

NTU 2020 Spring Machine Learning

Self-Supervised Learning

Chi-Liang Liu, Hung-Yi Lee

NTU

Previously on this Course

- Supervised Learning

- Given: a dataset $\mathcal{D} = \{(\mathbf{x}, \mathbf{y})_i\}_{i=1}^N$ and
a loss function $\ell : \hat{\mathcal{Y}} \times \mathcal{Y} \rightarrow \mathbb{R}, (\hat{\mathbf{y}}, \mathbf{y}) \rightarrow \ell(\hat{\mathbf{y}}, \mathbf{y})$

Goal: $\min_{\theta} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [\ell(f_{\theta}(\mathbf{x}), \mathbf{y})]$

- Works well when labeled data is abundant.
- Learn useful representation with the supervision.

- Problem:

Can we learn useful representation without the supervision?

Why Self-Supervised Learning?



Labeled Data



Unlabeled Data

Slide: Thang Luong

Why Self-Supervised Learning?

- ▶ **“Pure” Reinforcement Learning (cherry)**

- ▶ The machine predicts a scalar reward given once in a while.

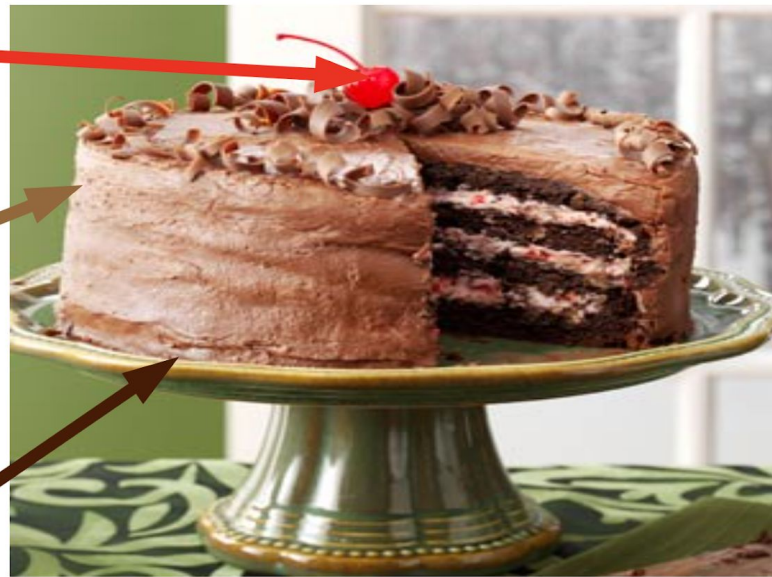
- ▶ **A few bits for some samples**

- ▶ **Supervised Learning (icing)**

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

- ▶ **Self-Supervised Learning (cake génoise)**

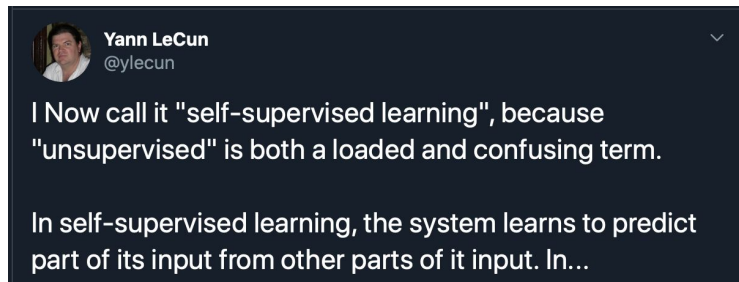
- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**



Yann LeCun's cake

What is Self-Supervised Learning?

