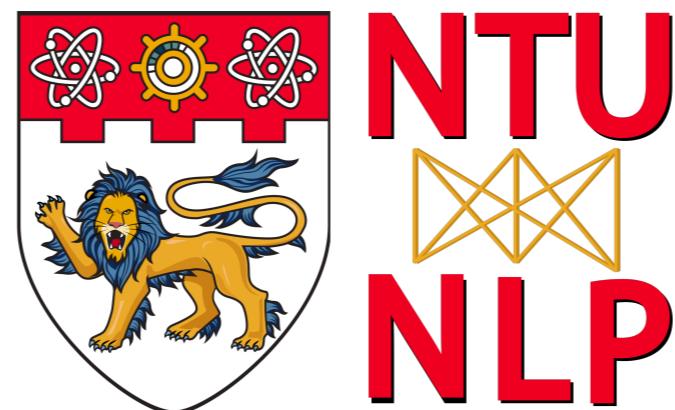


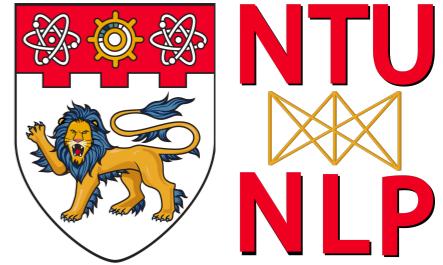
Deep Learning for Natural Language Processing

Shafiq Joty



Lecture 12: Adversarial NLP

Where we are



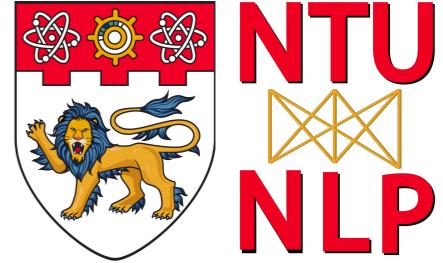
Models/Algorithms

- Linear models
- Feed-forward Neural Nets (FNN)
- Window-based methods
- Convolutional Nets
- Recurrent Neural Nets
- Recursive Neural Nets

NLP tasks/applications

- Word meaning
- Language modelling
- Sequence tagging
- Sequence encoding
- Parsing
- Hierarchical encoding

Where we are



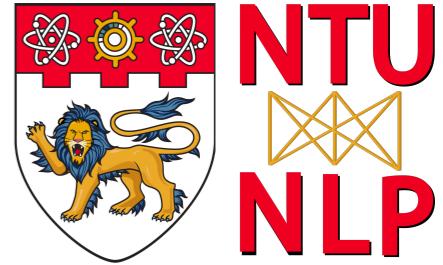
Models/Algorithms

- Seq2Seq
 - + Attention
 - + Subword
- Seq2Seq Variants
- Transformer Seq2Seq

NLP tasks/applications

- Machine Translation
- Summarization
- Parsing
- Dialogue generation

Where we are

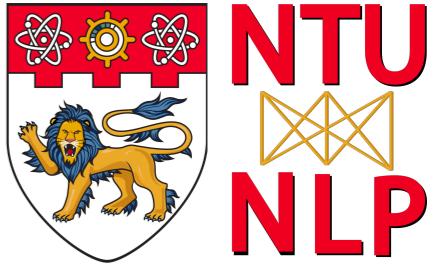


Models/Algorithms

NLP tasks/applications

- Self-supervised pretraining
- Multilingual NLP
- Robust & adversarial NLP
- Ethics, Bias in NLP

Today's Plan: Adversarial NLP



Adversarial Nets

- Generative adversarial nets (GANs)
- Domain adversarial nets (DANs)

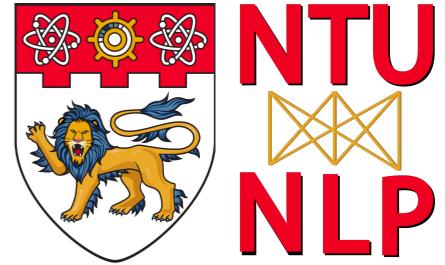
NLP tasks/applications

- Machine Translation
- Domain Transfer

Adversarial Examples & Training

- Adversarial attacks
- Defense
- Robust NLP

Adversarial Nets



Adversarial nets have been used in **two ways**:

① **Generative modeling:** GANs [Goodfellow et al., 2014].

- Generate real-like image samples from random vectors
- One neural net is pit against another:
 - (*i*) Generator, (*ii*) Discriminator
- Other Examples: image-to-image translation, Unsup. NMT

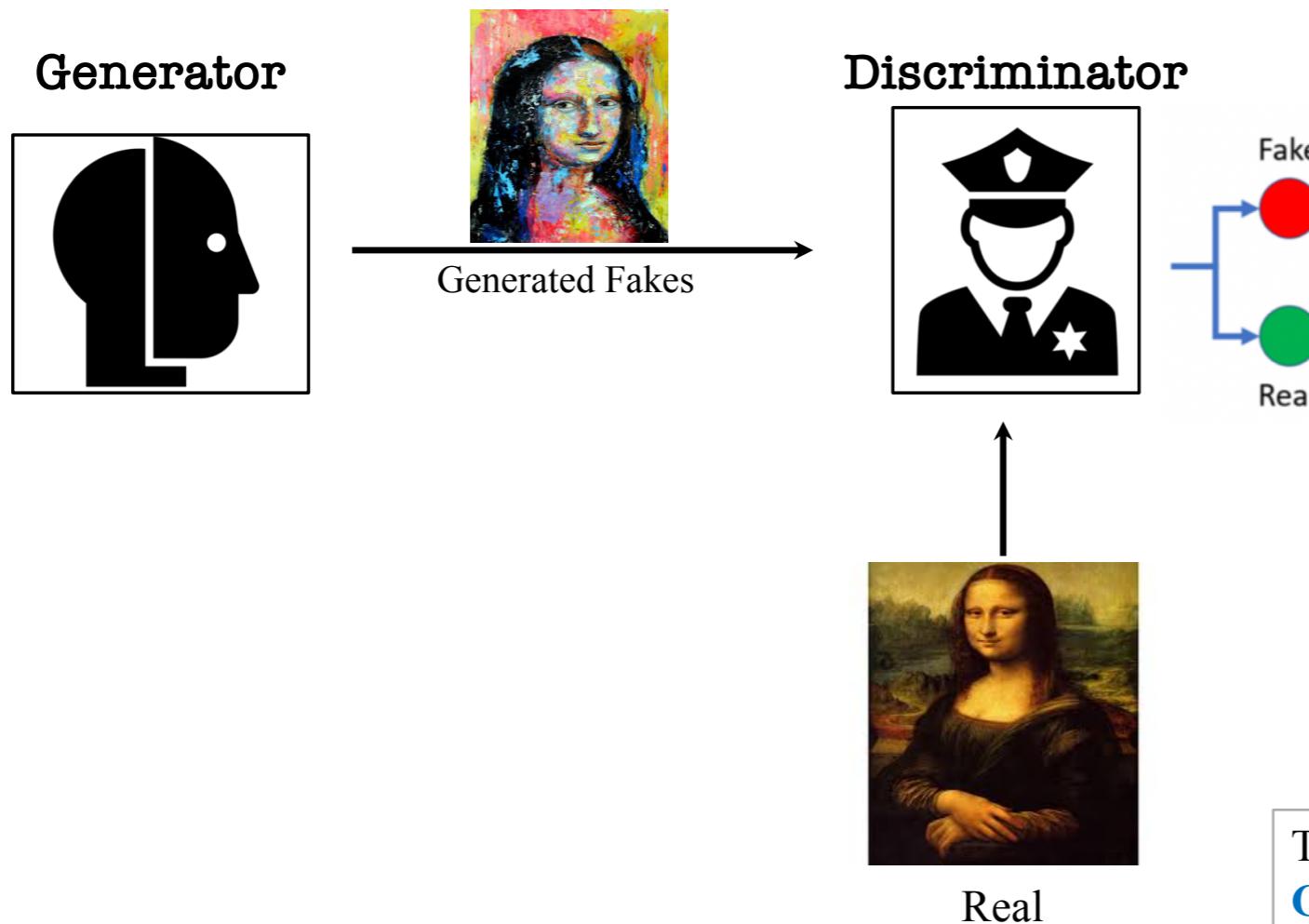
② **Transfer learning/domain adaptation:**

DANNs [Ganin et al., 2016]

- Map samples from two domains into a common feature space
- Generally a **three-player** game:
 - (*i*) Encoder, (*ii*) Classifier (*iii*) Discriminator
- Examples: MNIST \Rightarrow USPS, X-lingual NER

Generative Adversarial Nets

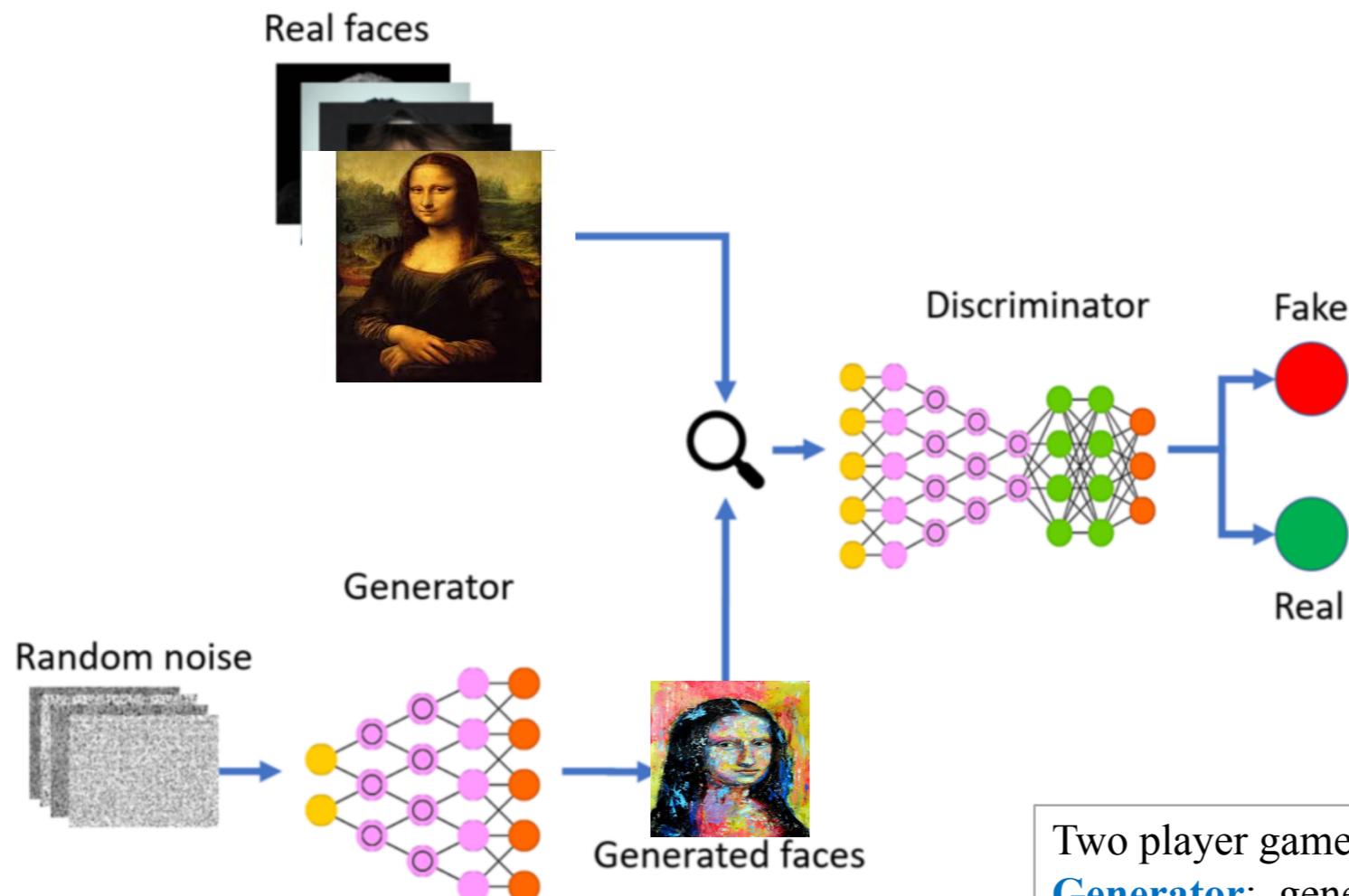
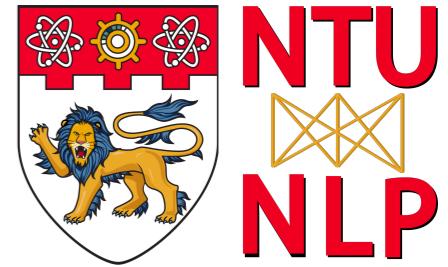
[Goodfellow et al., 2014]



Two player game, where
Generator: generates **fake** examples
Discriminator: differentiate between **real** and **fake**

