

01 - Introduction

1. Network : $f = W X$

2. Loss function: Softmax / Hinge a class
Optimal \Rightarrow zero centered
saturates / well enough

3. Activation function: Sigmoid / tanh / ReLU / Maxout
saturated / saturated / ideal ReLU continue the momentum

4. Optimization: GD / SGD / GD with Momentum / RMSprop / Adam
learning rate velocity second momentum
mini batch

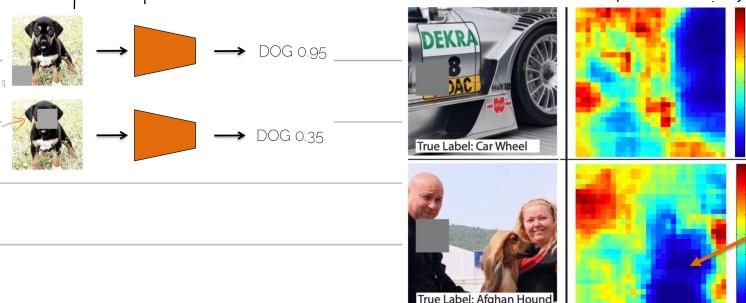
5. Training: Learning rate α / Over- and underfitting / Split data

6. Regularization: Low Validation error / High training error
Data augmentation / Early stopping / Bagging and Embedding / Dropout

7. Image task: Image filters / CNN / Pooling
LeNet / AlexNet / VGGNet / ResNet / Inception layer / GoogleNet

8. Training trick: Setup visualization / Overfit a sample \rightarrow 0 loss, 100% correct
Iteration and training time check
Small network \rightarrow debug \rightarrow bigger Network / Sample
 \hookrightarrow change only one time

9. Visualization:
① Image space
pick a unit \rightarrow Find maximize unit activation of image patches
② Occlusion experiment
Block different part in image \rightarrow check network catch the importance
 \Rightarrow Map / point-wise / classification probability



③ t-SNE All image Representation

Motivation: visualize the last FC layer

Nearest neighbor visualization based on L_2

\Rightarrow Map high-dim embedding into 2D Map

02 - Siamese and Similarity Learning

~~Si~~

1. ML \Rightarrow Classification, Regression,

Similarity Learning

Comparison

Ranking

com'parision

new ✓

Problem of Similarity Learning

Scalability:

retrain model for different people

2. Similarity Function: $d(A, B) > \tau \Rightarrow \text{not same} / d(A, B) < \tau \Rightarrow \text{same}$

3. Siamese NN: Siamese Network share weights $\Rightarrow f(A)$ and $f(B) \dots f(N)$
 Siamese encoding of image

4. Contrastive Loss: $d(A, B) = \|f(A) - f(B)\|^2 = L(A, B) \Leftrightarrow \text{positive}$

comparison $\max(0, m^2 - \|f(A) - f(B)\|^2) = L(A, B) \Leftrightarrow \text{negative}$

already far away \rightarrow not pull them together

positive samples $\rightarrow L(A, B) = y^* \|f(A) - f(B)\|^2 + (1-y^*) \max(0, m^2 - \|f(A) - f(B)\|^2) \Leftrightarrow$ CE
 to same points positive pair negative pair 2 input 1 input

$$\|f(A) - f(B)\|^2 \leq \|f(A) - f(N)\|^2$$

5 Triplet Loss: $L(A, P, N) = \max(0, \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + m)$

Ranking
 Encourage margin between cluster \rightarrow Hard negative mining training with hard case
 Train a few epochs \Rightarrow hard case $d(A, P) \approx d(A, N) \Rightarrow$ refine training

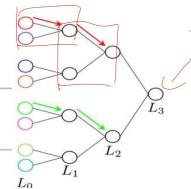
6 Improving Similarity learning long time / local minimum

- Loss:
 - Contrastive vs. triplet loss
- Sampling:
 - Choosing the best triplets to train with, sample the space wisely = diversity of classes + hard cases
- Ensembles:
 - Why not using several networks, each of them trained with a subset of triplets?
- Can we use a classification loss for similarity learning?

Sampling : Hierarchical Tree

leave \Leftrightarrow image class

recursively merge \Rightarrow root

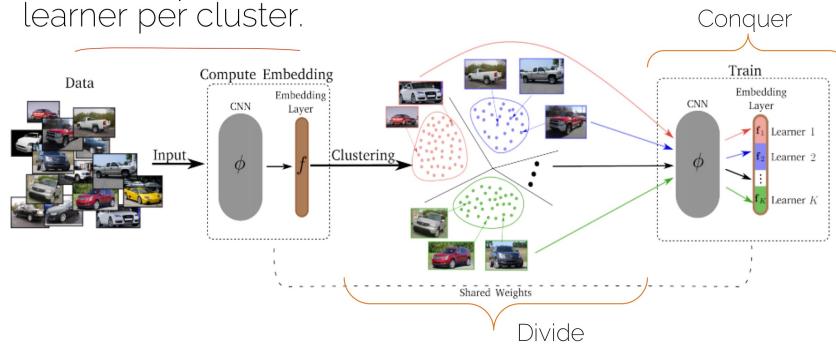


$$\text{distance: } d(p, q) = \frac{1}{n_p n_q} \sum_{i \in p, j \in q} \|r_i - r_j\|$$

- ① Randomly select l' nodes at 0 level learn discriminative
- ② select $m-1$ nearest classes based on d features from similar classes
- ③ t images per class are randomly collected

$$\mathcal{L}_M = \frac{1}{2Z_M} \sum_{T^z \in \mathcal{T}^M} [\|\mathbf{x}_a^z - \mathbf{x}_p^z\| - \|\mathbf{x}_a^z - \mathbf{x}_n^z\| + \alpha_z]_+$$

Ensemble: divide space into k cluster, each class set a learner learner per cluster.



Classification loss : For similarity learning

Tips and tricks

7 Application in Vision : Siamese on MNIST

Image correspondence : Object recognition / ...

3D correspondence

Image retrieval

Self-supervised Learning : Learn from video

Automatically generate label

Optical flow

\hookrightarrow FlowNet with CNN

\hookrightarrow Siamese Network / Correlation layer

\downarrow
Find image correspondence

03 - AE and VAE

1. ML : Supervised Learning : Label / Find a mapping / Classification, Regression

Unsupervised Learning : No label / Find data structure / Clustering

2. AE : Learning a lower-dimensional feature representation from unlabeled data

Reconstruction Loss L_1, L_2

Latent space $z \Rightarrow \dim(z) < \dim(x)$

3. AE for pre-training : Large unlabeled data / Small labeled data

AE learn the features present

Unsupervised train AE \rightarrow Supervised train encoder with label

4. AE for pixel-wise prediction FCN \rightarrow SegNet (convolution + layer / Upsampling + convolution)

Transposed convolution : Unpooling + Convolution

UpSampling : Interpolation / Interpolation + Convs / DeconvNet
 ↓
 Maxpooling

UNet : Skip connections = pass the low level informations

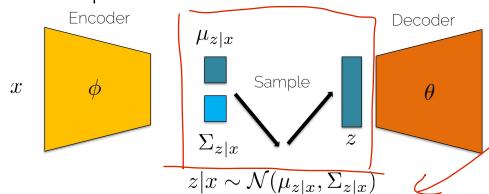
Convolution = pass the high level informations