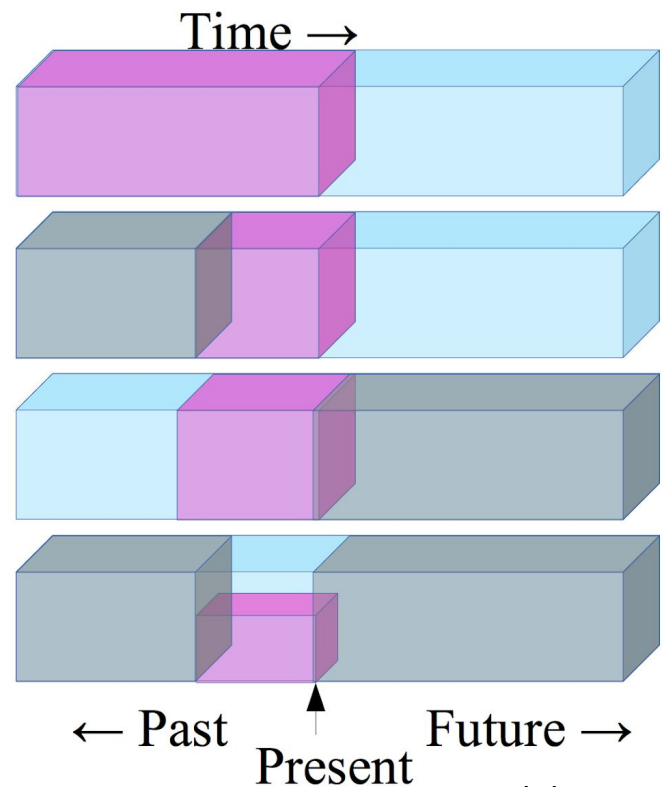
- A version of unsupervised learning where data provides the supervision.



- In general, withhold some part of the data and the task a neural network to predict it from the remaining parts.
- Goal: Learning to represent the world before learning tasks.

Self-Supervised Learning= Filling the Blanks

- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the occluded from the visible
- ▶ **Pretend there is a part of the input you don't know and predict that.**