Generative Adversarial Nets

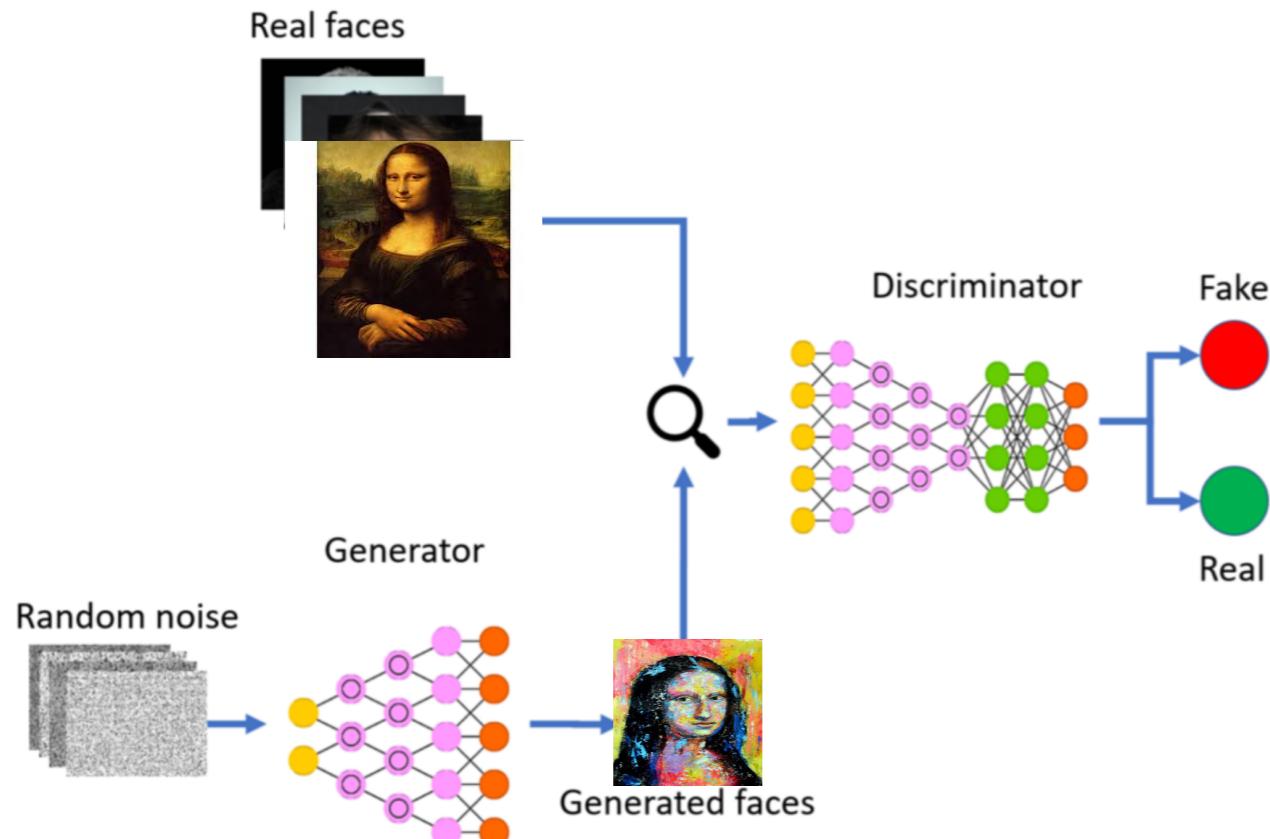
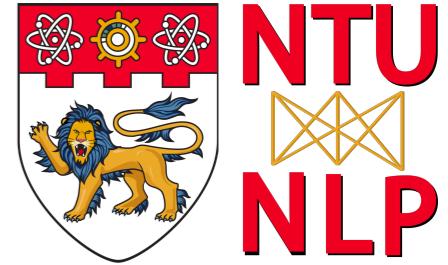
[Goodfellow et al., 2014]



Two player game, where
Generator: generates **fake** examples
Discriminator: differentiate between **real** and **fake**

Generative Adversarial Nets

[Goodfellow et al., 2014]



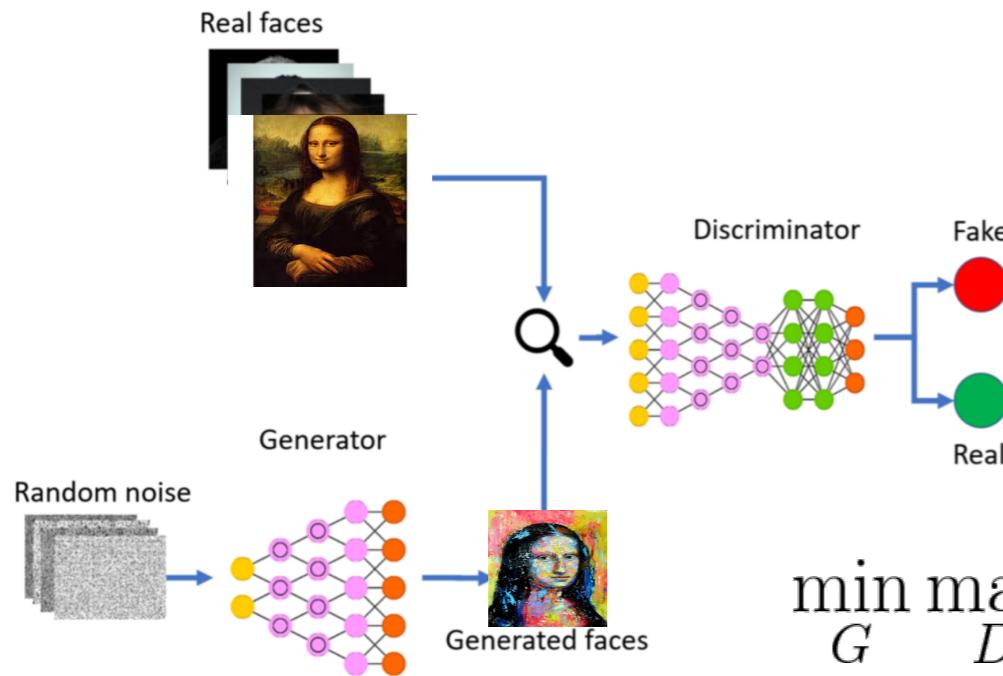
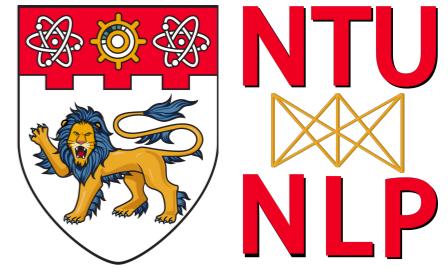
$$\min_G \max_D V(D, G)$$

$$= \mathbb{E}_{q(\mathbf{x})}[\log(D(\mathbf{x}))] + \mathbb{E}_{p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$$

$D(x)$: the probability that x came from the data rather than generator

Generative Adversarial Nets

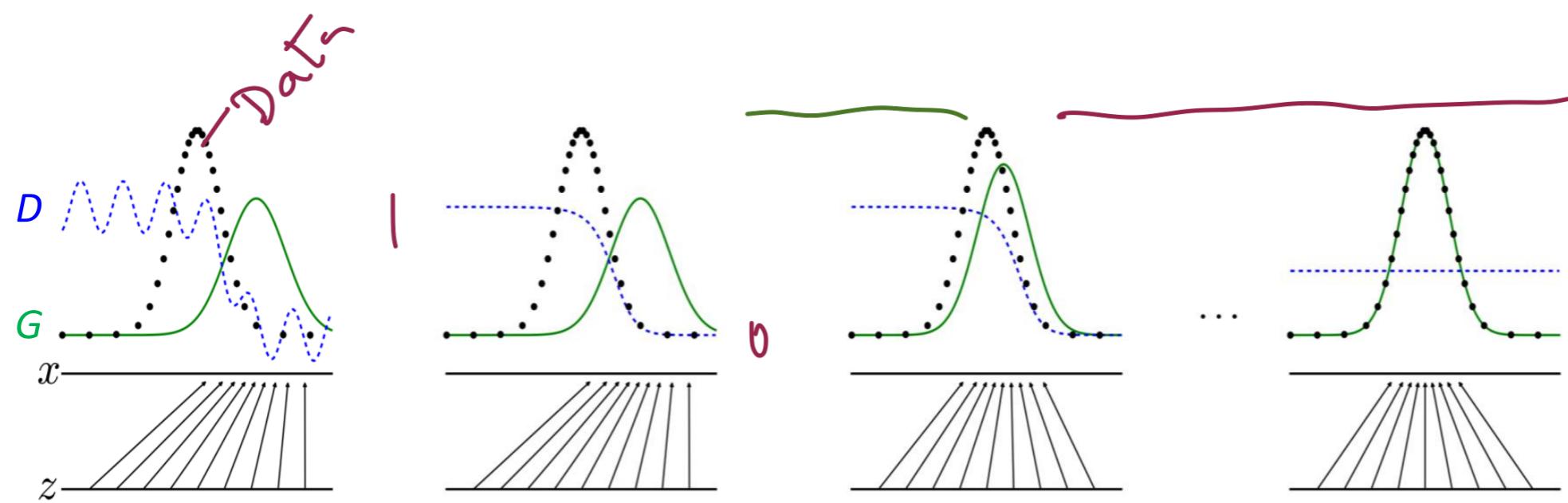
[Goodfellow et al., 2014]



$$\min_G \max_D V(D, G)$$

$$= \mathbb{E}_{q(\mathbf{x})}[\log(D(\mathbf{x}))] + \mathbb{E}_{p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$$

$D(x)$: the probability that x came from the data rather than generator



Generative Adversarial Nets

[Goodfellow et al., 2014]

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

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

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(\mathbf{x}^{(i)}) + \log (1 - D(G(\mathbf{z}^{(i)}))) \right].$$

end for

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

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(\mathbf{z}^{(i)}))).$$

end for

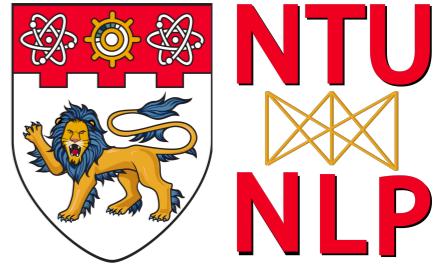
The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Discriminator

Generator

Generative Adversarial Nets

[Goodfellow et al., 2014]

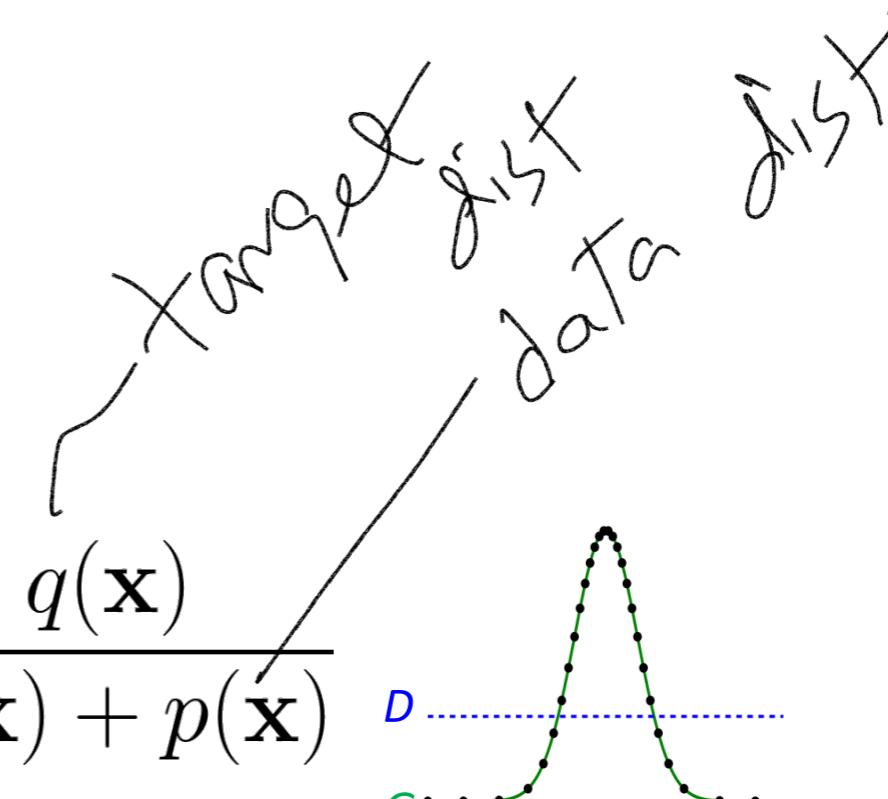


GAN Equilibrium

- Global optimality
 - Discriminator
 - Generator

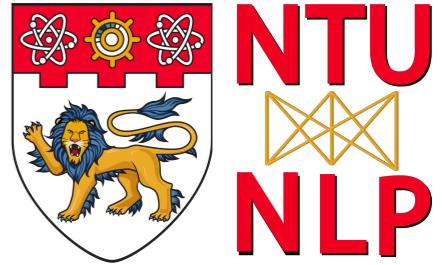
$$D^*(\mathbf{x}) = \frac{q(\mathbf{x})}{q(\mathbf{x}) + p(\mathbf{x})}$$

