR features



Train feature encoder on **ImageNet** (entire training set) using SimCLR.

Freeze feature encoder, train a linear classifier on top with labeled data.

Semi-supervised learning on SimCLR features

		Label	Label fraction			
Method	Architecture	1%	10%			
		Top 5				
Supervised baseline	ResNet-50	48.4	80.4			
Methods using other label-propagation:						
Pseudo-label	ResNet-50	51.6	82.4			
VAT+Entropy Min.	ResNet-50	47.0	83.4			
UDA (w. RandAug)	ResNet-50	-	88.5			
FixMatch (w. RandAug)	ResNet-50	-	89.1			
S4L (Rot+VAT+En. M.)	ResNet-50 (4 \times)	-	91.2			
Methods using representation learning only:						
InstDisc	ResNet-50	39.2	77.4			
BigBiGAN	RevNet-50 $(4\times)$	55.2	78.8			
PIRL	ResNet-50	57.2	83.8			
CPC v2	ResNet-161(*)	77.9	91.2			
SimCLR (ours)	ResNet-50	75.5	87.8			
SimCLR (ours)	ResNet-50 (2 \times)	83.0	91.2			
SimCLR (ours)	ResNet-50 $(4\times)$	85.8	92.6			

Table 7. ImageNet accuracy of models trained with few labels.

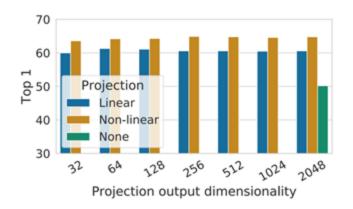
Train feature encoder on **ImageNet** (entire training set) using SimCLR.

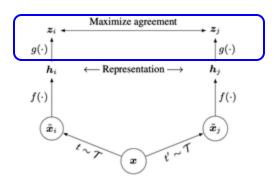
Finetune the encoder with 1% / 10% of labeled data on ImageNet.

Source: Chen et al.,

2020

SimCLR design choices: projection head





Linear / non-linear projection heads improve representation learning.

A possible explanation:

- contrastive learning objective may discard useful information for downstream tasks
- representation space z is trained to be invariant to data transformation.
- by leveraging the projection head g(·), more information can be preserved in the h representation space

SimCLR design choices: large batch size

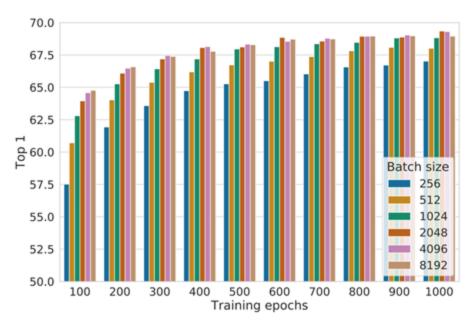


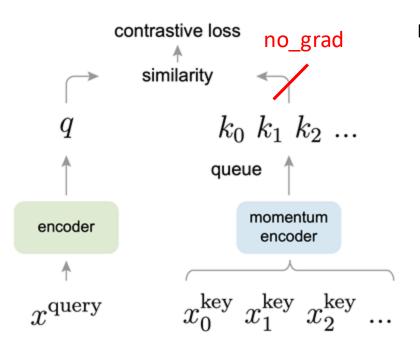
Figure 9. Linear evaluation models (ResNet-50) trained with different batch size and epochs. Each bar is a single run from scratch. 10

Large training batch size is crucial for SimCLR!

Large batch size causes large memory footprint during backpropagation: requires distributed training on TPUs (ImageNet experiments)

Source: Chen et al.,

Momentum Contrastive Learning (MoCo)

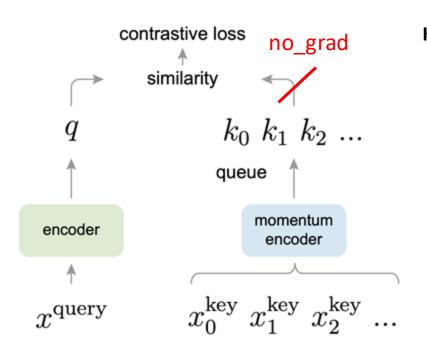


Key differences to SimCLR:

- Keep a running queue of keys (negative samples).
- Compute gradients and update the encoder only through the queries.
- Decouple min-batch size with the number of keys: can support a large number of negative samples.

Source: <u>He et al., 2020</u>

Momentum Contrastive Learning (MoCo)



Key differences to SimCLR:

- Keep a running queue of keys (negative samples).
- Compute gradients and update the encoder only through the queries.
- Decouple min-batch size with the number of keys: can support a large number of negative samples.
- The key encoder is slowly progressing through the momentum update rules:

$$\theta_{\mathbf{k}} \leftarrow m\theta_{\mathbf{k}} + (1-m)\theta_{\mathbf{q}}$$

Source: He et al., 2020

"MoCo V2"

Improved Baselines with Momentum Contrastive Learning

Xinlei Chen Haoqi Fan Ross Girshick Kaiming He Facebook AI Research (FAIR)

A hybrid of ideas from SimCLR and MoCo:

- From SimCLR: non-linear projection head and strong data augmentation.
- From MoCo: momentum-updated queues that allow training on a large number of negative samples (no TPU required!).