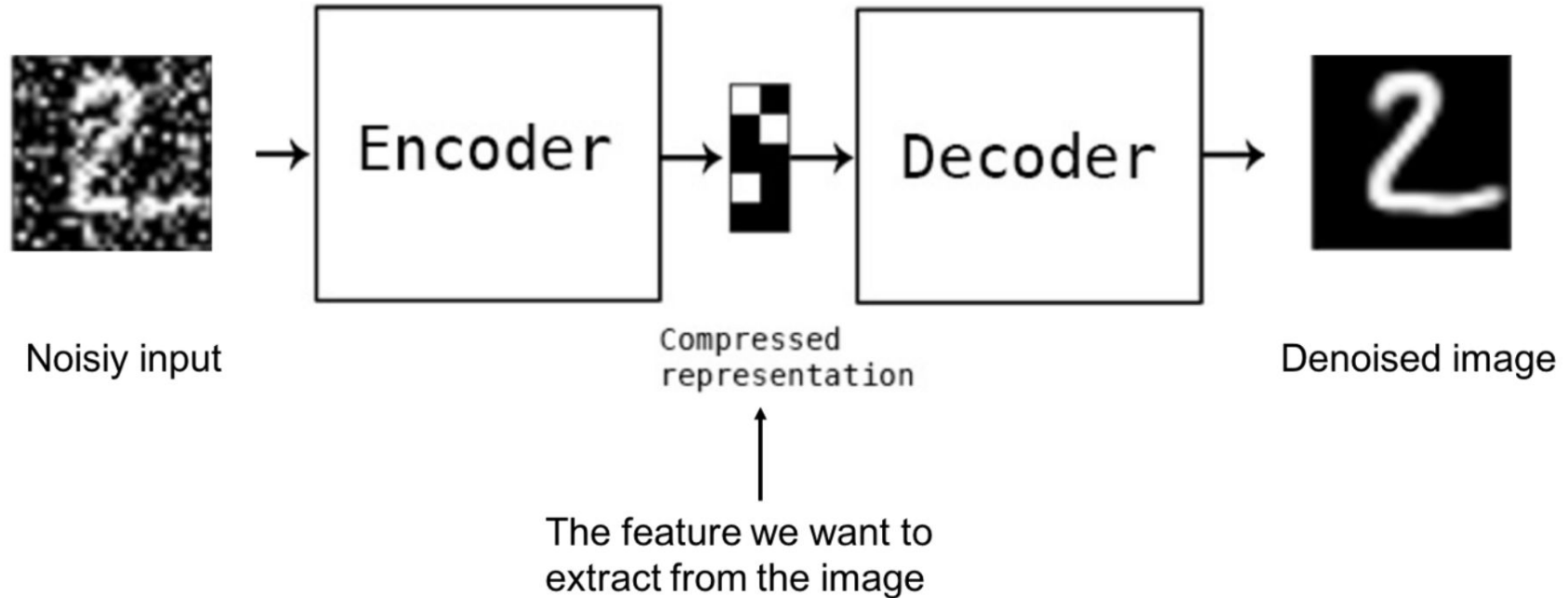


Slide: LeCun

Methods of Self-Supervised Learning

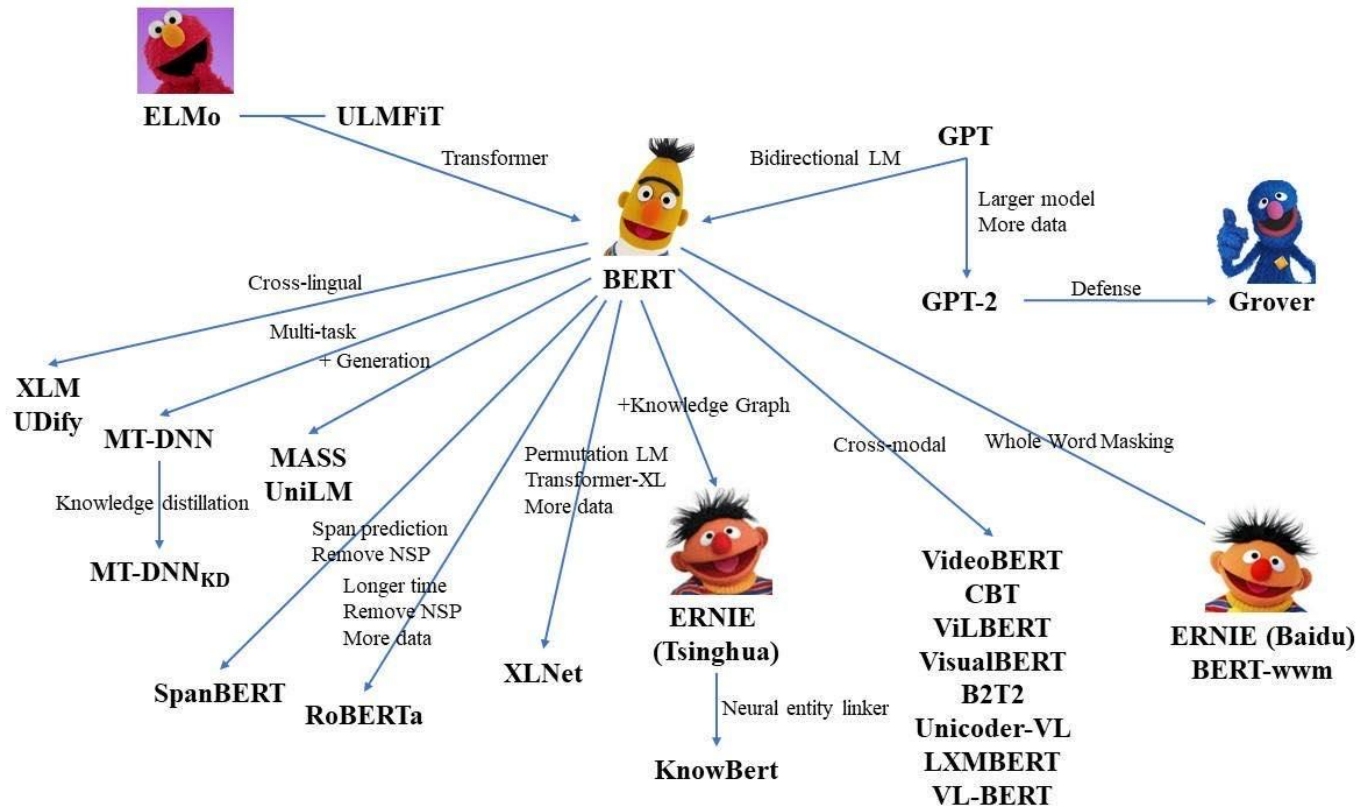
- Reconstruct from a corrupted (or partial) data
 - Denoising Autoencoder
 - Bert-Family (Text)
 - In-painting (Image)
- Visual common sense tasks
 - Jigsaw puzzles
 - Rotation
- Contrastive Learning
 - word2vec
 - Contrastive Predictive Coding (CPC)
 - SimCLR

Denoising AutoEncoder



Slide: CS294-158

BERT-Family



Language Model

- A statistical **language model** is a **probability distribution** over sequences of words. Given such a sequence, say of length m , it assigns a probability $P(w_1, \dots, w_m)$ to the whole sequence.
ex. $P(\text{"Is it raining now?"}) > P(\text{"Is it raining yesterday?"})$

- How to Compute?