$$G^*(\mathbf{z}) \text{ s.t. } p(\mathbf{z}) = q(\mathbf{x})$$

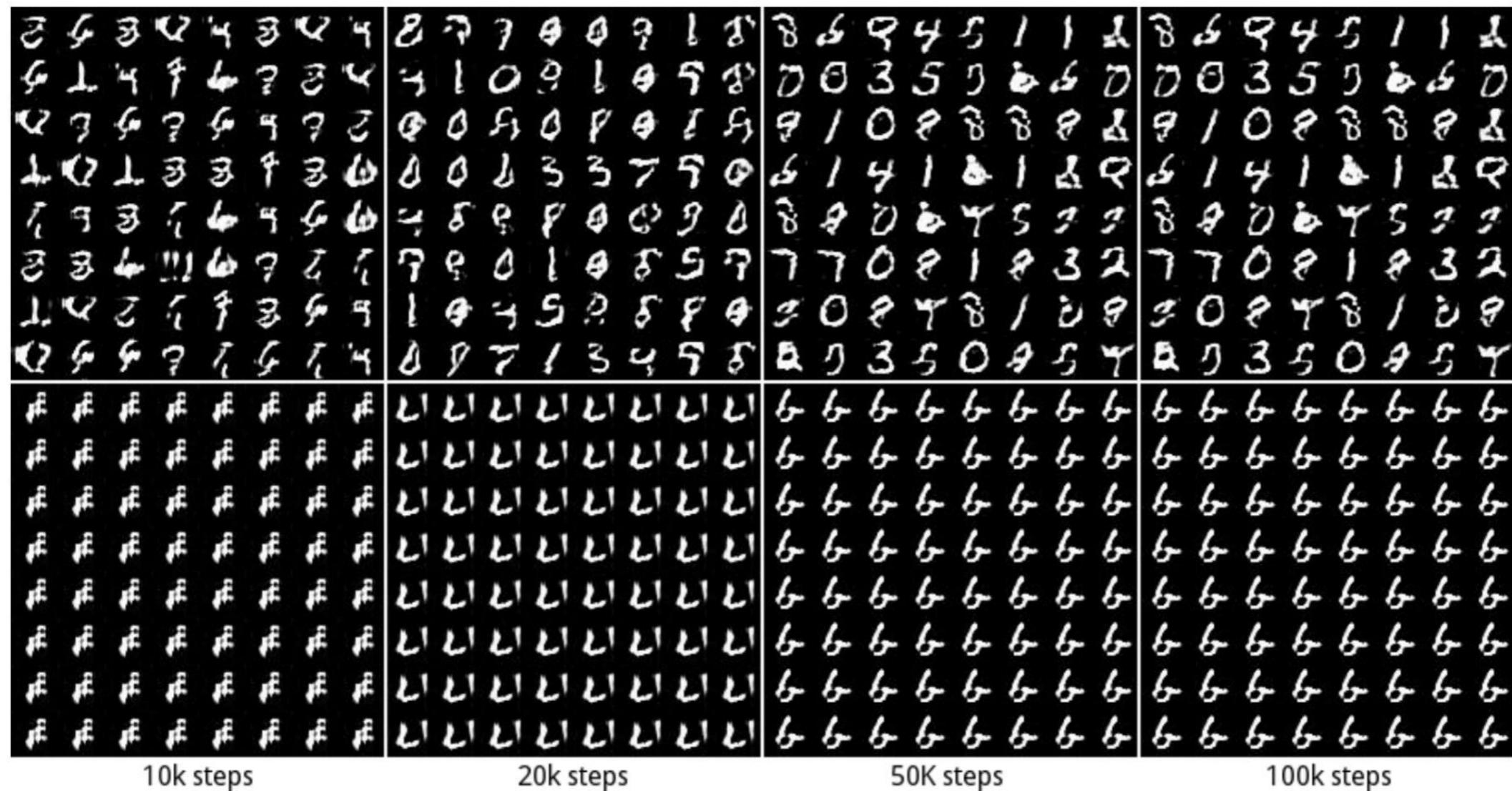


target dist
data dist
Mode
dist

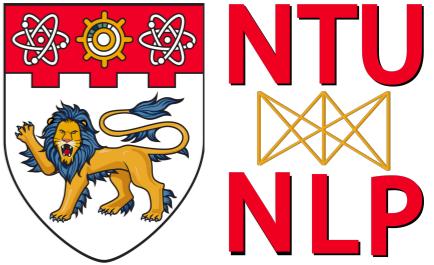
Major Issues of GAN Training



- Mode Collapse (unable to produce diverse samples)



Major Issues of GAN Training

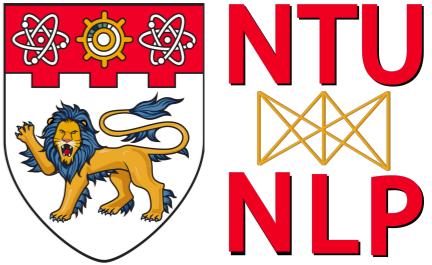


- Unstable Adversarial Training
 - We are dealing with two networks / learners / agents
 - Should we update them at the same rate?
- The discriminator might overpower the generator.

Particularly for NLP

- GANs were originally designed for images, which is continuous — $D(G(\mathbf{z}))$
- But text is discrete — can't back-propagate

Major Issues of GAN Training



Particularly for NLP

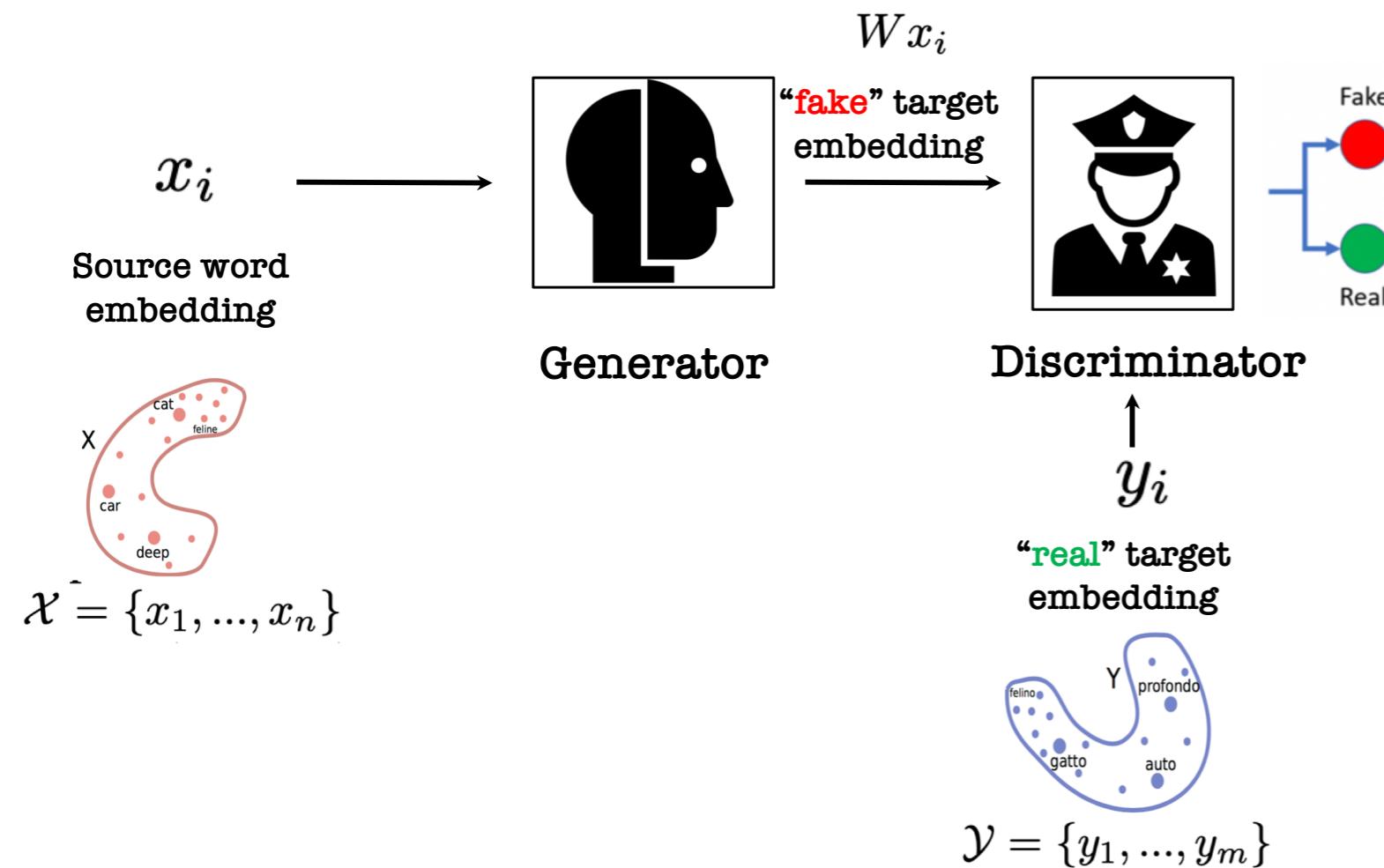
- GANs were originally designed for images, which is continuous — $D(G(\mathbf{z}))$
- But text is discrete — can't back-propagate

Some workaround

- Use it in **embedding** space
- Use Reinforcement learning with **policy gradient**

Applications of GANs in NLP

Unsupervised Word Translation



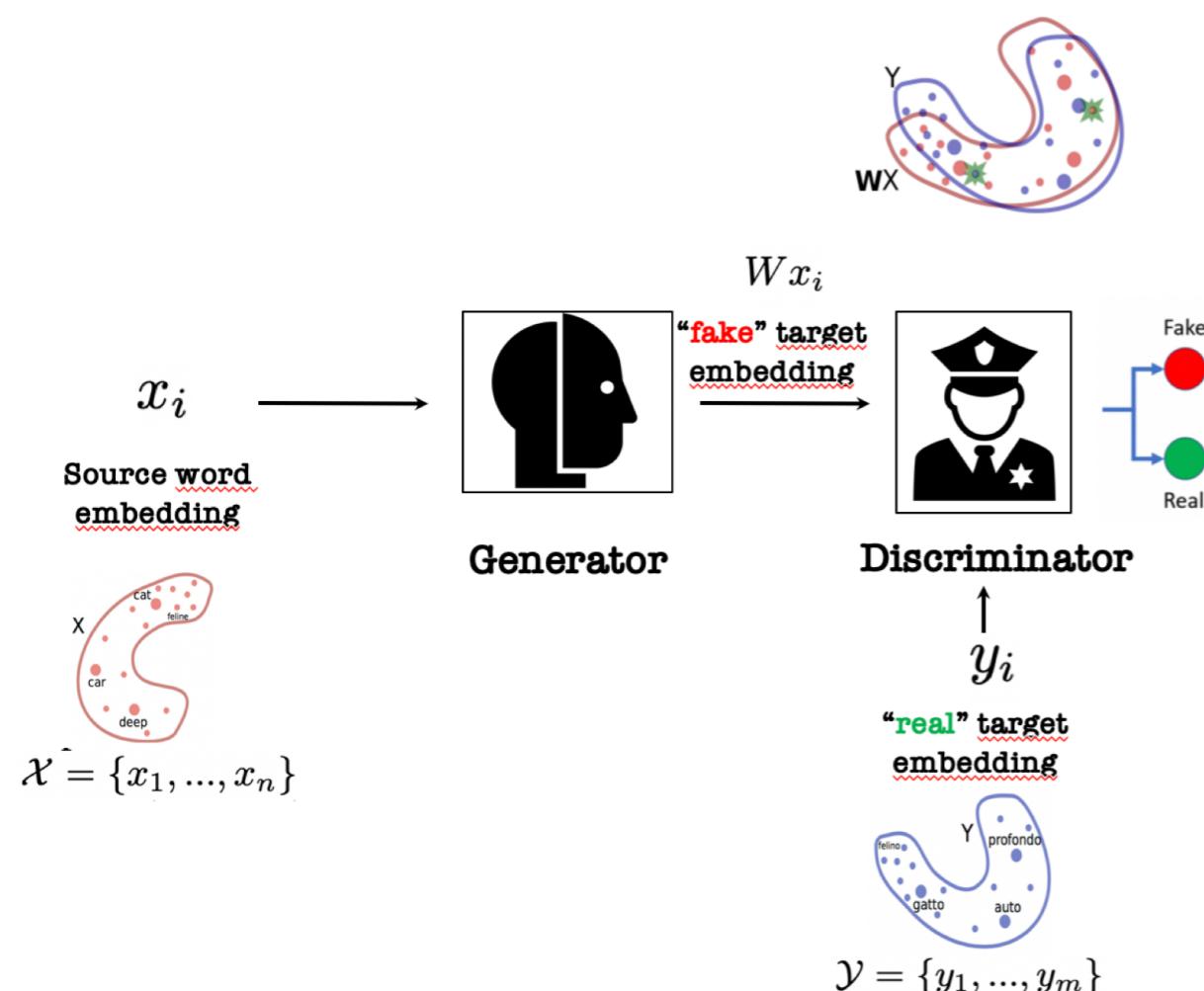
- Generator needs to **match distribution** of target language in order to **fool discriminator** consistently.
- **Hypothesis:** Best way to do this is to align words with their translations.