5. AE in Vision : ① Semantic segmentation (SegNet)

② Depth estimation (Stereo camera / Monocular depth)

③ Image Super Resolution (Low Resolution \rightarrow high / Learning residual)

6. VAE : Sample from the latent distribution to generate new output



Latent space is a Gaussian Distribution

Loss function :

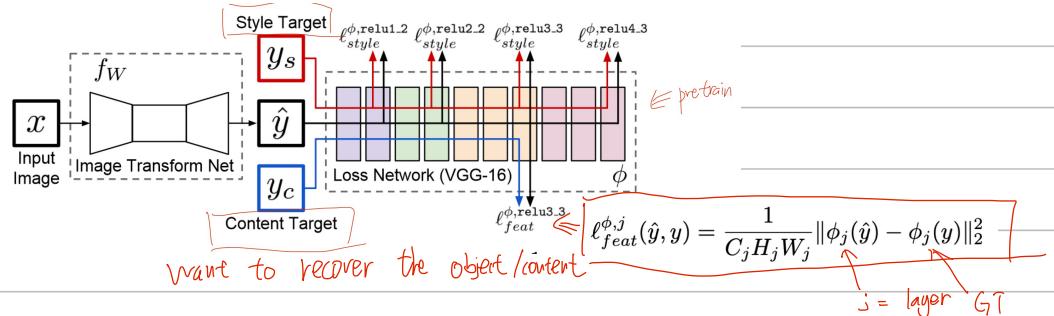
$$\begin{aligned} &= E_z [\log p_\theta(x_i|z)] - E_z \left[\log \frac{q_\phi(z|x_i)}{p_\theta(z)} \right] + E_z \left[\log \frac{q_\phi(z|x_i)}{p_\theta(z|x_i)} \right] \\ &= E_z [\log p_\theta(x_i|z)] - KL(q_\phi(z|x_i)||p_\theta(z)) + KL(q_\phi(z|x_i)||p_\theta(z|x_i)) \end{aligned}$$

Tower bound $KL \geq 0$
 $= L(x_i, \phi, \theta)$

Optimize $\phi^*, \theta^* = \arg \max \sum_{i=1}^n L(x_i, \phi, \theta)$

7. Image synthesis : Semantic Segmentation image \rightarrow Real image

~~perceptual loss~~ : L_2 loss will penalize realistic result
 Content loss measure the content of the image
 Feature representation similarity



Style transfer: Content image + Style image = Style Transfer

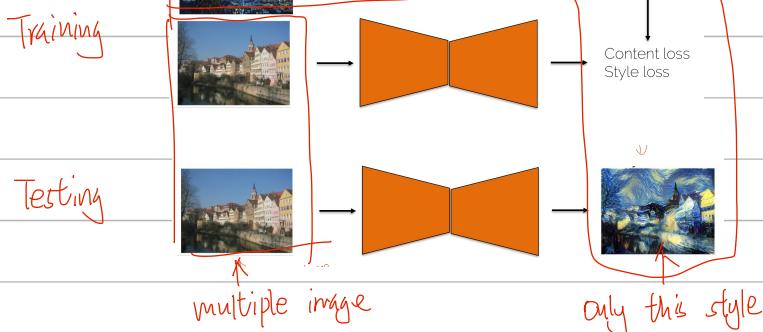
$$\ell_{style}^{\phi,j}(\hat{y}, y) = \|G_j^\phi(\hat{y}) - G_j^\phi(y)\|_F^2$$

Gram Matrix

$$G_j^\phi(x)_{c,c'} = \frac{1}{C_j H_j W_j} \sum_{h=1}^{H_j} \sum_{w=1}^{W_j} \phi_j(x)_{h,w,c} \phi_j(x)_{h,w,c'}$$

Capture the information about which features tend to activate together
 Loss preserve the stylistic features but not content

Fast transfer:



04 - Representation Learning

encode priors about data distributions

1. Representation Learning : Transform the raw data into a representation

- Smoothness: close inputs map to close outputs
- Compactness: input dimension >> output dimension
- Robustness: features are insensitive to input noise
- Abstraction and invariances -> problem driven

Convert the observation in the real world → mathematical form

→ Feature vector ⇒ Classification / Reconstruction / Generation
 ↳ Handcrafted attribute / Binary / Embedding vector

Supervised approach : Res Net + FC + Classifier

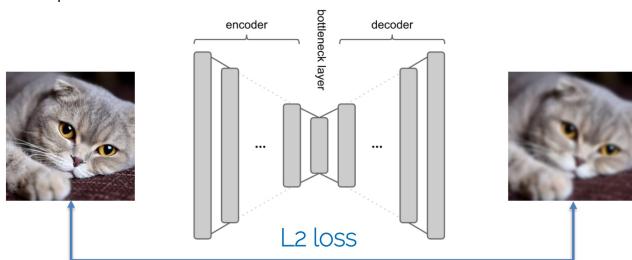
↑
Feature vector ↑
cat

Unsupervised approach : Clustering (k-means) ← mean vector

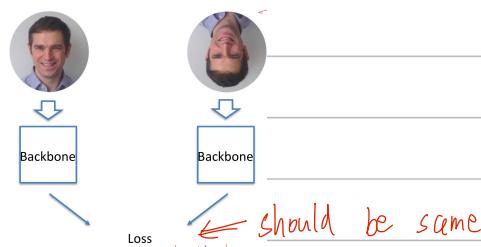
Self-supervised approach : Data provide supervision

Loss force NN to learn features

Expensive to obtain annotations / Most data are unlabeled



Self-supervised by Augmentation :



2. DINO : Self-distillation with no labels

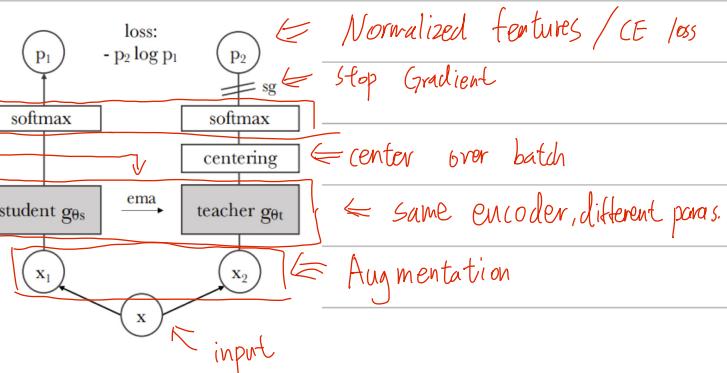
Teacher paras. are updated with an exponential moving average of student