- n-gram model

$$P(w_1, \dots, w_m) = \prod_{i=1}^m P(w_i \mid w_1, \dots, w_{i-1}) \approx \prod_{i=1}^m P(w_i \mid w_{i-(n-1)}, \dots, w_{i-1})$$

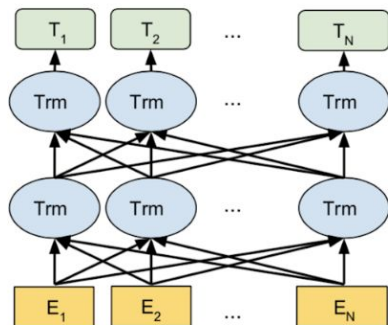
- Neural Network

$$P(w_t \mid \text{context}) \forall t \in V$$

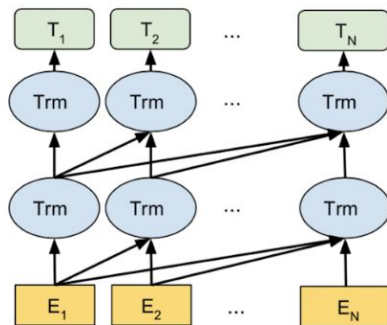
Wiki: Language Model

ELMO & GPT & BERT

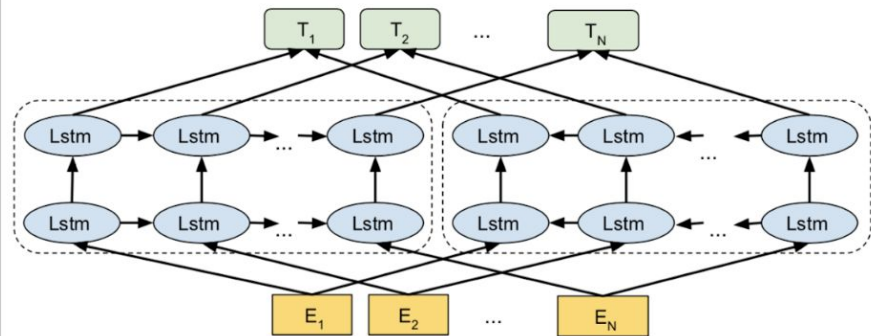
BERT (Ours)