Generator: projects source word embedding x_i into the target language using linear mapper $W \Rightarrow Wx_i$

Discriminator: differentiate between “fake” projected embeddings and “true” target language embeddings y_i

Applications of GANs in NLP

Unsupervised Word Translation



Objective

$\theta_D \Rightarrow$ Discriminator parameters

$P_{\theta_D}(\text{source} = 1|z) \Rightarrow$ Probability that z is the mapping of a source embedding

Discriminator Loss

$$\mathcal{L}_D(\theta_D|W) = -\frac{1}{n} \sum_{i=1}^n \log P_{\theta_D}(\text{source} = 1|Wx_i) - \frac{1}{m} \sum_{i=1}^m \log P_{\theta_D}(\text{source} = 0|y_i)$$

An embedding in X

An embedding in Y

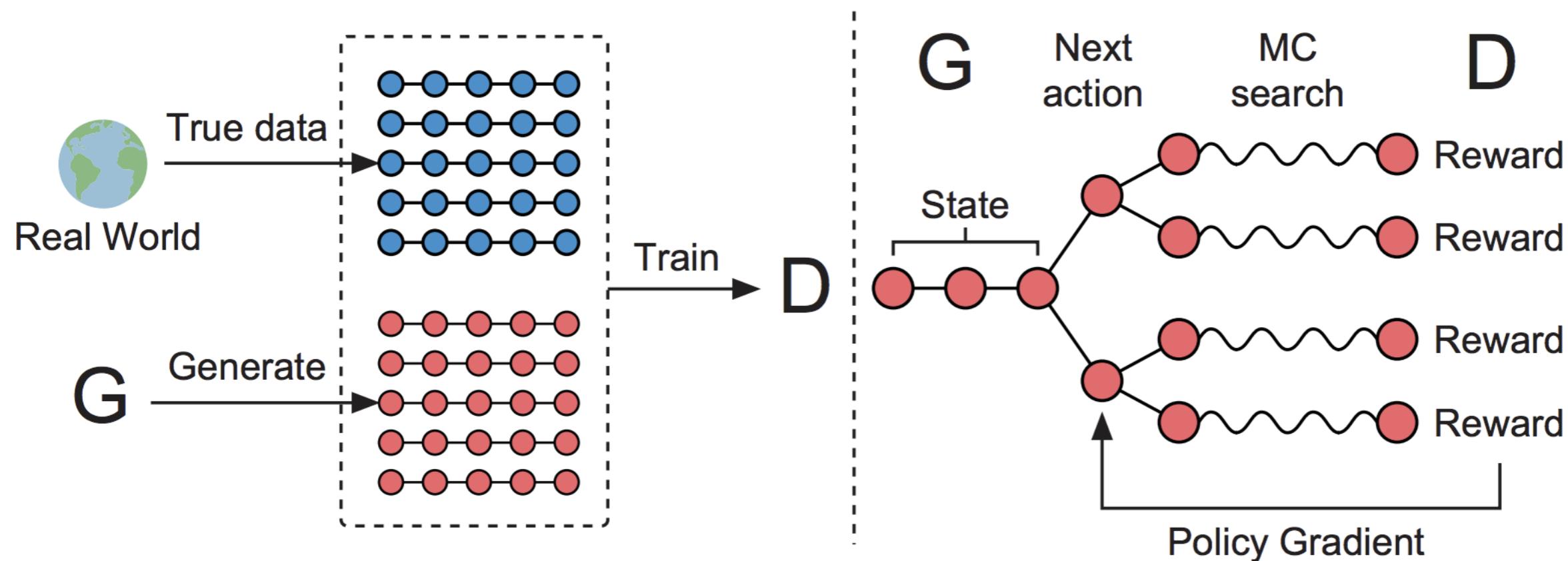
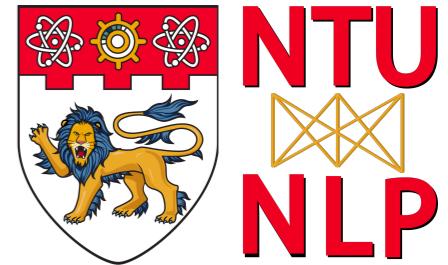
Maximize probability of predicting correct source

Generator/Mapper Loss

$$\mathcal{L}_W(W|\theta_D) = -\frac{1}{n} \sum_{i=1}^n \log P_{\theta_D}(\text{source} = 0|Wx_i) - \frac{1}{m} \sum_{i=1}^m \log P_{\theta_D}(\text{source} = 1|y_i)$$

Maximize probability of fooling discriminator

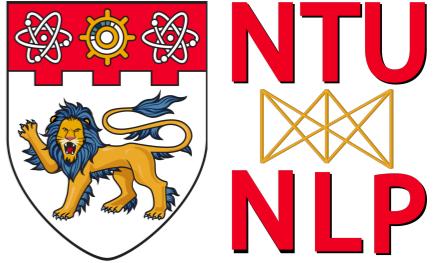
SeqGAN: policy gradient for generating sequences



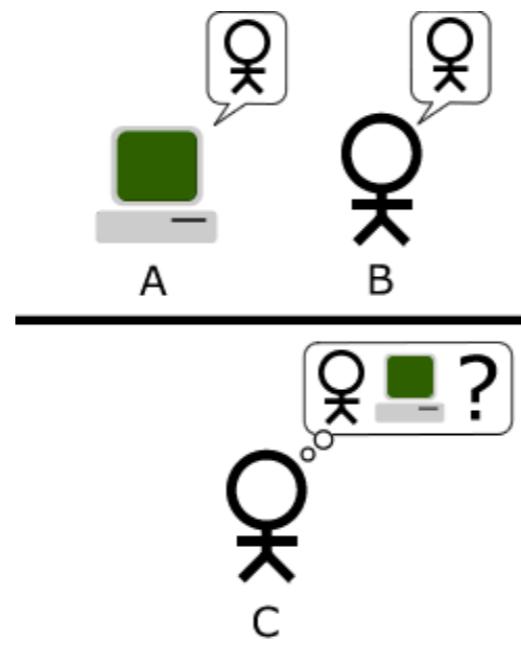
The discriminator is often learning a metric.

(Yu et al., 2017)

GAN for Dialogue



Turing Test



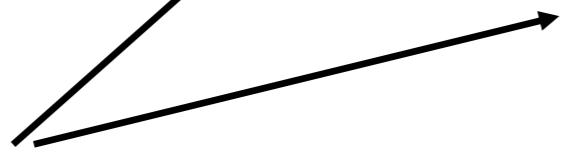
Reward for Dialogue

How old are you ?



A human evaluator/ judge

I'm 25.



I don't know what you are
talking about



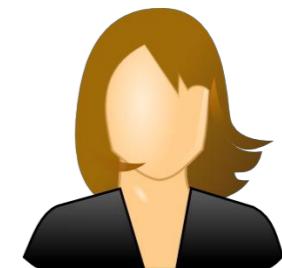
Reward for Dialogue

How old are
you ?



P= 90% human generated

I'm 25.



I don't know what you are
talking about

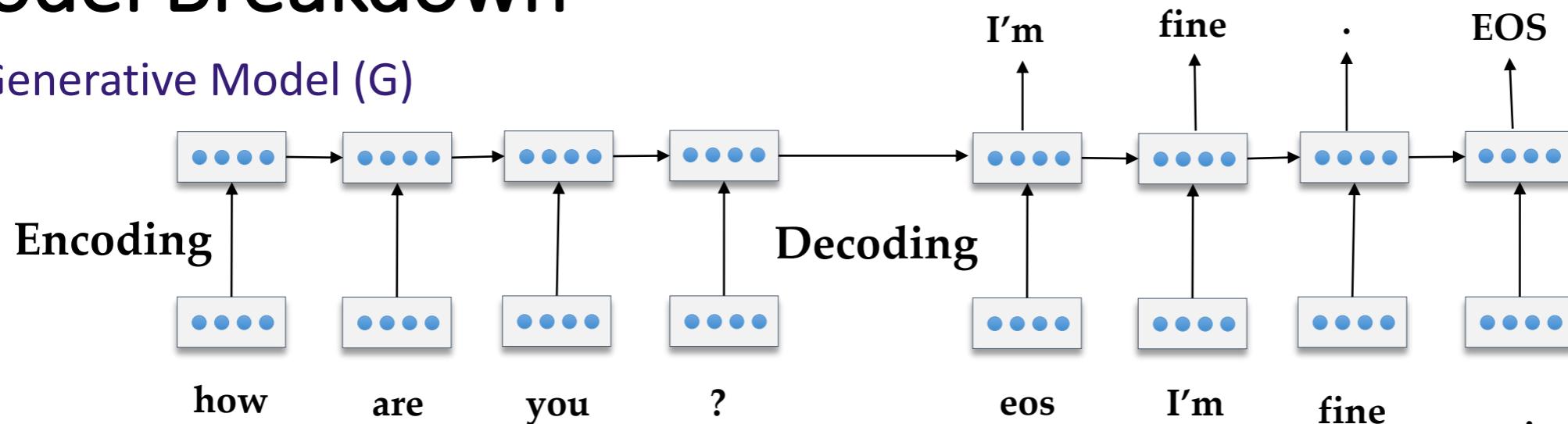
P= 10% human generated



Reward for Dialogue

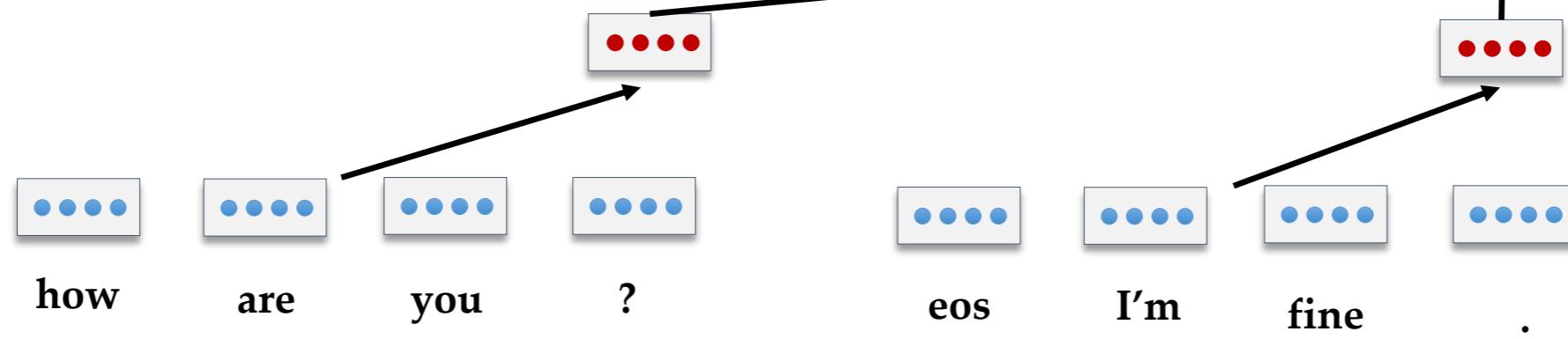
Model Breakdown

Generative Model (G)

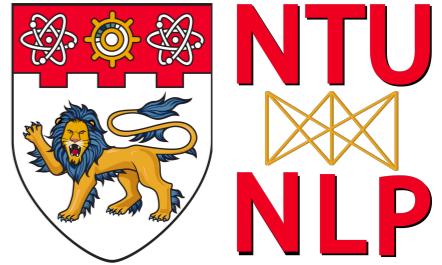


Discriminative Model (D)

Reward
 $P = 90\% \text{ human generated}$

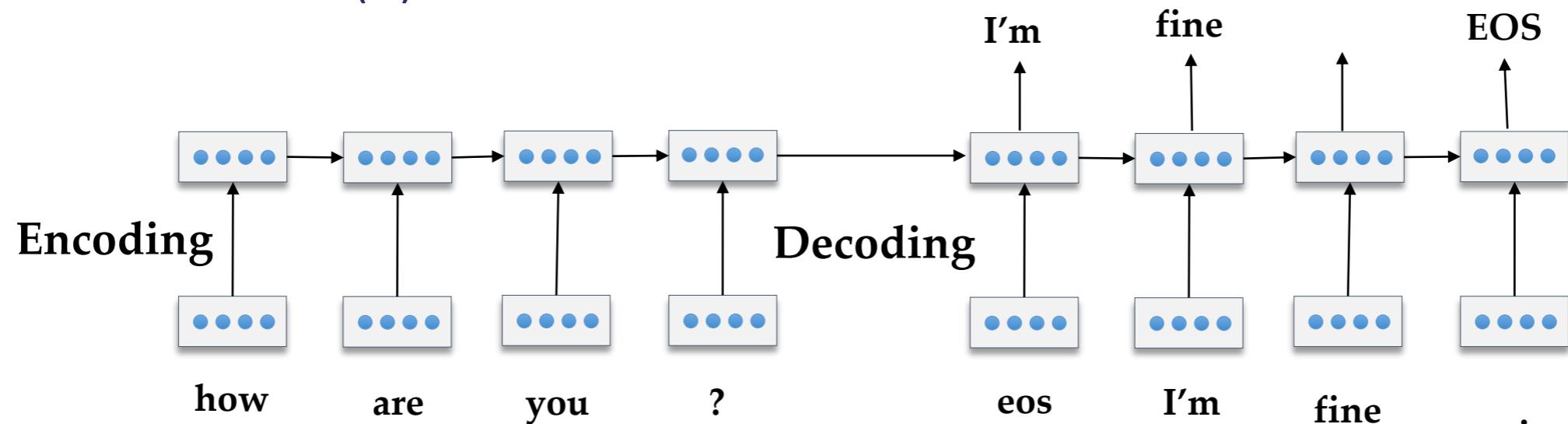


Reward for Dialogue



Policy Gradient

Generative Model (G)



REINFORCE Algorithm (William, 1992)

$$J = E[R(y)]$$

Reward for Dialogue

For number of training iterations **do**

- . **For** $i=1, D$ -steps **do**
- . Sample (X, Y) from real data
- . Sample $\hat{Y} \sim G(\cdot | X)$
- . Update D using (X, Y) as positive examples and (X, \hat{Y}) as negative examples.
- . **End**

Update the Discriminator

For $i=1, G$ -steps **do**

- Sample (X, Y) from real data
- Sample $\hat{Y} \sim G(\cdot | X)$
- Compute Reward r for (X, \hat{Y}) using D .
- Update G on (X, \hat{Y}) using reward r
- Teacher-Forcing: Update G on (X, Y)

Update the Generator

End

End

The discriminator forces the generator to produce correct responses

Reward for Dialogue

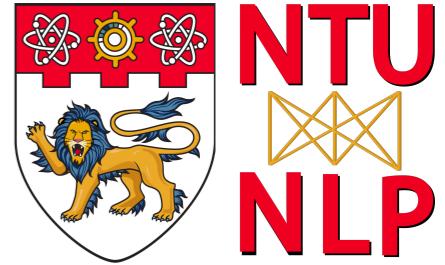


Human Evaluator

vs a vanilla generation model

Adversarial Win	Adversarial Lose	Tie
62%	18%	20%

Adversarial Nets



Adversarial nets have been used in **two ways**:

① **Generative modeling:** GANs [Goodfellow et al., 2014].