temperature softmax

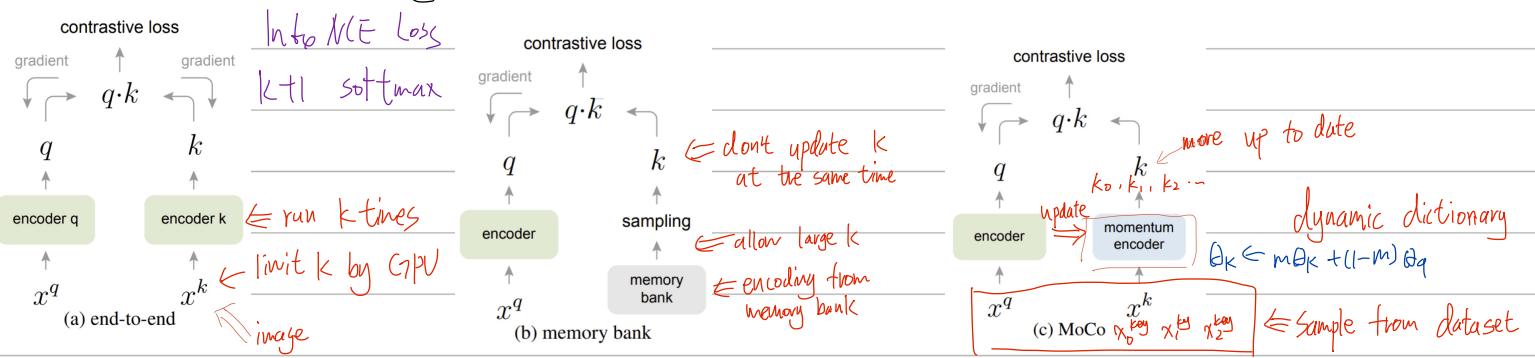
Only update students

Teachers is built from previous iterations of students

Students are for Feature extraction

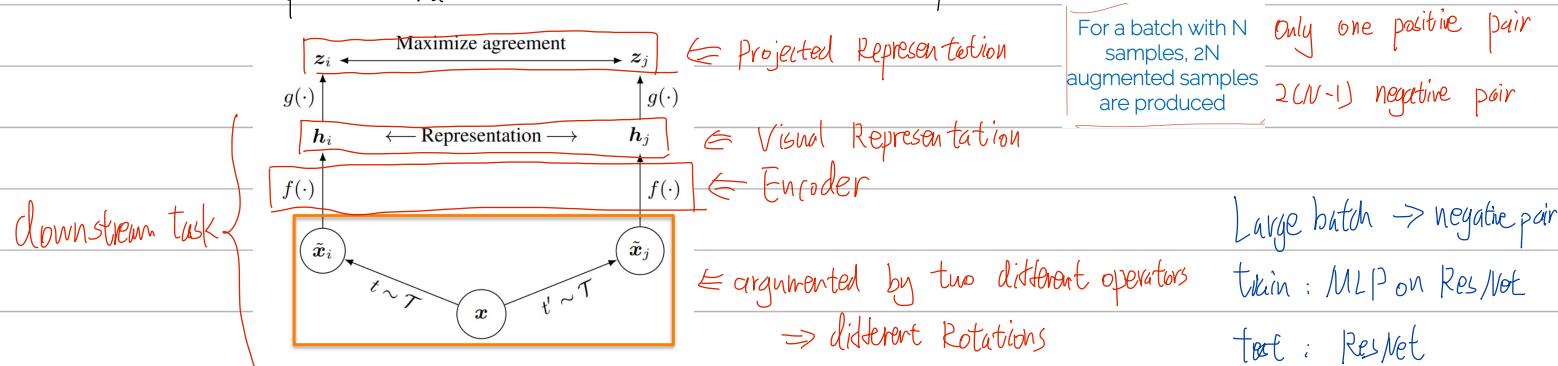


3. Contrastive Learning : Learning embedding space, similar samples pairs are close



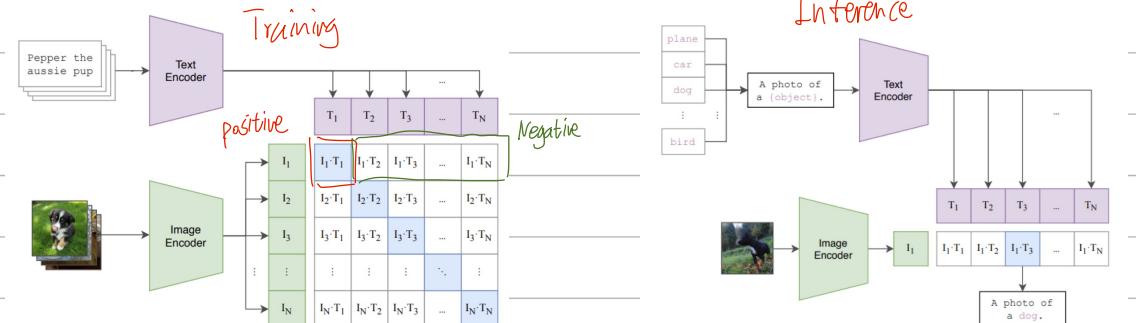
4. SimCLR : Simple Framework for Contrastive Learning of Visual Representation

Supervised via a contrastive loss in the latent space

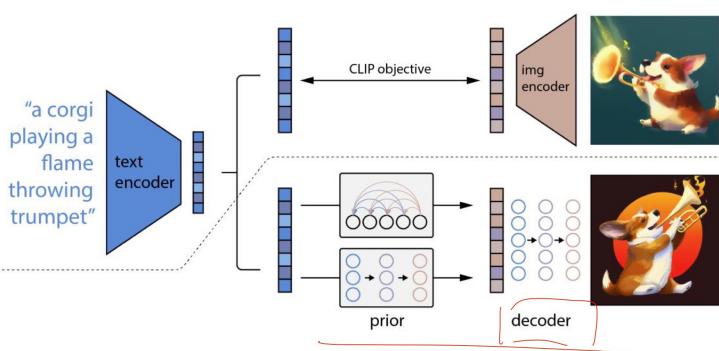


5. Multi-Modal Representation Learning : image \leftrightarrow text pair

CLIP: Contrastive Language - Image Pre-train



Application : text \Rightarrow image (DALL-E₂)



Training

Inference

DS - Sequence Models

1. Sequence Modelling : Text : Classification / Translation / Generation

Image : Classification / Reconstruction

Image patch

2 RNN : Recurrent Neural Network \Rightarrow short sequence / long-term dependency issue

LSTM : Long-short Term Memory \Rightarrow long-term sequence / still not for extreme long

All words are not equal important

3 Attention : Soft attention : Attend each part / sum (w) = 1 / Deterministic and differentiable

Hard attention : Attend one part / stochastic and non-differentiable / Monte-Carlo \rightarrow gradient