OpenAI GPT

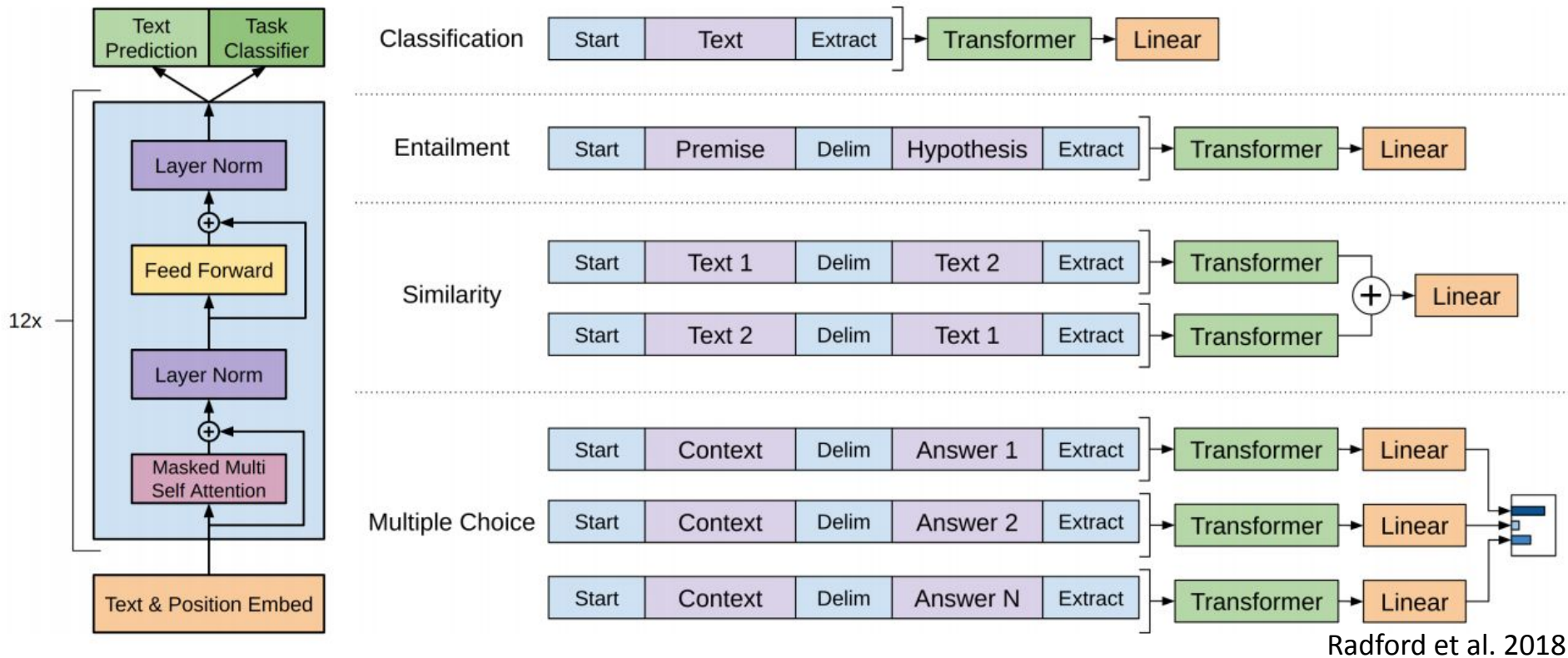


ELMo



Devlin et al. 2018

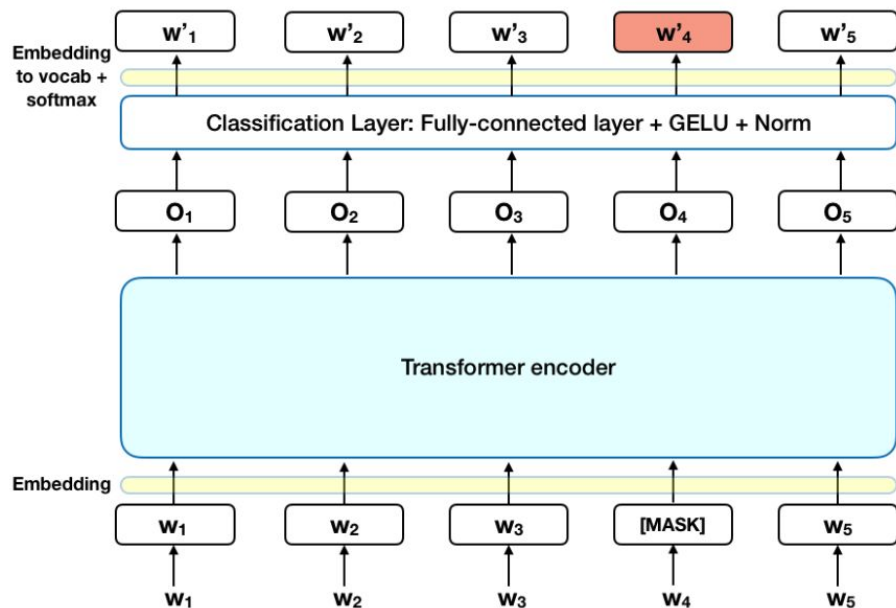
GPT - Language Model



Radford et al. 2018

BERT - Masked LM

Masked LM:



Predict 15% of the tokens in the input.

80% replaced with a masked token

10% replaced with a random word

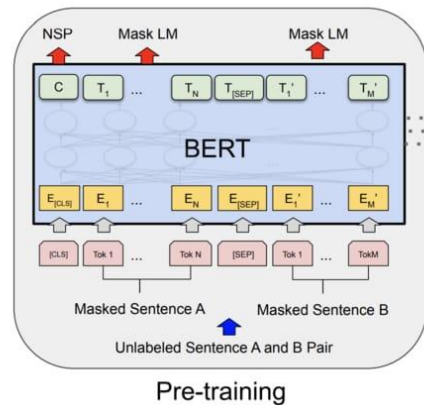
10% remain the same

ex.

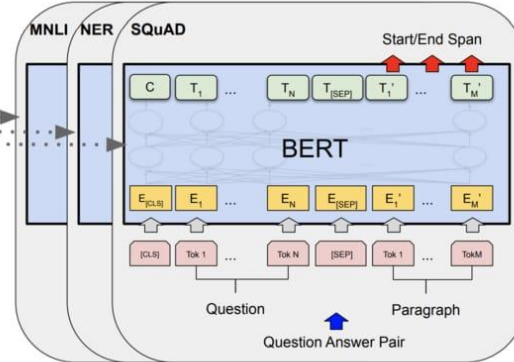
Input: The [mask] sat on the [mask] .

Output: The cat sat on the mat .