- Generate real-like image samples from random vectors
- One neural net is pit against another:
 - (i) Generator, (ii) Discriminator
- Other Examples: image-to-image translation, Unsup. NMT

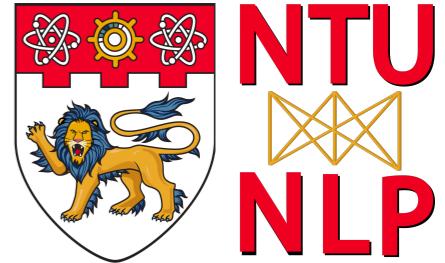


② **Transfer learning/domain adaptation:**

DANNs [Ganin et al., 2016]

- Map samples from two domains into a common feature space
- Generally a **three-player** game:
 - (i) Encoder, (ii) Classifier (iii) Discriminator
- Examples: MNIST \Rightarrow USPS, X-lingual NER

Domain Adversarial Nets



Training Data:

$$\begin{aligned}\mathcal{D}_s &= \{(\mathbf{x}_i, y_i)\}_{i=1}^N \\ \mathcal{D}_t &= \{\mathbf{x}_i\}_{i=1}^M\end{aligned}$$

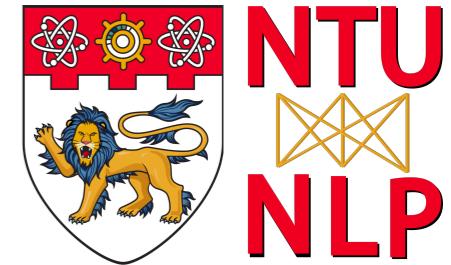
- \mathcal{D}_s denotes source (labeled) dataset.
- \mathcal{D}_t denotes target (unlabeled) dataset.
- Distribution (domain) shift from \mathcal{D}_s to \mathcal{D}_t

Examples

NLP Sentiment analysis (Movie \Rightarrow Book), NER (Eng \Rightarrow Ger)

CV Image classification: MNIST \Rightarrow USPS

Domain Adversarial Nets

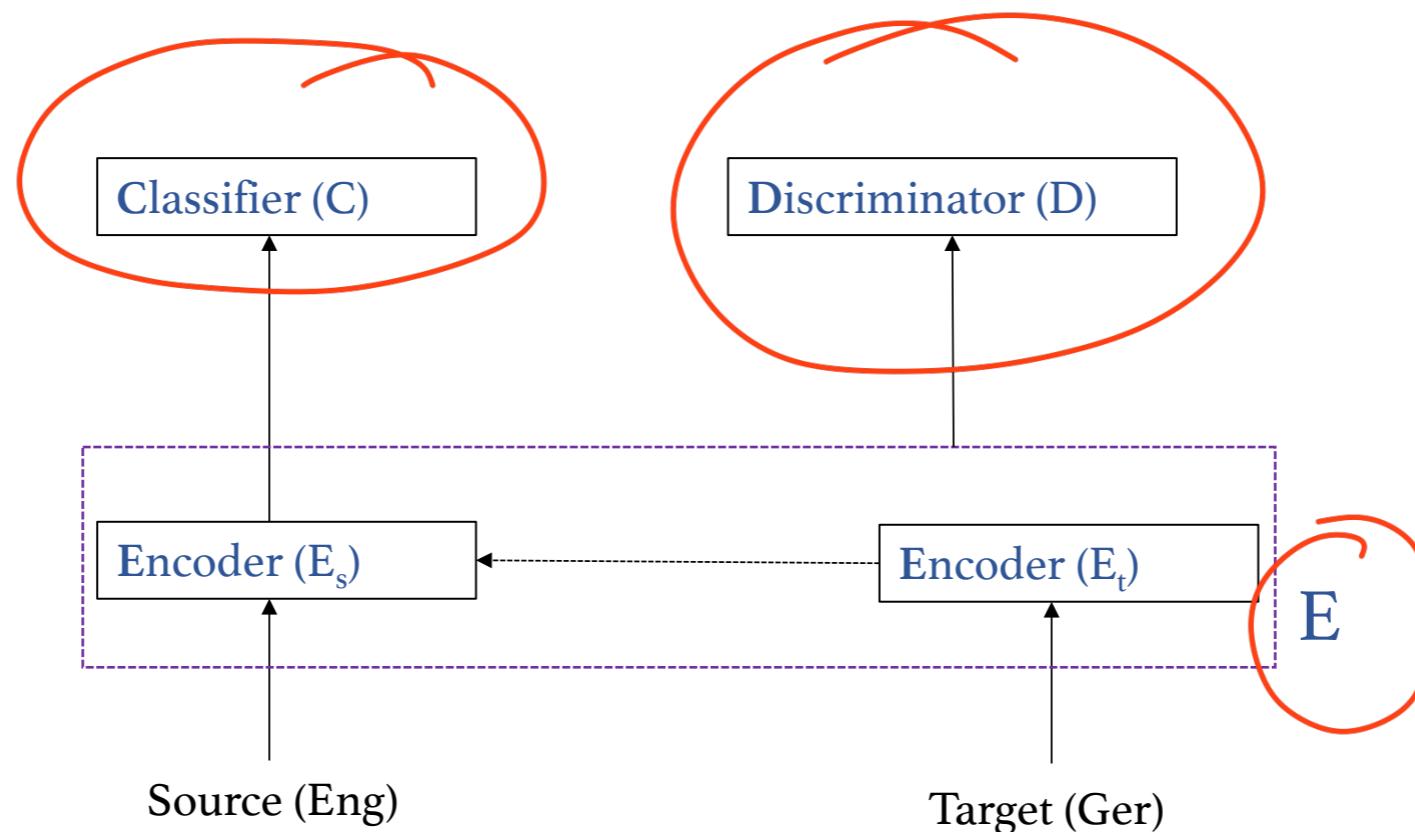


DANN Approach: Learn a common feature space for \mathcal{D}_s and \mathcal{D}_t

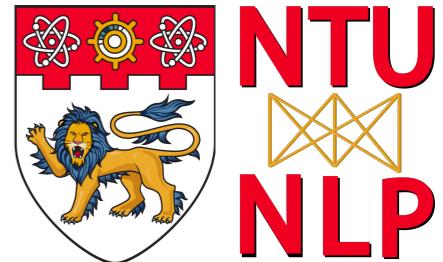
DANN Framework

Three players

(i) Encoder; (ii) Classifier; (iii) Discriminator



Domain Adversarial Nets



Design Choices

Three players

- (i) Encoder; (ii) Classifier; (iii) Discriminator

Encoder

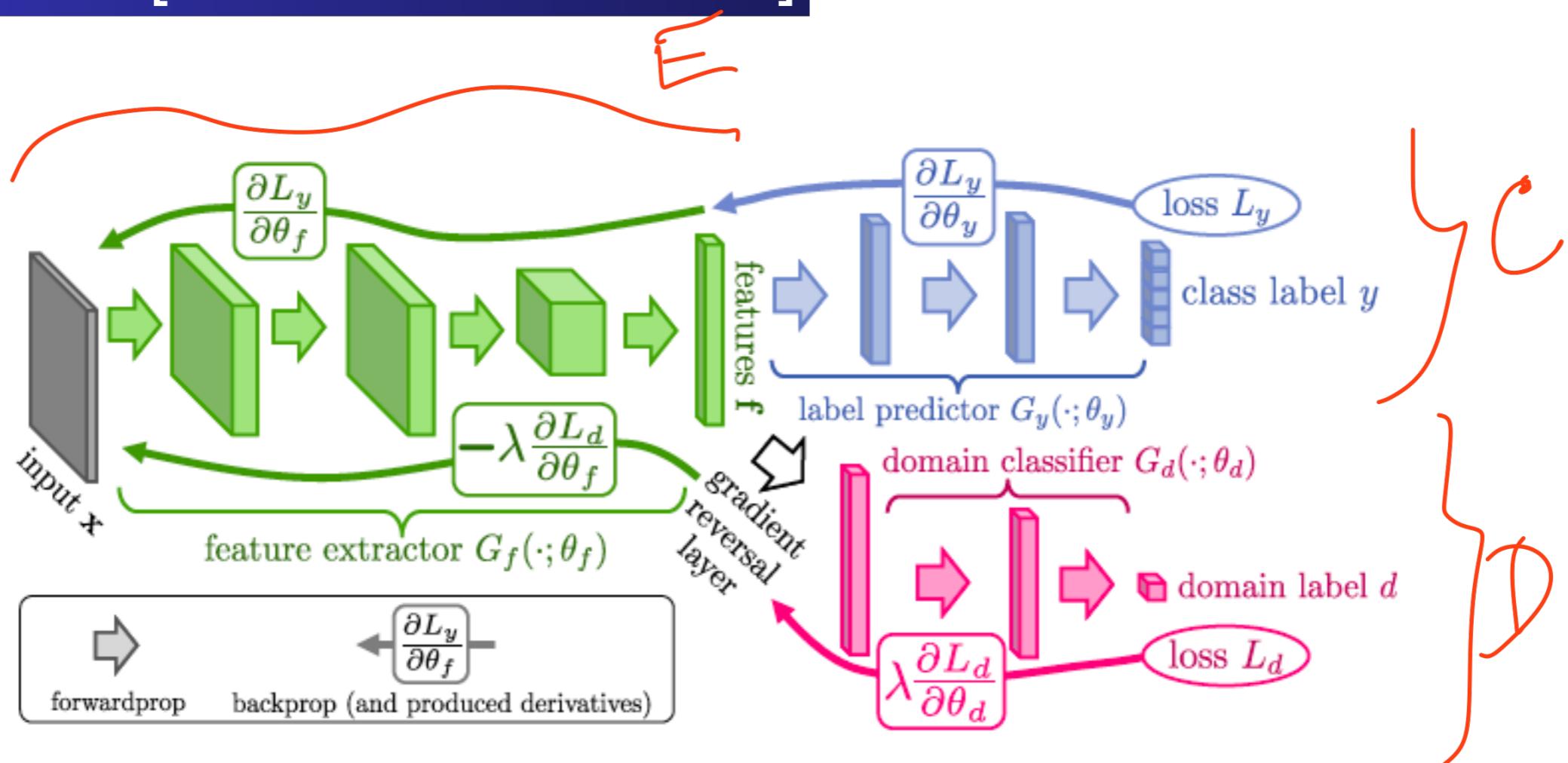
- Shared [Ganin et al., 2016] vs. Separate [Tzeng et al., 2017]
- Train Encoders concurrently or in steps.
- Generative [Volpi et al., 2017] vs. Discriminative [Tzeng et al., 2017]

Adversary Loss

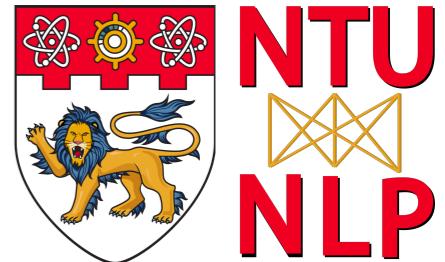
- Flip gradient [Ganin et al., 2016]
- GAN loss [Tzeng et al., 2017]

Domain Adversarial Nets

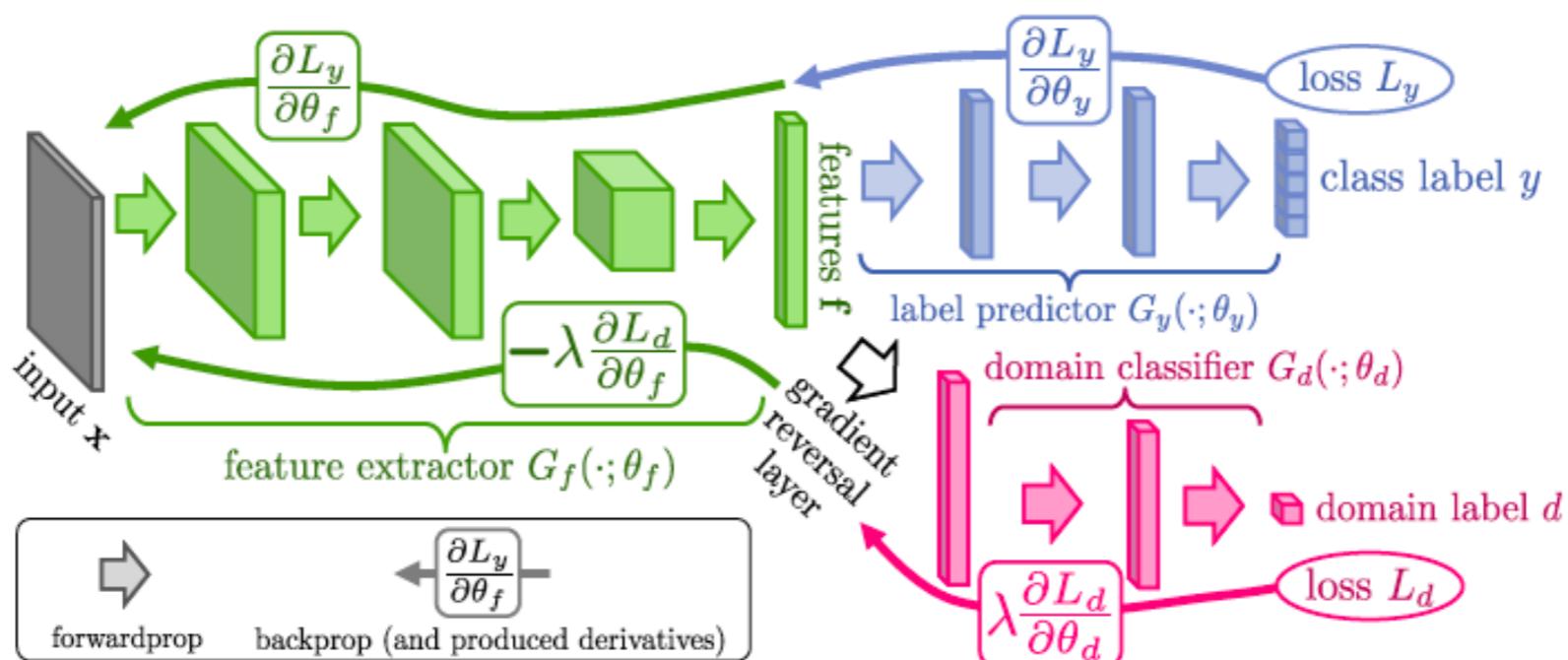
DANN [Ganin et al., 2016]