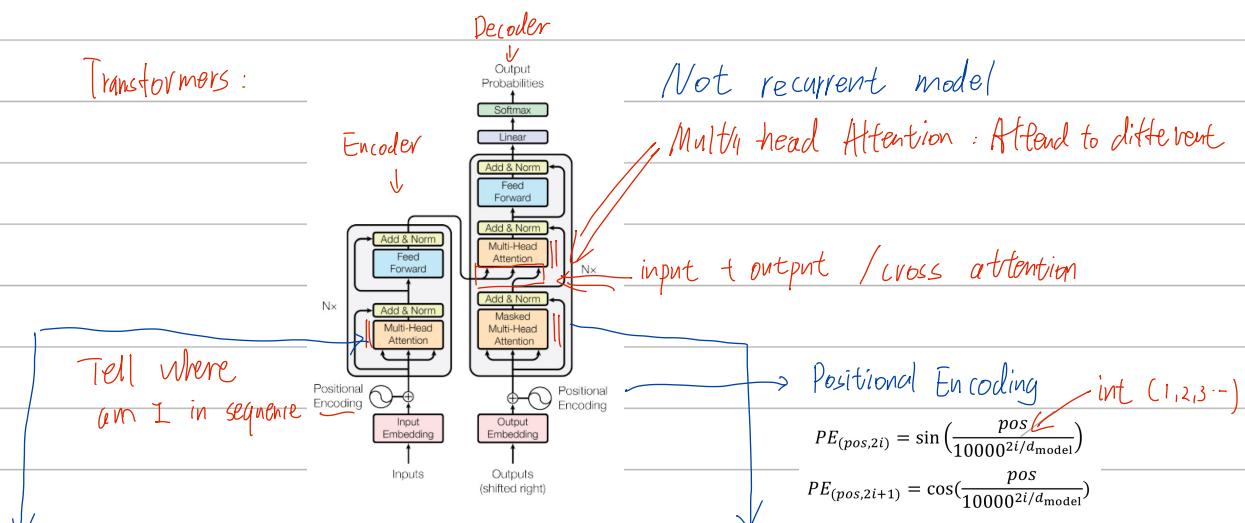
Self-attention : Attend input signal itself

Cross-attention : Attend other input signal as side information

4. Transformer in Language: Attention : Attention (Q, K, V) = softmax $\left(\frac{QK^T}{\sqrt{d_k}} \right) V$

tell how important the Location is
match the look-up table
tell which word I should look up

Transformers :



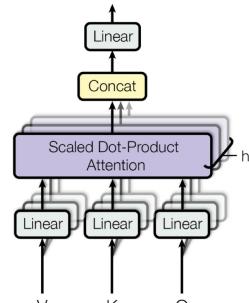
Scale dot-product attention

$$\text{Attention } (Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

Triangular masking for unidirectional modelling : attend past

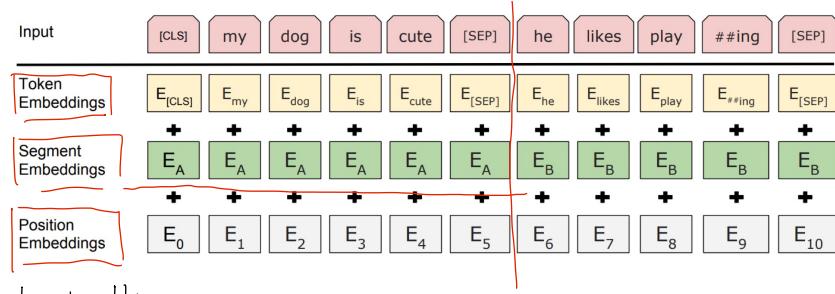
Full masking for bidirectional modelling : attend past and future

Multi-Head Attention



Different heads attend to different parts

BERT : Big Transformer as text encoder



Two unsupervised task :

- MLM
- NSP

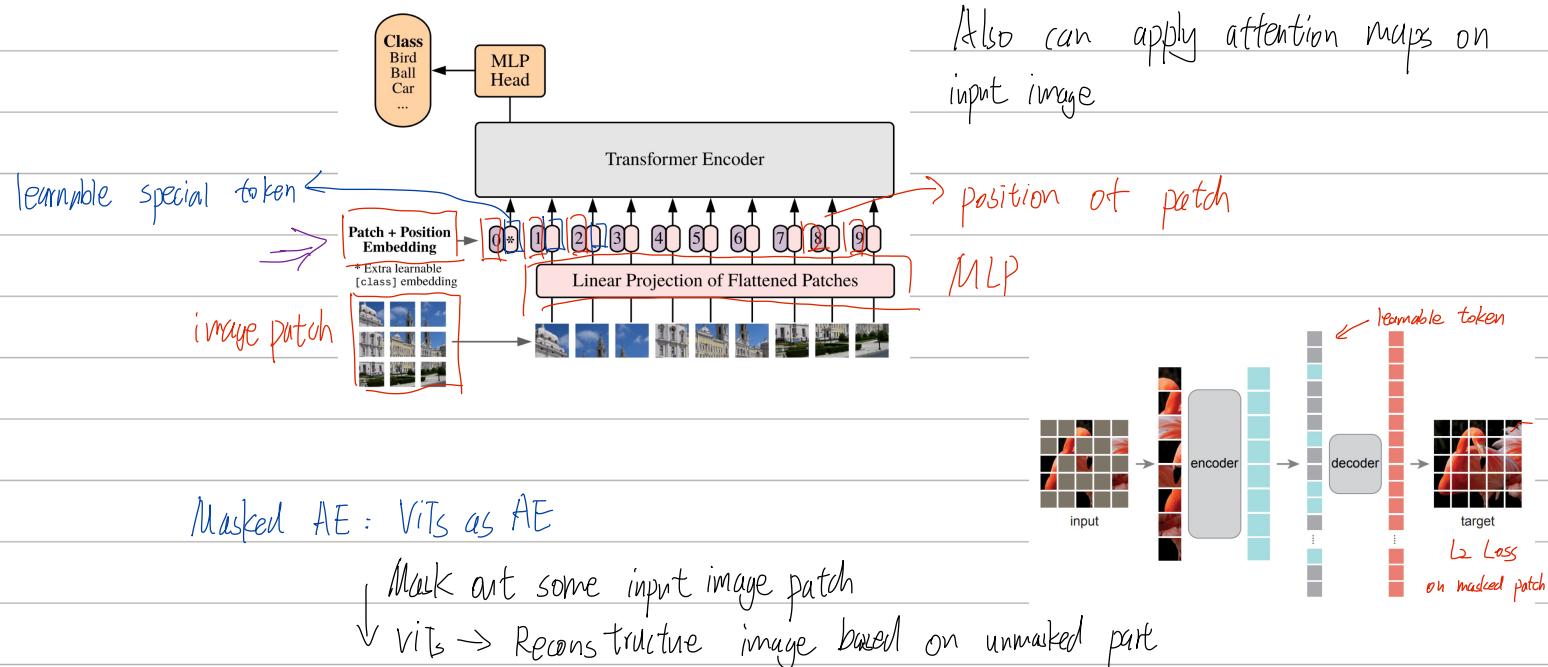
Masked Language Modelling (MLM)

- ↓ Randomly mask out some input word
- ↓ Predict the mask word

Next Sentence Prediction (NSP)

- Take two sentence A and B
- 70% real B and 30% random B
- ↓ Predict B based on A