Source: Chen et al.,

MoCo vs. SimCLR vs. MoCo V2

	unsup. pre-train			ImageNet	VOC detection			
case	MLP	aug+	cos	epochs	acc.	AP ₅₀	AP	AP ₇₅
supervised					76.5	81.3	53.5	58.8
MoCo v1				200	60.6	81.5	55.9	62.6
(a)	✓			200	66.2	82.0	56.4	62.6
(b)		✓		200	63.4	82.2	56.8	63.2
(c)	✓	✓		200	67.3	82.5	57.2	63.9
(d)	✓	✓	✓	200	67.5	82.4	57.0	63.6
(e)	✓	\checkmark	✓	800	71.1	82.5	57.4	64.0

Table 1. **Ablation of MoCo baselines**, evaluated by ResNet-50 for (i) ImageNet linear classification, and (ii) fine-tuning VOC object detection (mean of 5 trials). "**MLP**": with an MLP head; "**aug+**": with extra blur augmentation; "**cos**": cosine learning rate schedule.

Key takeaways:

 Non-linear projection head and strong data augmentation are crucial for contrastive learning.

Source: Chen et al.,

2020

MoCo vs. SimCLR vs. MoCo V2

	unsup. pre-train					ImageNet
case	MLP	aug+	cos	epochs	batch	acc.
MoCo v1 [6]				200	256	60.6
SimCLR [2]	✓	✓	✓	200	256	61.9
SimCLR [2]	✓	✓	✓	200	8192	66.6
MoCo v2	✓	✓	✓	200	256	67.5
results of longer unsupervised training follow:						
SimCLR [2]	✓	✓	✓	1000	4096	69.3
MoCo v2	✓	✓	✓	800	256	71.1

Table 2. MoCo vs. SimCLR: ImageNet linear classifier accuracy (ResNet-50, 1-crop 224×224), trained on features from unsupervised pre-training. "aug+" in SimCLR includes blur and stronger color distortion. SimCLR ablations are from Fig. 9 in [2] (we thank the authors for providing the numerical results).

Key takeaways:

- Non-linear projection head and strong data augmentation are crucial for contrastive learning.
- Decoupling mini-batch size with negative sample size allows MoCo-V2 to outperform SimCLR with smaller batch size (256 vs. 8192).

Source: Chen et al.,

2020

MoCo vs. SimCLR vs. MoCo V2

mechanism	batch	memory / GPU	time / 200-ep.
MoCo	256	5.0G	53 hrs
end-to-end	256	7.4G	65 hrs
end-to-end	4096	93.0G [†]	n/a

Table 3. **Memory and time cost** in 8 V100 16G GPUs, implemented in PyTorch. †: based on our estimation.

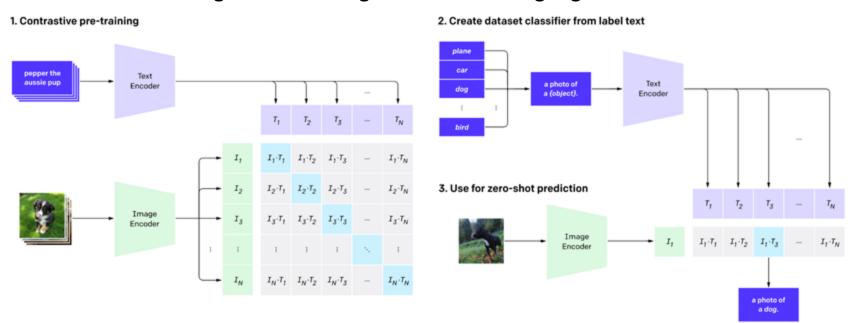
Key takeaways:

- Non-linear projection head and strong data augmentation are crucial for contrastive learning.
- Decoupling mini-batch size with negative sample size allows MoCo-V2 to outperform SimCLR with smaller batch size (256 vs. 8192).
- ... all with much smaller memory footprint! ("end-to-end" means SimCLR here)

Source: Chen et al.,

Other examples

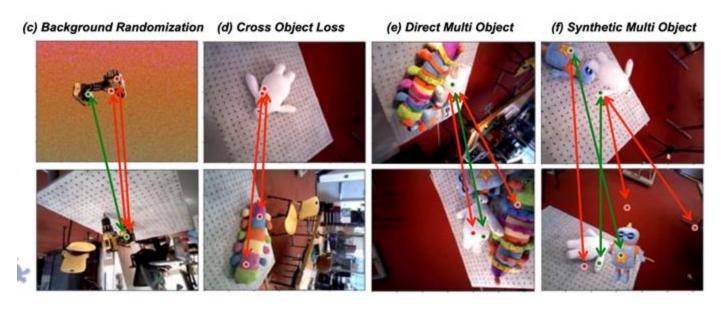
Contrastive learning between image and natural language sentences



CLIP (Contrastive Language-Image Pre-training) Radford et al., 2021

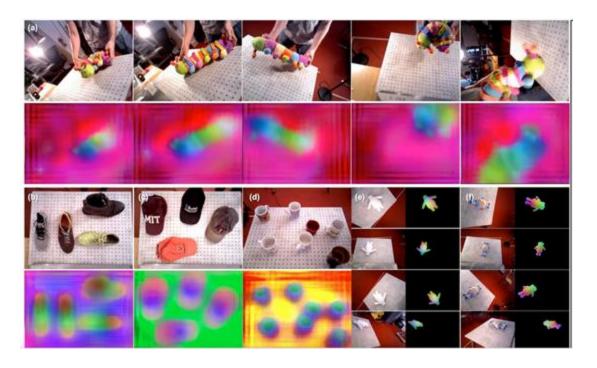
Other examples

Contrastive learning on pixel-wise feature descriptors



Dense Object Net, Florence et al., 2018

Other examples



Dense Object Net, Florence et al., 2018