Domain Adversarial Nets



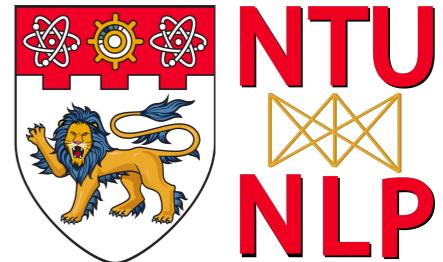
DANN [Ganin et al., 2016]



Choices:

- Shared **encoder**. Train source/target concurrently.
 - **Discriminatively** trained source.
 - **Flip gradient** for **adversary**.

Domain Adversarial Nets



DANN [Ganin et al., 2016]

- **Discriminator:**

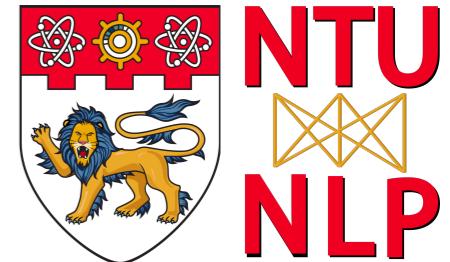
$$-\mathbb{E}_{x_s \sim X_s} \log D(E(x_s)) - \mathbb{E}_{x_t \sim X_t} \log (1 - D(E(x_t))).$$

- **Flip Gradient Adversary (to Encoder):**

⇒ **Encoder (for source):** $\nabla_{\theta_e} \mathbb{E}_{x_s \sim X_s} \log D(E(x_s))$

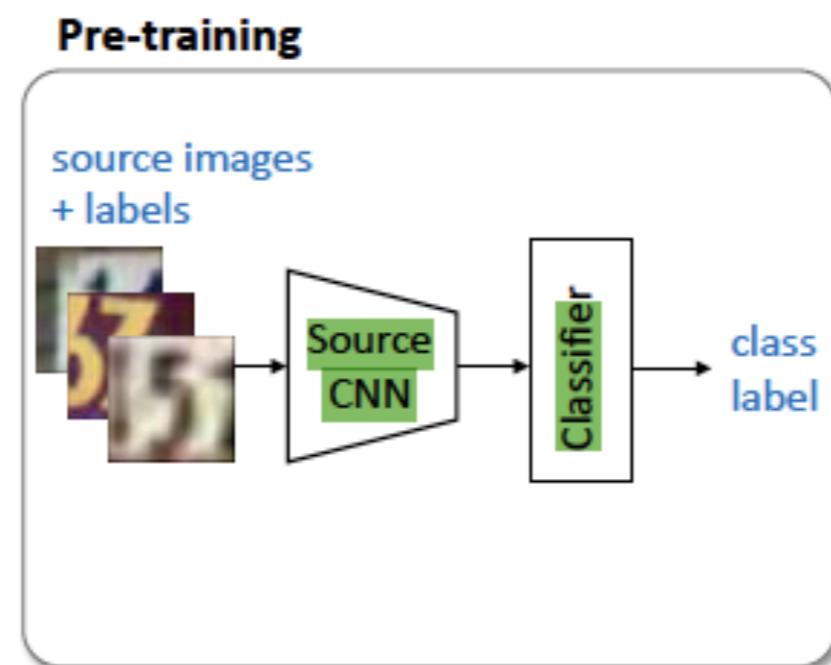
⇒ **Encoder (for target):** $\nabla_{\theta_e} \mathbb{E}_{x_t \sim X_t} \log (1 - D(E(x_t)))$

Domain Adversarial Nets

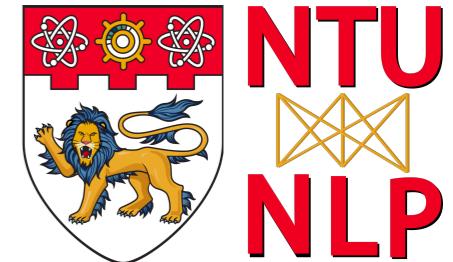


ADDA [Tzeng et al., 2017]

Step 1: Pre-training on Source

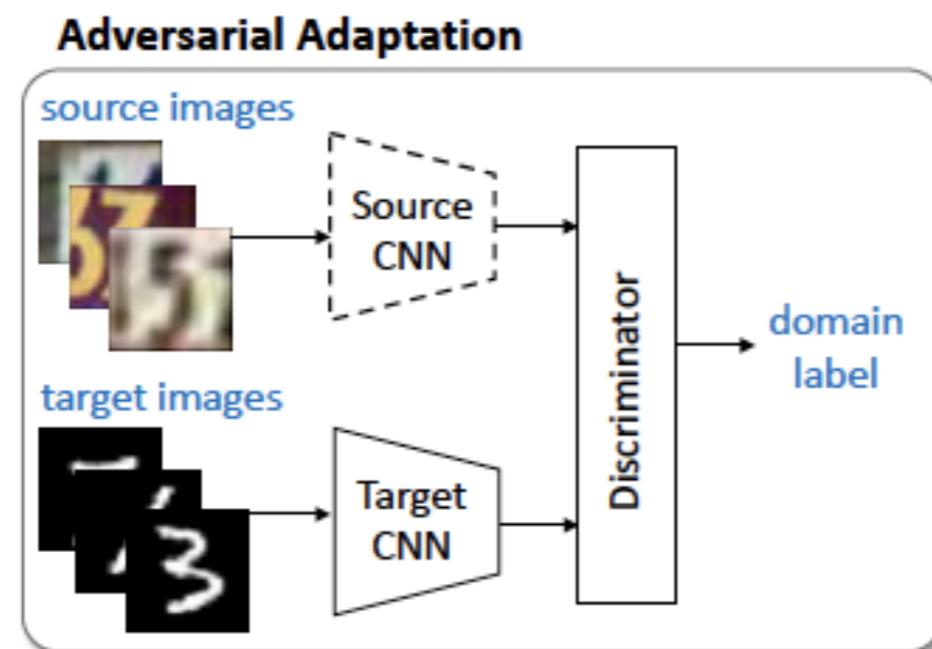


Domain Adversarial Nets

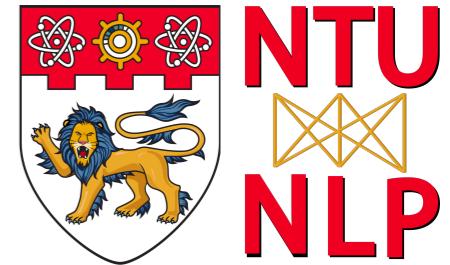


ADDA [Tzeng et al., 2017]

Step 2: Adapt Target towards (fixed) Source

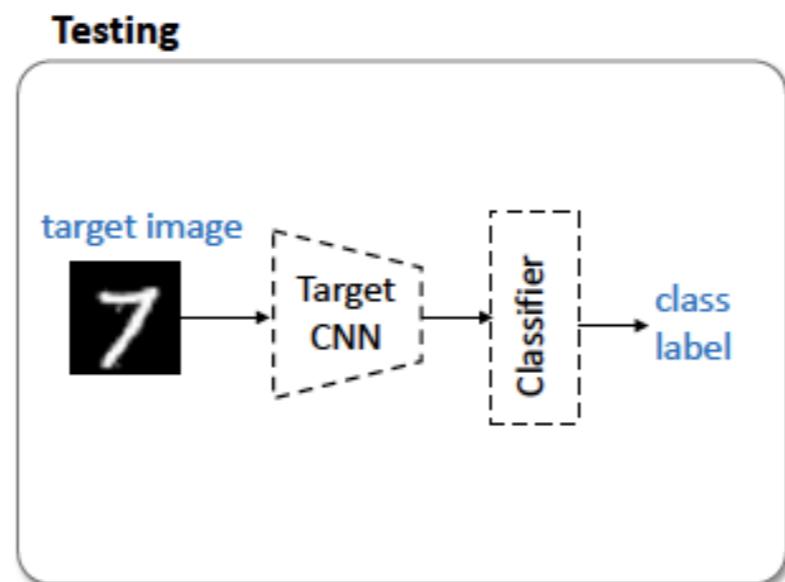


Domain Adversarial Nets

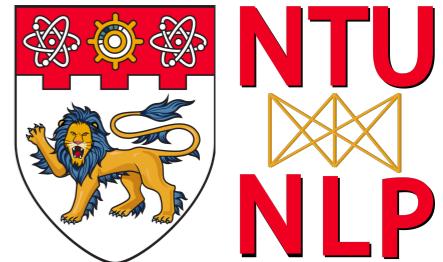


ADDA [Tzeng et al., 2017]

Step 3: Testing



Domain Adversarial Nets

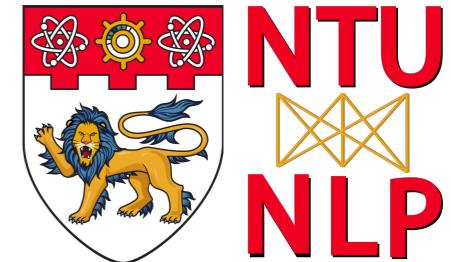


ADDA [Tzeng et al., 2017]

ADDA Design Choices:

- Separate **encoder** for source and target.
- Train source/target separately.
- **Discriminatively** trained source.
- **GAN loss** for **adversary**.

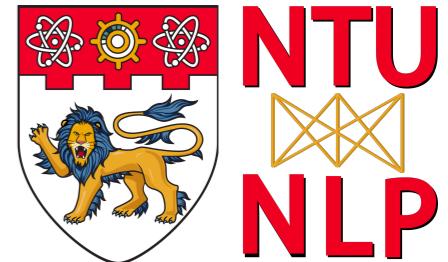
Domain Adversarial Nets



ADDA [Tzeng et al., 2017]

- **Discriminator:**
– $\mathbb{E}_{x_s \sim X_s} \log D(E(x_s)) - \mathbb{E}_{x_t \sim X_t} \log (1 - D(E(x_t)))$.
- **GAN Adversary:**
⇒ **Encoder (for target):** - $\nabla_{\theta_e} \mathbb{E}_{x_t \sim X_t} \log (D(E(x_t)))$
- Recall Adversary for Flip Gradient [Ganin et al., 2016]
⇒ **Encoder (for source):** $\nabla_{\theta_e} \mathbb{E}_{x_s \sim X_s} \log D(E(x_s))$
⇒ **Encoder (for target):** $\nabla_{\theta_e} \mathbb{E}_{x_t \sim X_t} \log (1 - D(E(x_t)))$