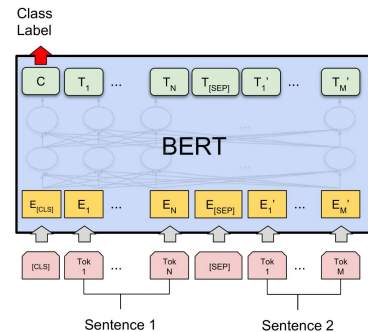
BERT - Pipeline



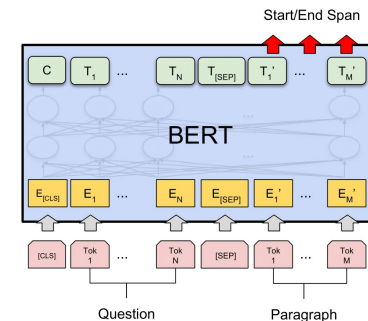
Pre-training



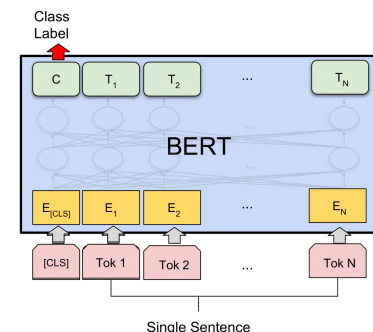
Fine-Tuning



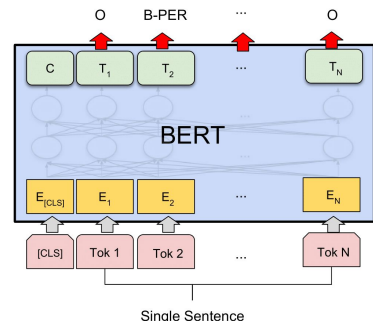
(a) Sentence Pair Classification Tasks:
MNLI, QQP, QNLI, STS-B, MRPC,
RTE, SWAG



(c) Question Answering Tasks:
SQuAD v1.1



(b) Single Sentence Classification Tasks:
SST-2, CoLA



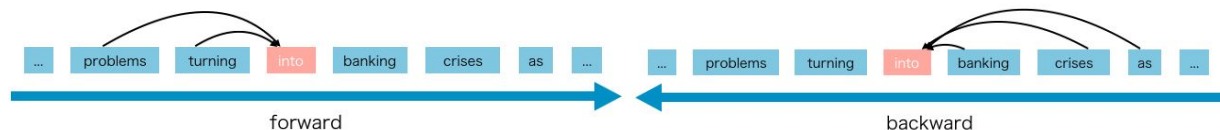
(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

Devlin et al. 2018

ARLM vs AELM

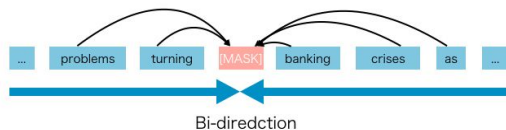
- Autoregressive Language Model (ARLM)

- Pro: Does not rely on data corruption
- Con: Can only be forward or backward



- Autoencoding Language Model (AELM a.k.a MaskedLM)

- Pro: Can use bidirectional information (Context Dependency)
- Con: Pretrain-Finetune discrepancy (Input Noise), Independence Assumption



Yang et al. 2019

XLNet - Permutation LM

Permutation LM:

Step 1. Permutation

Step 2. Autoregressive

ex.

Given Sequence $[x_1, x_2, x_3, x_4]$

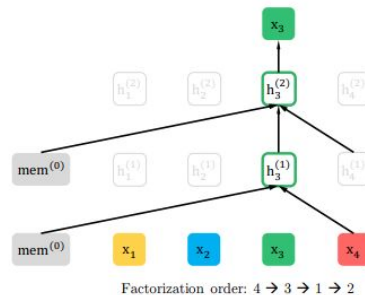
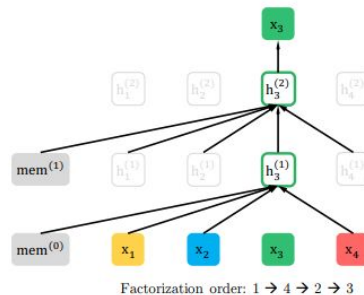
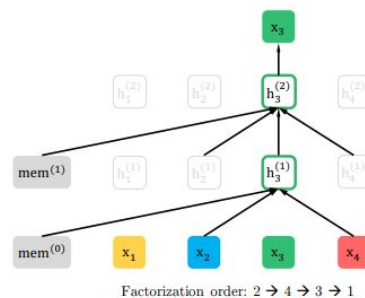
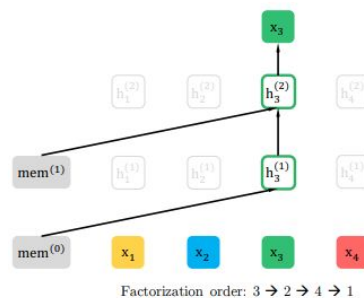
Step 1. $[x_1, x_2, x_3, x_4] \rightarrow [x_2, x_4, x_3, x_1]$