5 Transformer in CV: Vision Transformers (ViTs)



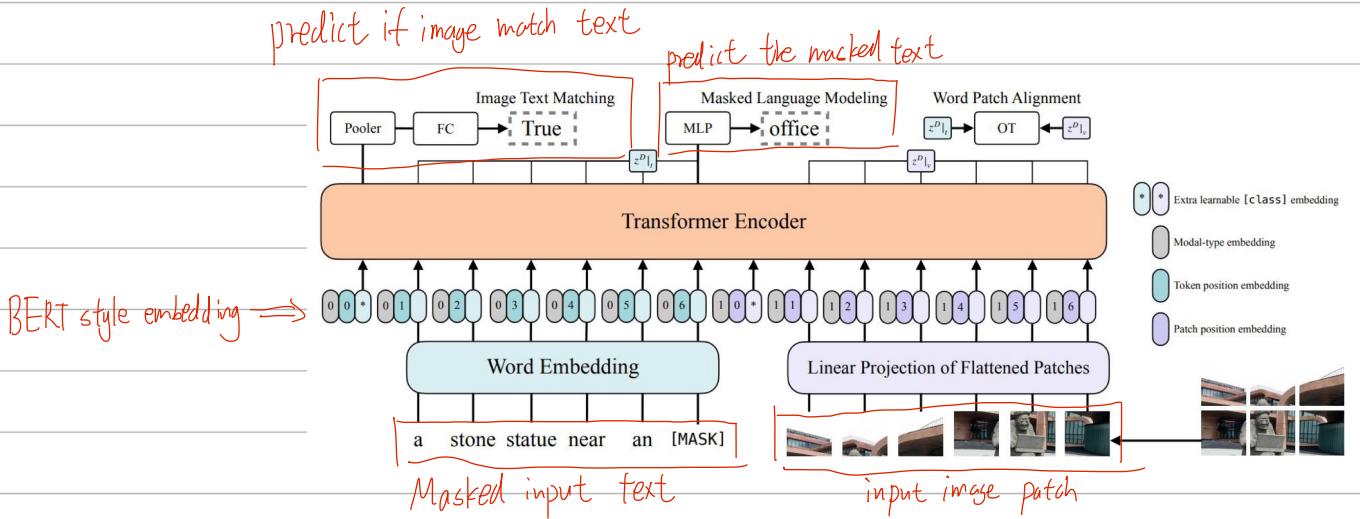
6. Transformer in Multimodal Learning: Vision - Language Transformers (ViLT)

Without supervision

Input: concatenate image patches with text sequences

Pre-training with two self-supervised objectives

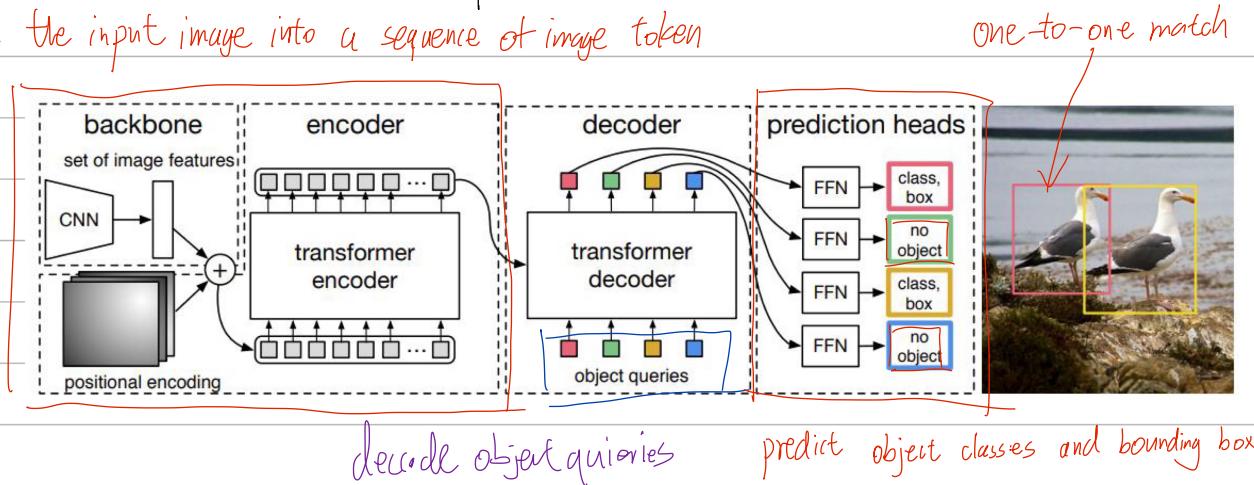
|| \Rightarrow { Image text matching
Masked Language modelling



7. Transformer in Object Detection: Detection Transformer (DETR)

on top of CNN

encode the input image into a sequence of image token



Matching via Hungarian Algorithm

Find prediction-GT assignments that minimize this cost

$$\hat{\sigma} = \arg \min_{\sigma \in \mathfrak{S}_N} \sum_i^N \left[\mathbb{1}_{\{c_i \neq \emptyset\}} \hat{p}_{\sigma(i)}(c_i) + \mathbb{1}_{\{c_i \neq \emptyset\}} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\sigma(i)}) \right]$$

Predicted class probability ("no object" class excluded)

L1 loss and generalized IoU loss

$$\mathcal{L}_{\text{Hungarian}}(y, \hat{y}) = \sum_{i=1}^N \left[-\log \hat{p}_{\hat{\sigma}(i)}(c_i) + \mathbb{1}_{\{c_i \neq \emptyset\}} \mathcal{L}_{\text{box}}(b_i, \hat{b}_{\hat{\sigma}(i)}) \right]$$

Predicted class probability (including the "no object" class)

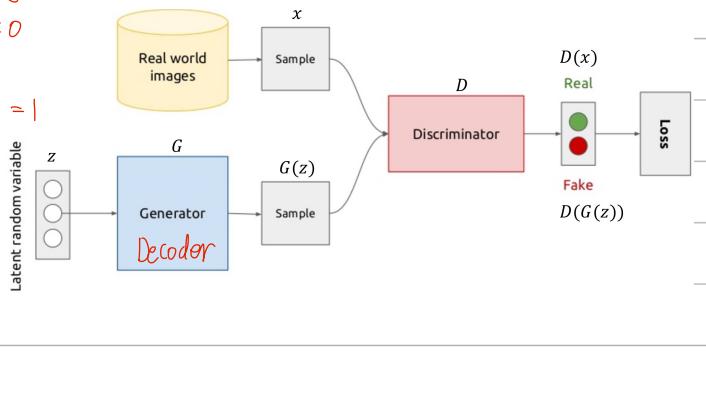
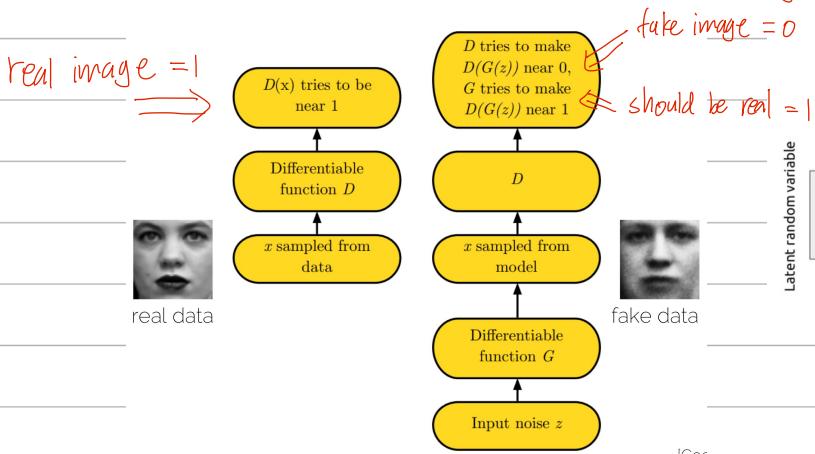
Hungarian Loss \Rightarrow
(after bipartite matching)