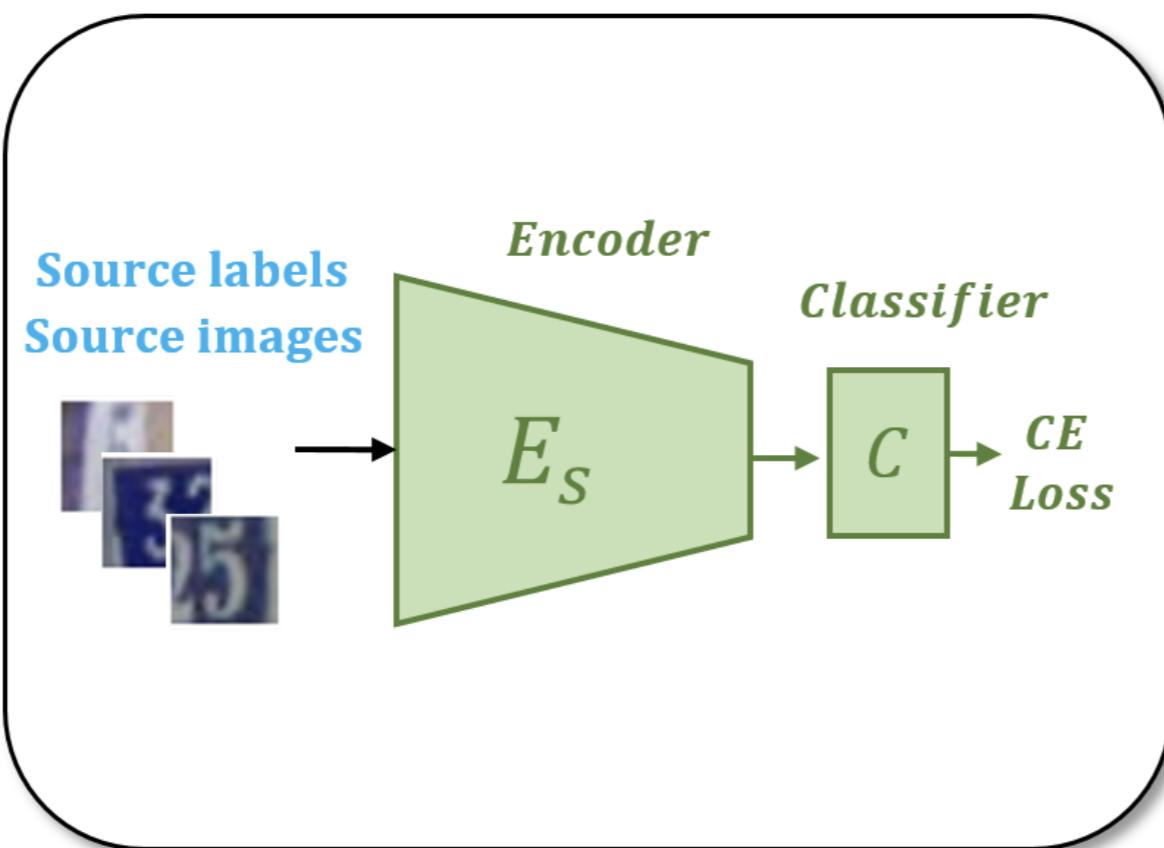
Domain Adversarial Nets



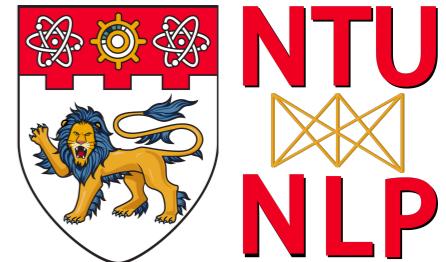
Feature Augmented DANN [Volpi et al., 2017]

Step 0: Pre-training on Source

Step 0: training E_s, C



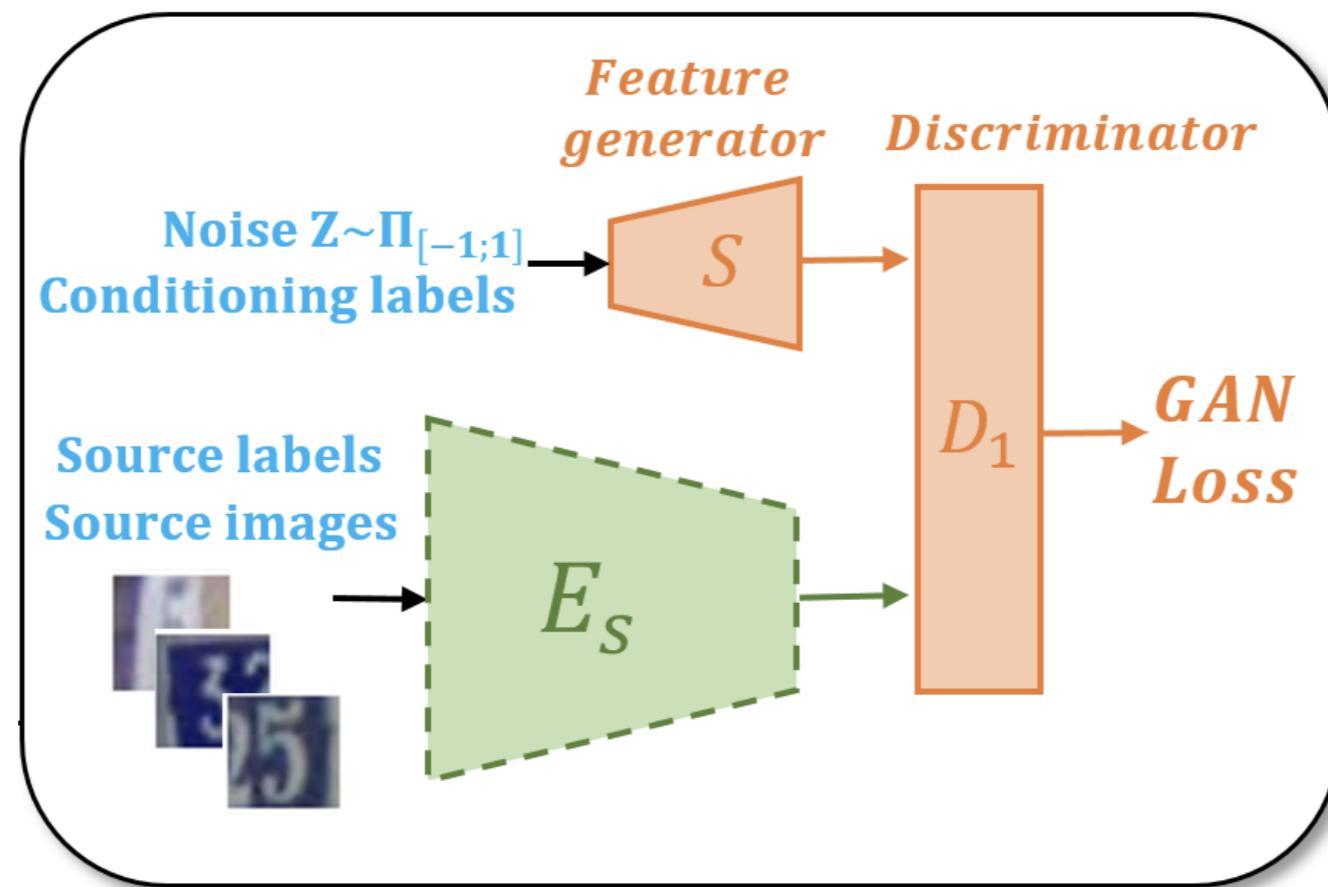
Domain Adversarial Nets



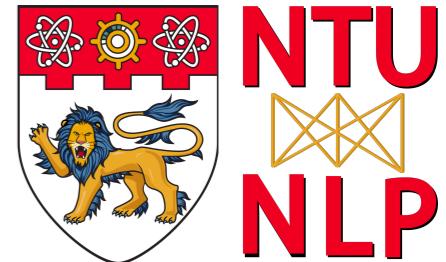
Feature Augmented DANN [Volpi et al., 2017]

Step 1: Feature Generation with CGAN

Step 1: training S, D_1



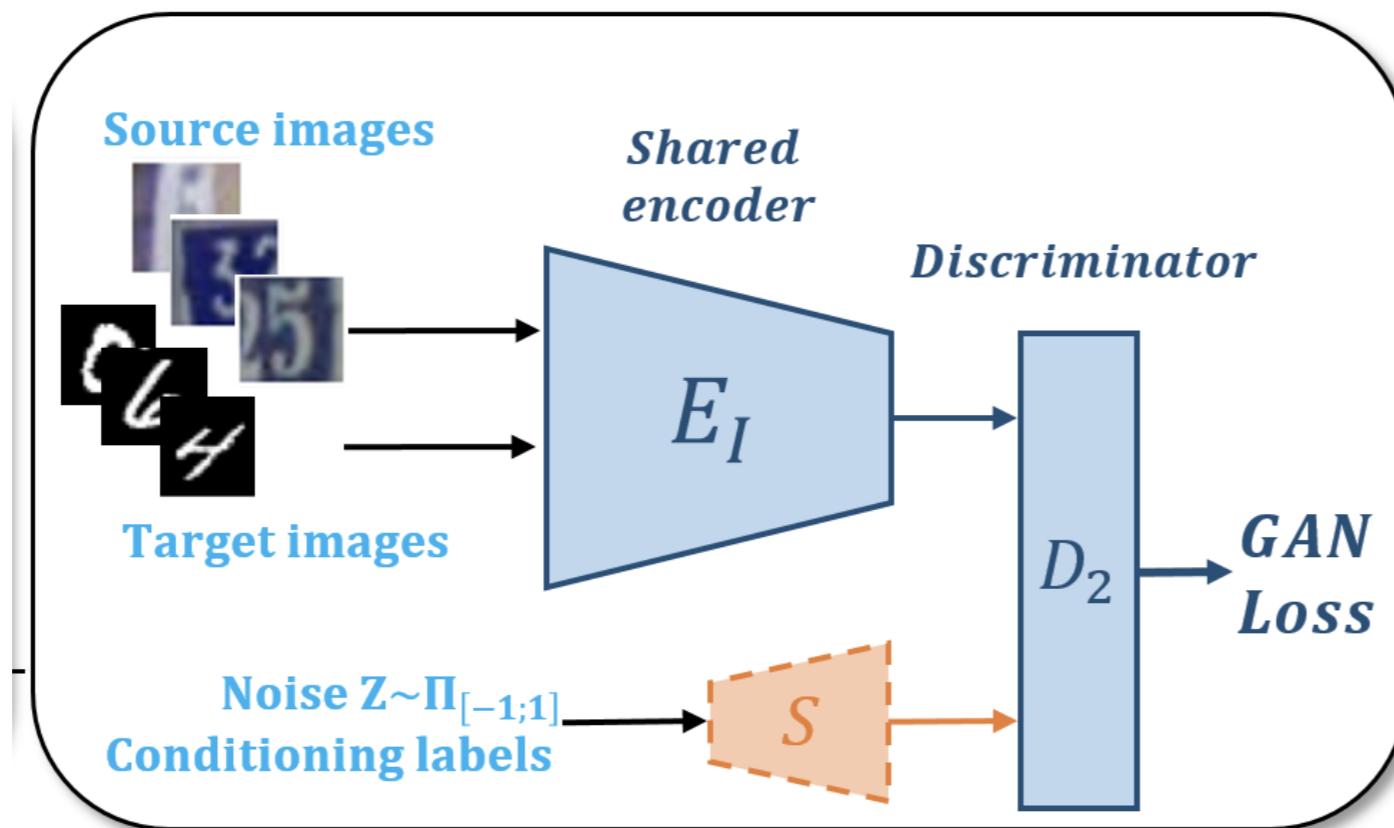
Domain Adversarial Nets



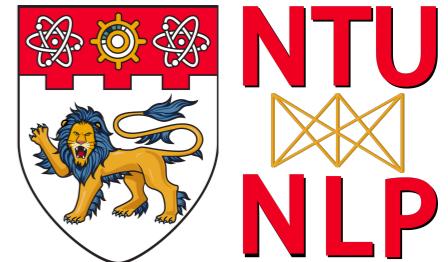
Feature Augmented DANN [Volpi et al., 2017]

Step 2: Feature Adaptation

Step 2: training E_I, D_2

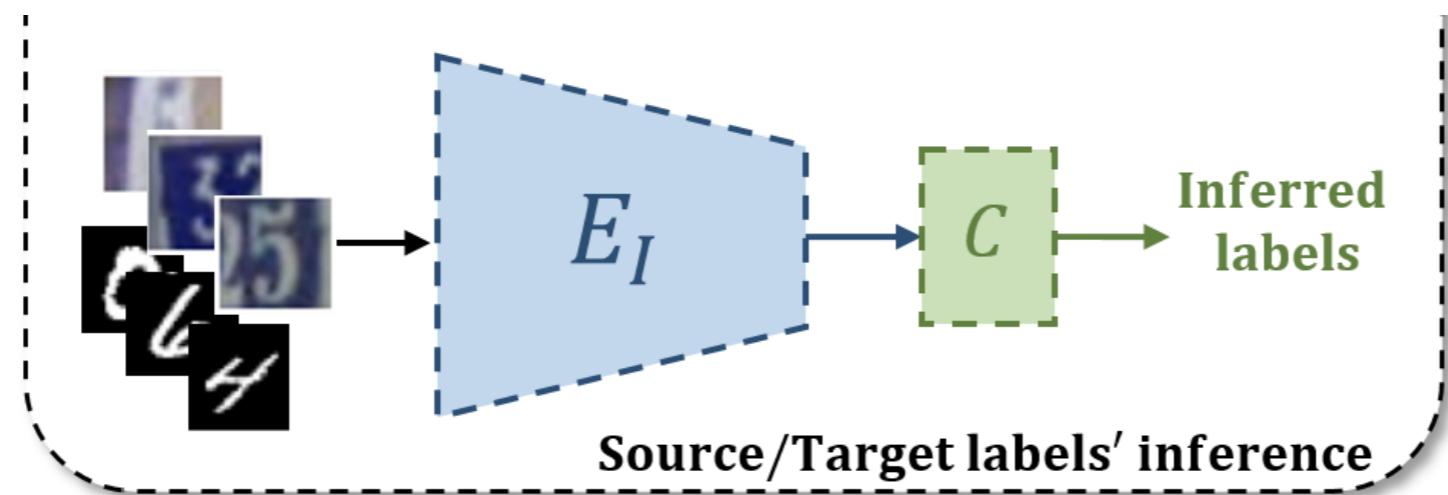


Domain Adversarial Nets

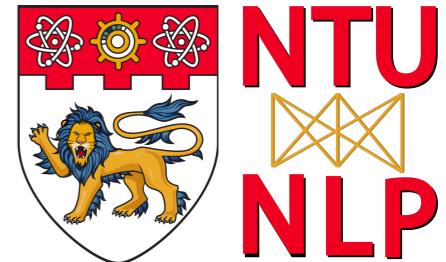


Feature Augmented DANN [Volpi et al., 2017]

Step 3: Testing



Domain Adversarial Nets

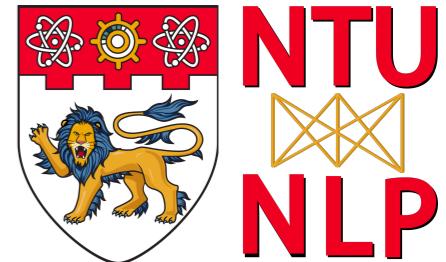


Feature Augmented DANN [Volpi et al., 2017]

FA-DANN Design Choices:

- Same shared **encoder** for source and target.
- Train source/target separately.
- **Discriminatively** trained source + **Generative** features.
- **GAN loss** for **adversary**.