06 - GANs - 01

1. Generative Models : Implicit density \rightarrow Direct \rightarrow GAN

2. Generative Adversarial Networks : Decoder in AE as Generative Model

Instead of using L_2 , want to learn a loss function



Loss Function

Discriminator loss

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

binary cross entropy

Generator loss

$$J^{(G)} = -J^{(D)}$$

G minimize probability that D is correct
Equilibrium is saddle point of discriminator loss
 D supervise for G $\Rightarrow 0.5$

Loss Function (Heuristic Method)

Discriminator loss

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

Generator loss

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z}} \log D(G(\mathbf{z}))$$

G maximize the log probability of D being mistaken

G can still learn even D rejects all samples

Final Loss

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log (1 - D(G(\mathbf{z})))].$$

* $L(D)$ should never touch zero \Rightarrow no gradient for G ~~if~~ ~~if~~

Loss curve should keep stable after some iterations

* Adaptive Schedule \Rightarrow while $\text{loss_discriminator} > t_d$:

train discriminator

while $\text{loss_generator} > t_g$:

train generator

* Balance needed ! weak $D \Rightarrow$ No good gradient
weak $G \Rightarrow$ D always right

3 Mode Collapse: \Rightarrow no recover problem (restart training)

$$\boxed{\min_G \max_D V(G, D) \neq \max_D \min_G V(G, D)}$$

D in inner loop \Rightarrow convergence to correct dist

G in inner loop \Rightarrow convergence to one sample

More modes \Rightarrow small recovery rate e.g. GAN \rightarrow Face

Large latent space \Rightarrow more mode collapse

Evaluation metrics: ① Human Evaluation: Visualization / train curves

② Inception Score (IS): Measure saliency and diversity

① Train a classifier

② Train a image generation model

③ Test the generation model by classifier

Saliency: check whether generated images can be classified with high confidence

Diversity: check whether obtain samples from all classes

③ Frechet Inception Distance (FID)

calculate the feature distance between the real and synthetic distribution

④ Strong D \Rightarrow good G

...

4. GAN Hack: ① Normalized Input: between -1 and 1 / Tanh as the last layer

② Sampling: spherical z / Gaussian Distribution

③ Batch Norm:

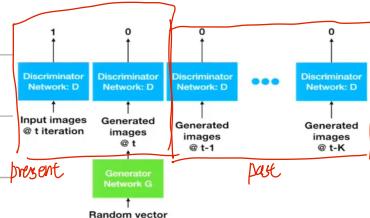
④ ADAM: SGD for D / ADAM for G

⑤ One-sided Label smoothing: prevent D give too large gradient to G

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{x \sim p_{\text{data}}} \log D(x) - \frac{1}{2} \mathbb{E}_z \log (1 - D(G(z)))$$

Some value smaller than 1; e.g., 0.9

⑥ Historical Generator Batches: D also use historical generated images



⑦ Avoid Sparse Gradient: Leaky ReLU / ...

⑧ D noisy: Due to SGD \Rightarrow exponential average of weights

4 Loss Function: Heuristic is standard / Loss alone will not affect the Network

① EBGAN: D is an AE

Reconstruction error of $D(x)$ should be low

Penalize D if reconstruction error below a value m

$$\rightarrow D(x) = \|Dec(Enc(x)) - x\|$$

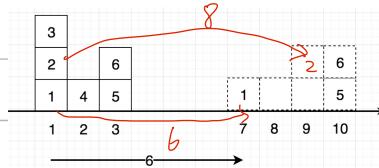
$$\mathcal{L}_D(x, z) = D(x) + [m - D(G(z))]^+$$

$$\mathcal{L}_G(z) = D(G(z))$$

where $[u]^+ = \max(0, u)$

② BEGAN: Measure difference in data distribution of real and generated images

③ WGAN:



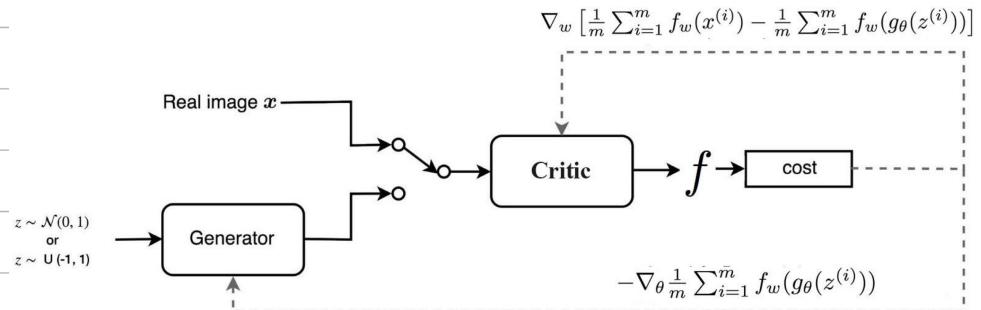
Two Expectations should match

$$W(\mathbb{P}_r, \mathbb{P}_\theta) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] - \mathbb{E}_{x \sim \mathbb{P}_\theta}[f(x)]$$

1-Lipschitz function: upper bound between densities

$$|f(x_1) - f(x_2)| \leq |x_1 - x_2|.$$

Network



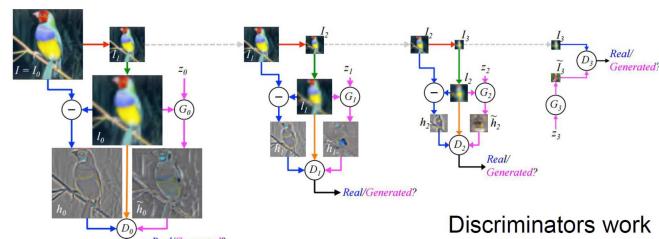
Comparison

| | Discriminator/Critic | Generator |
|------|---|---|
| GAN | $\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D(x^{(i)}) + \log (1 - D(G(z^{(i)})))]$ | $\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m [\log (D(G(z^{(i)})))]$ |
| WGAN | $\nabla_w \frac{1}{m} \sum_{i=1}^m [f(x^{(i)}) - f(G(z^{(i)}))]$ | $\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m [f(G(z^{(i)}))]$ |



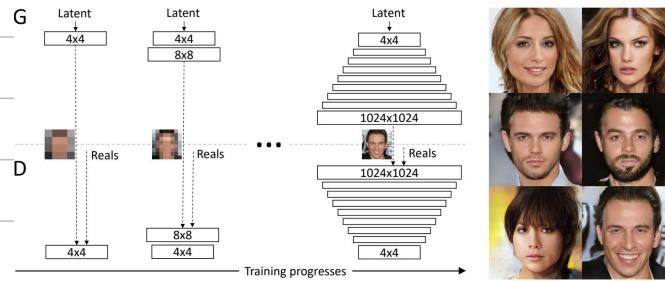
Gratest improvement \Rightarrow WGAN Loss can converge instead of stable terrible weight clipping

5. GAN Architecture : ① Multiscale GANs

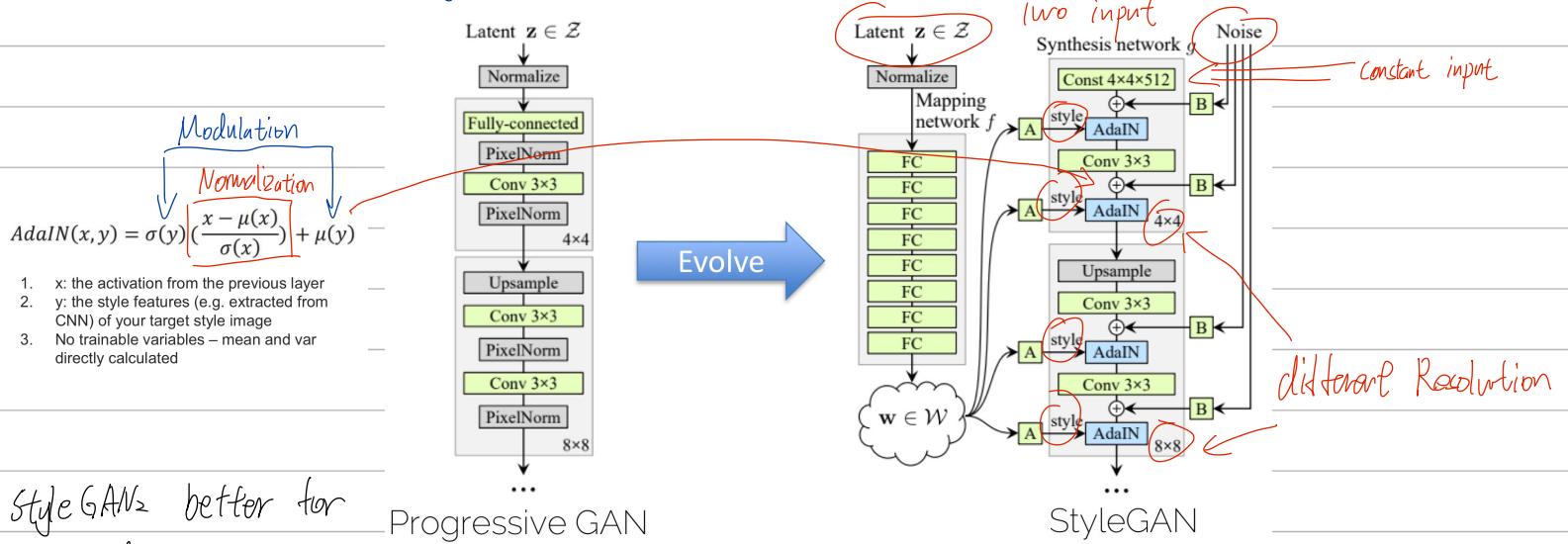


Discriminators work at every scale!

② Progressive Growing GANs

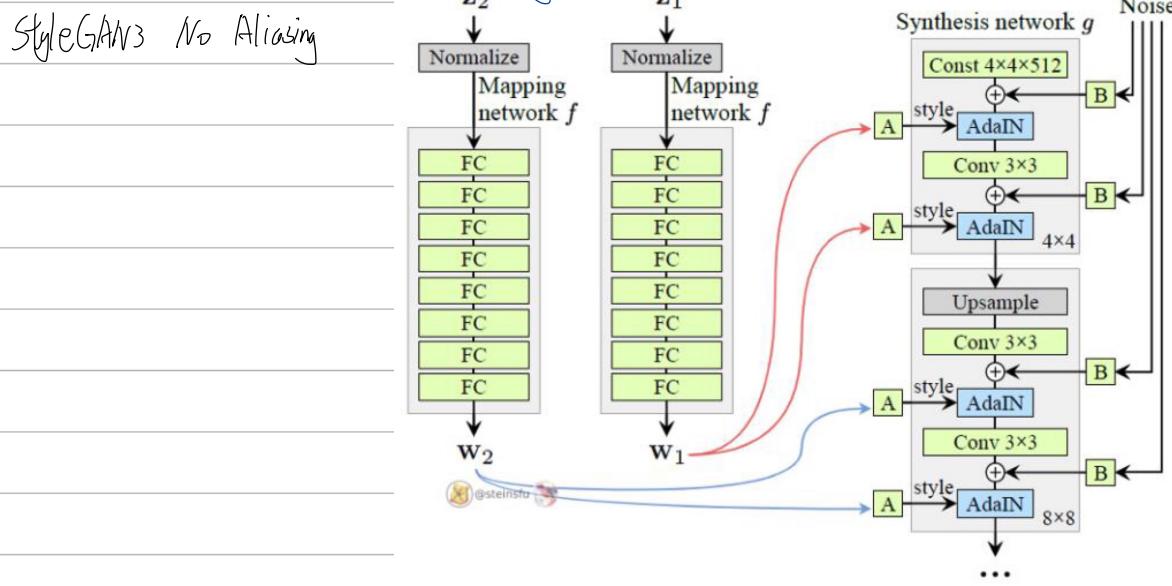


③ Style GAN



Style GANs better for limited data (No Argumentation)

Mixing Regularization

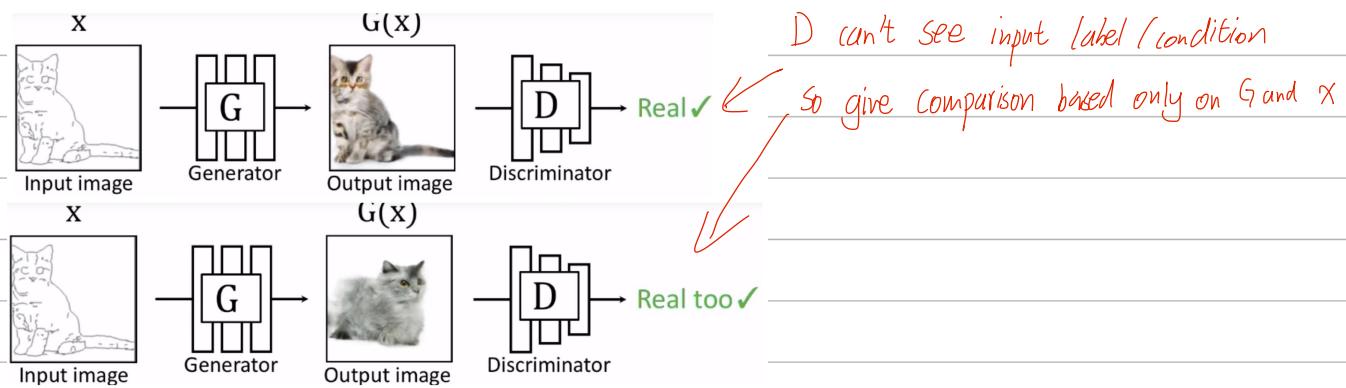


07 - GANs - 02

1. Conditional / Generative Adversarial Networks (cGAN):
 - ① Gain control of output
 - ② Modeling: add semantic meaning to latent space
 - ③ Domain transfer: Label on A \rightarrow transfer to B, train D on A, test B

2 GAN Manifold: $G(z_0) \rightarrow$ Linear interpolation in Z space: $G(z_0 + t \cdot (z_1 - z_0)) \rightarrow G(z_1)$

3 cGAN:



4 iGAN: ① Projection image on Manifold Input: real image x / output: latent vector z