Domain Adversarial Nets



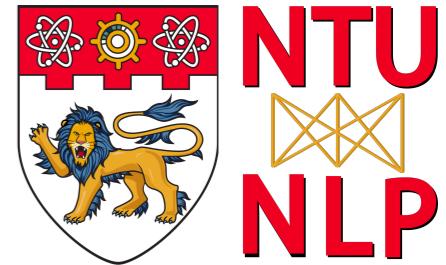
Results

- ADDA [Tzeng et al., 2017]

Method	MNIST → USPS	USPS → MNIST	SVHN → MNIST
	173 → 105	105 → 173	143 51 → 173
Source only	0.752 ± 0.016	0.571 ± 0.017	0.601 ± 0.011
Gradient reversal	0.771 ± 0.018	0.730 ± 0.020	0.739 [19]
Domain confusion	0.791 ± 0.005	0.665 ± 0.033	0.681 ± 0.003
CoGAN	0.912 ± 0.008	0.891 ± 0.008	did not converge
ADDA (Ours)	0.894 ± 0.002	0.901 ± 0.008	0.760 ± 0.018

Table 2: Experimental results on unsupervised adaptation among MNIST, USPS, and SVHN.

Domain Adversarial Nets

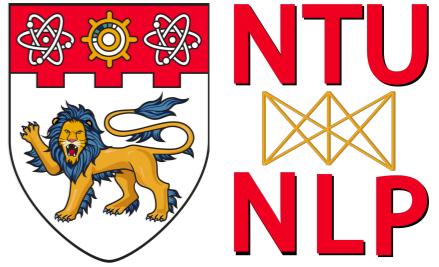


Results

- FA-DANN [Volpi et al., 2017]

	SVHN→MNIST	MNIST→USPS _{P1}	MNIST→USPS _{P2}	USPS→MNIST	SYN→SVHN	NYUD
Source	0.682	0.723	0.797	0.627	0.885	0.139
DANN [9, 10]	0.739	0.771 ± 0.018 [35]	-	0.730 ± 0.020 [35]	0.911	-
DDC [35]	0.681 ± 0.003	0.791 ± 0.005	-	0.665 ± 0.033	-	-
DSN [3]	0.827	-	-	-	0.912	-
ADDA [35]	0.760 ± 0.018	0.894 ± 0.002	-	0.901 ± 0.008	-	0.211
Tri [29]	0.862	-	-	-	0.931	-
DTN [33]	0.844*	-	-	-	-	-
PixelDA ** [2]	-	-	0.959	-	-	-
UNIT [18]	0.905*	-	0.960	-	-	-
CoGANs [19]	no conv. [35]	0.912 ± 0.008	0.957 [18]	0.891 ± 0.008	-	-
LS-ADDA	0.743 ± 0.028	0.914 ± 0.000	0.912 ± 0.003	0.910 ± 0.004	0.908 ± 0.004	no conv.
<i>Ours (DI)</i>	0.851 ± 0.026	0.914 ± 0.000	0.954 ± 0.002	0.879 ± 0.005	0.925 ± 0.002	0.287 ± 0.002
<i>Ours (DIFA)</i>	0.897 ± 0.020	0.923 ± 0.001	0.962 ± 0.002	0.897 ± 0.005	0.930 ± 0.002	0.313 ± 0.002
Target	0.992	0.999	0.999	0.975	0.913	0.468 [35]

More Examples of Adversarial Learning



Monet \leftrightarrow Photos



Monet \rightarrow photo

Zebras \leftrightarrow Horses



zebra \rightarrow horse

Summer \leftrightarrow Winter



summer \rightarrow winter

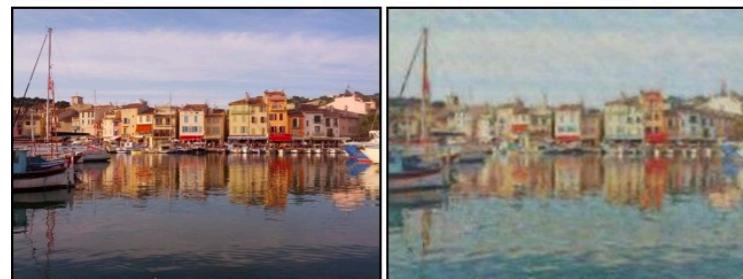


photo \rightarrow Monet



horse \rightarrow zebra



winter \rightarrow summer



Photograph



Monet



Van Gogh

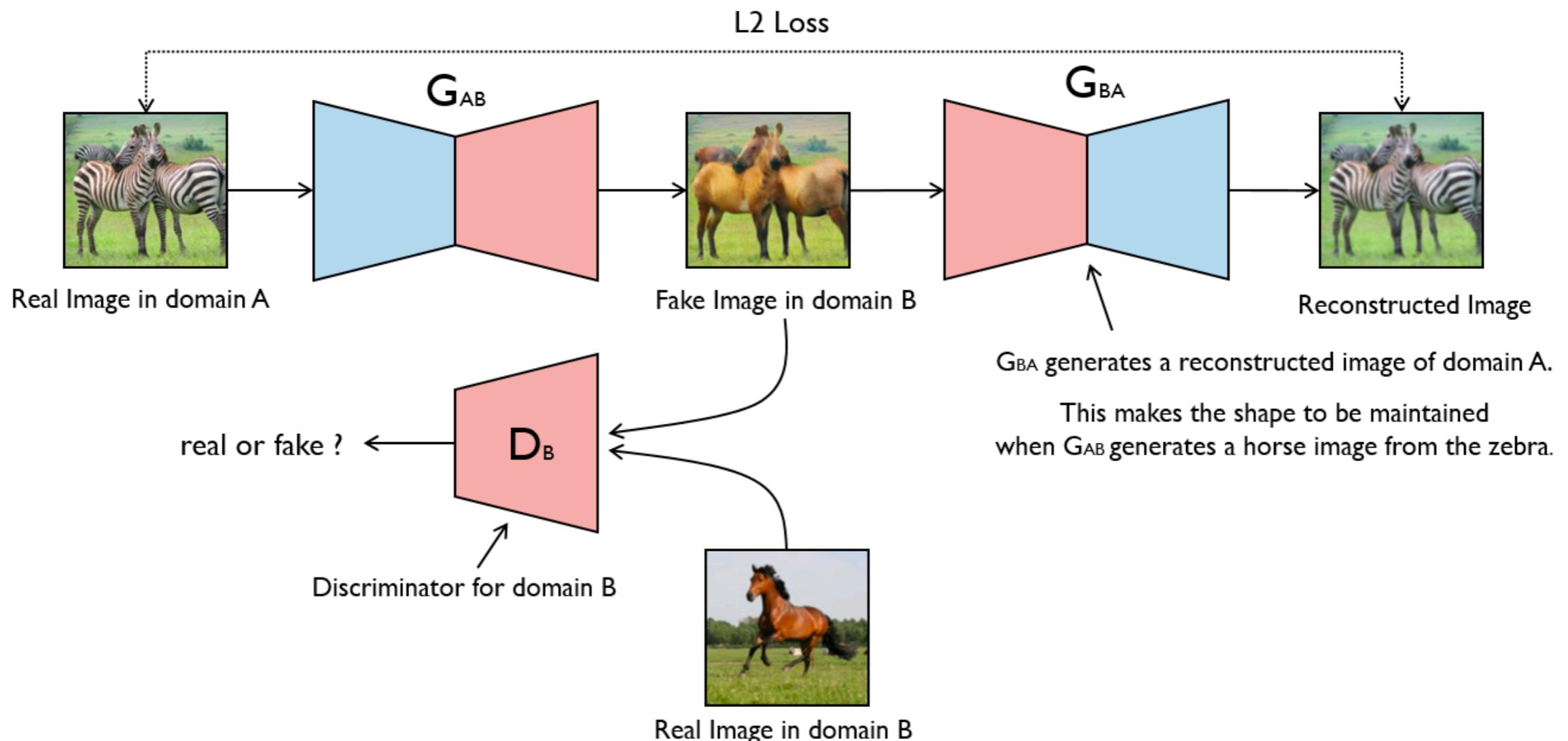
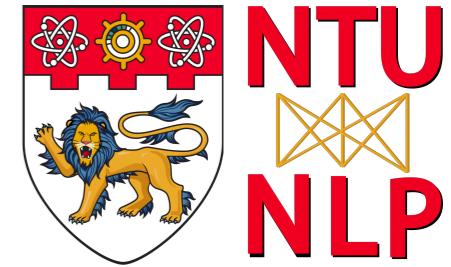


Cezanne

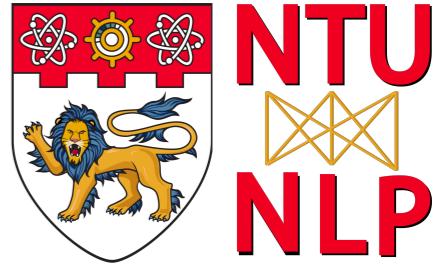


Ukiyo-e

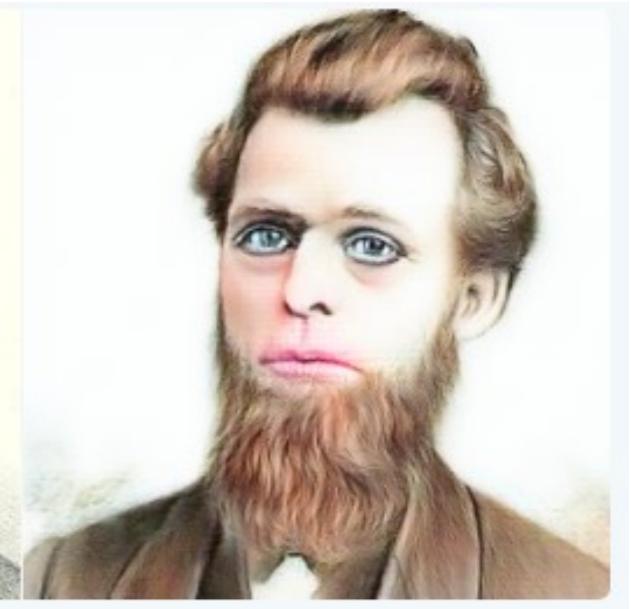
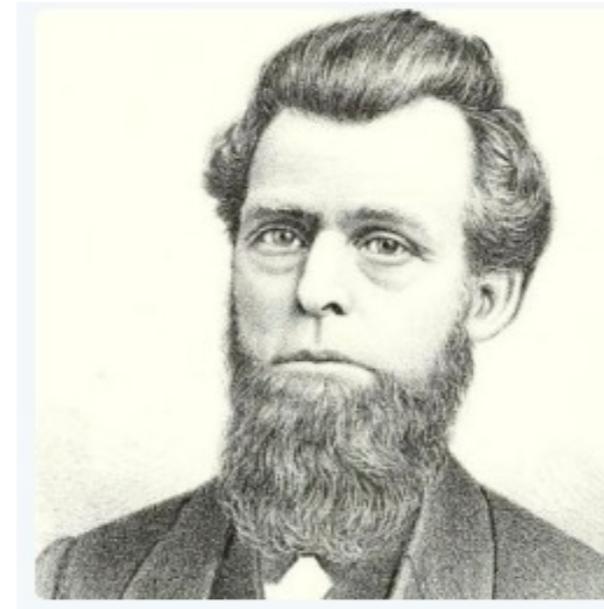
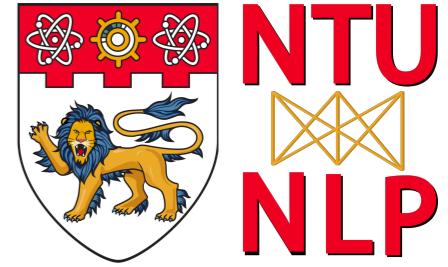
More Examples of Adversarial Learning



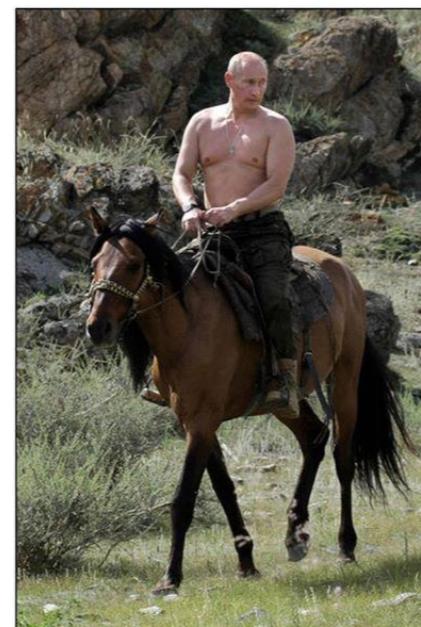
More Examples of Adversarial Learning



More Examples of Adversarial Learning



Failure case



Part 2: Adversarial Attacks