\hookrightarrow ① Optimization:

$$z^* = \arg \min \mathcal{L}(G(z), x^R)$$

Reconstruction loss L

Generative model $G(z)$

② Inverse Network:

Inverting Network $z = P(x)$

$$\theta_P^* = \arg \min_{\theta_P} \sum_{x_n^R} \mathcal{L}(G(P(x^R; \theta_P)), x^R)$$

Auto-encoder

with a fixed decoder G

Prof. Niessner

③ Hybrid: Inverse Network + Optimization

④ Manipulating the latent vector

Objective:
$$z^* = \arg \min_{z \in \mathbb{Z}} \left\{ \underbrace{\sum_g (\mathcal{L}_g(G(z), v_g) + \lambda_s \cdot \|z - z_0\|_2^2)}_{\text{data term}} \right\}$$

Guidance v_g



$G(z)$

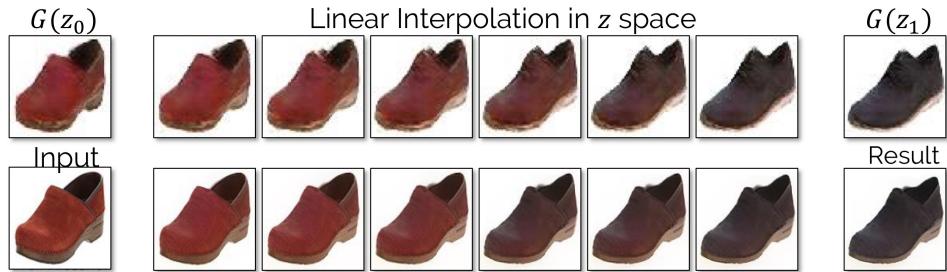
z_0

③ Edit Transfer

Motion (u, v) + Color ($A_{3 \times 4}$): estimate per-pixel geometric and color variation

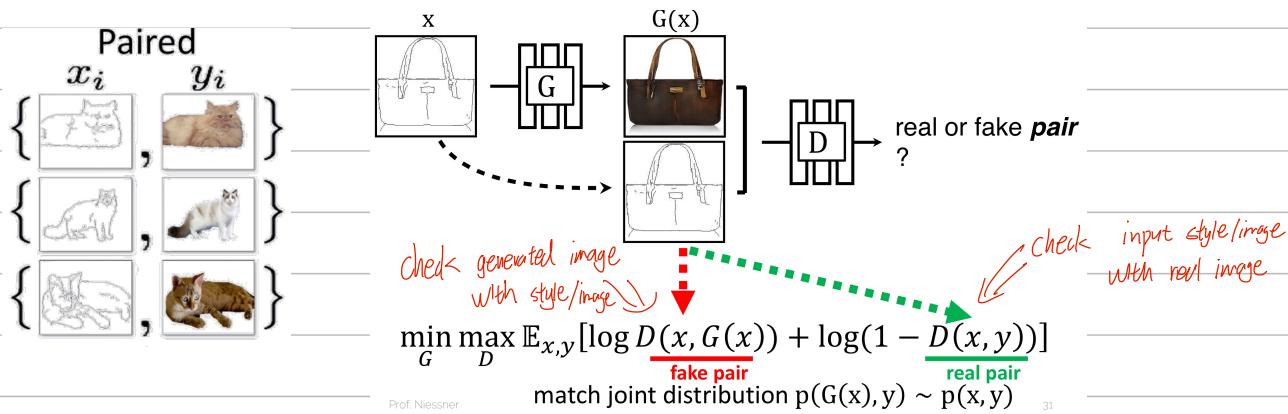
$$\iint \underbrace{\|I(x, y, t) - A \cdot I(x+u, y+v, t+1)\|^2}_{\text{data term}} + \underbrace{\sigma_s (\|\nabla u\|^2 + \|\nabla v\|^2)}_{\text{spatial reg}} + \underbrace{\sigma_c \|\nabla A\|^2}_{\text{color reg}} dxdy$$

(constraint motion) (constraint colour)



cGAN : Interactive GAN (User edit \Rightarrow generated (customized image))

5 Mapping problem : ① Pair : Pix2Pix



$$\Rightarrow L = L_{GAN} + \lambda L_1 \text{ more constraint}$$

- ↳ Unet preserve the structure
- ↳ cGAN tend to ignore random vector z
- ↳ L_1/L_2 for low frequency details
- ↳ D for high frequency details

Pix2Pix HD

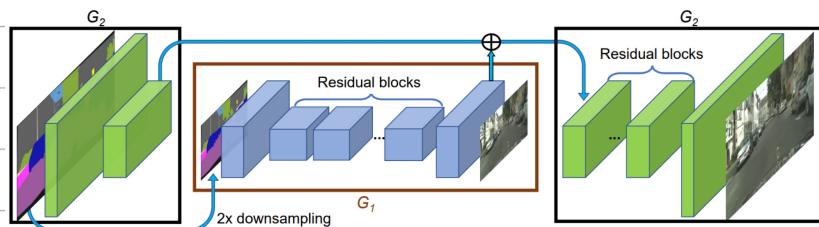
Expand Pix2Pix to multi-scale

Same G and D but on different resolutions \rightarrow large receptive fields

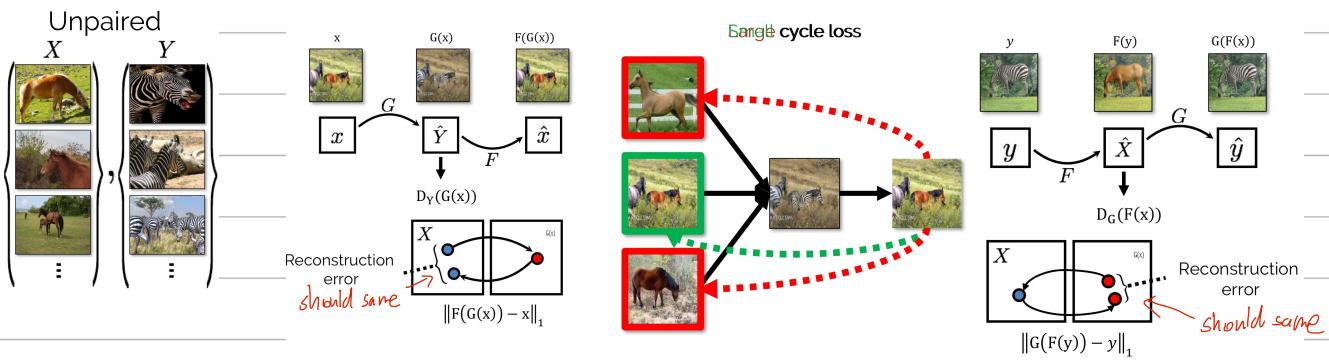
Use multi-scale D

\Rightarrow different receptive fields

$$\min_G \max_{D_1, D_2, D_3} \sum_{k=1,2,3} L_{GAN}(G, D_k)$$



② Unpaired: Cycle-Consistent Adversarial Network



GAN don't force output is as same as input
 \Rightarrow use F map output return to input

b CGAN Application : Image Manipulation / Image Generation / Image Restoration / Image Interpolation
 Image Understanding / Multi-modal Manipulation / 3D Aware Synthesis