Step 2. Given $[x_2]$ predict $[x_2, x_4]$

-> Given $[x_2, x_4]$ predict $[x_2, x_4, x_3]$

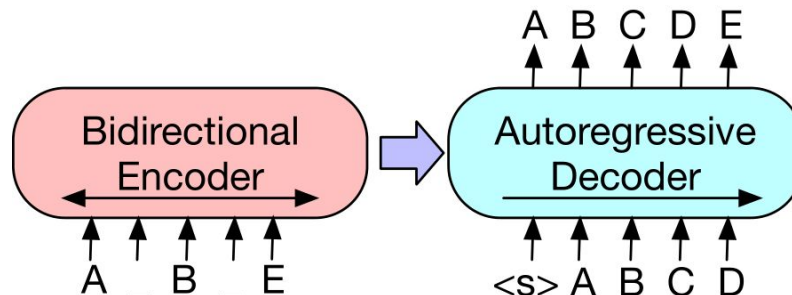
-> Given $[x_2, x_4, x_3]$ predict $[x_2, x_4, x_3, x_1]$

Yang et al. 2019

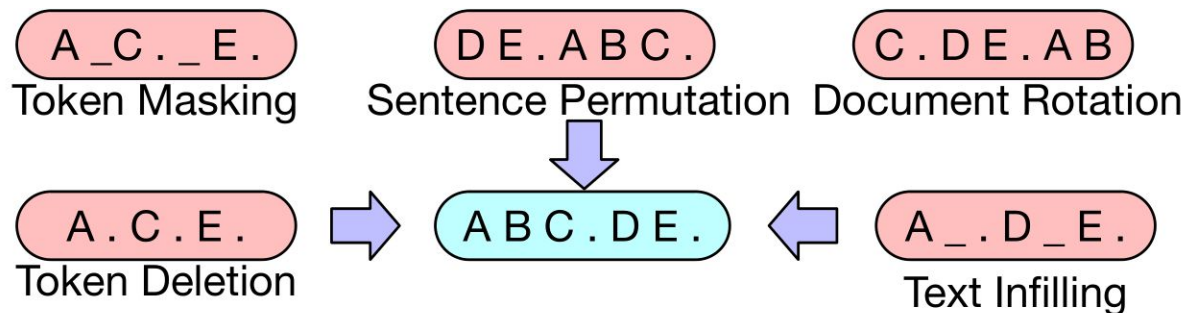


BART - Encoder & Decoder

- Previous: Fixed-length to Fixed-length
- BART:

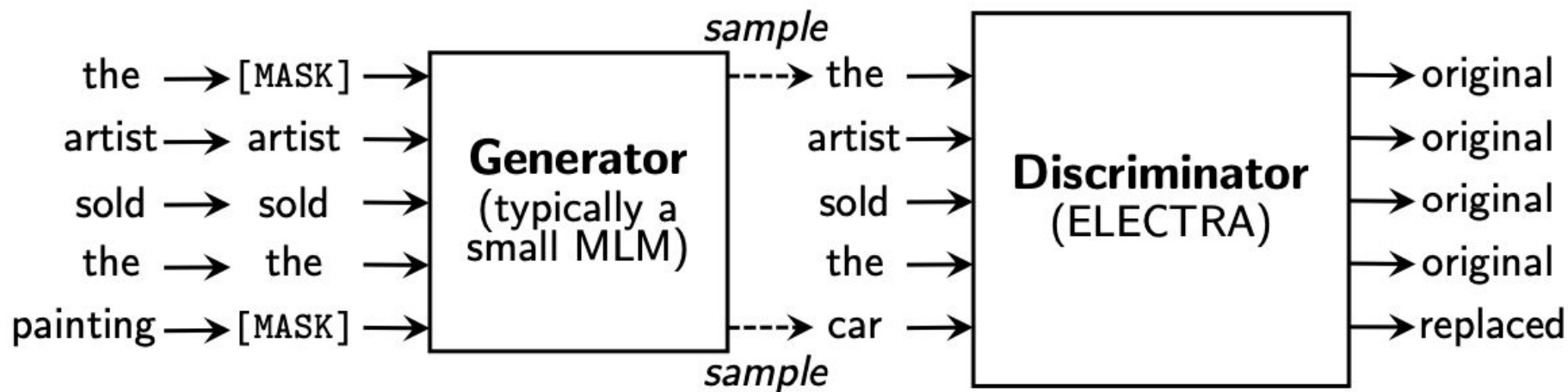


- Tasks:



Lewis et al. 2018

ELECTRA - Discriminator



Clark et al. 2018

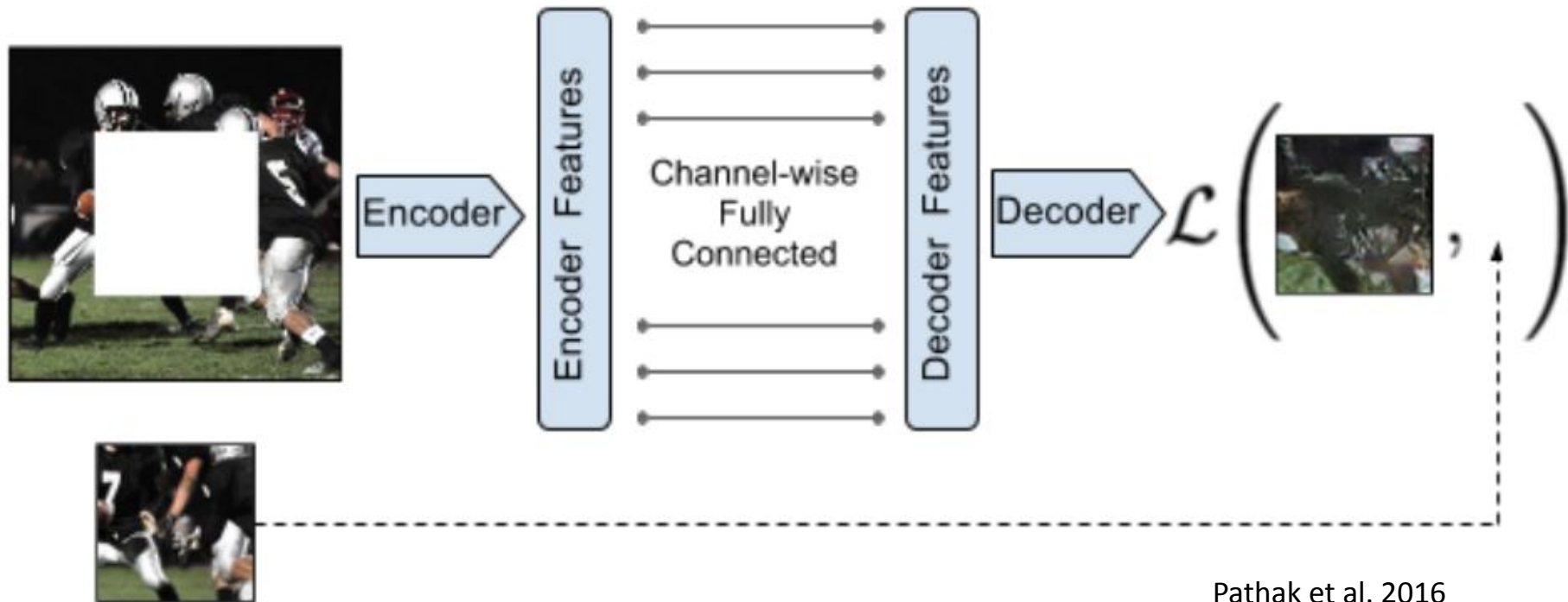
Others

- RoBERTa: A Robustly Optimized BERT Pretraining Approach
- ALBERT: A LITE BERT FOR SELF-SUPERVISED LEARNING OF LANGUAGE REPRESENTATIONS

Predict missing pieces

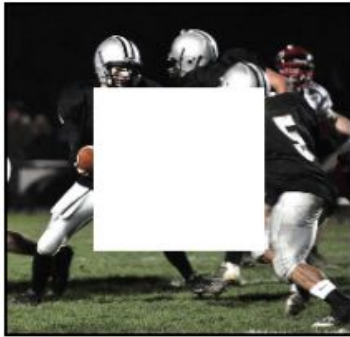


Pathak et al. 2016
Slide: CS294-158

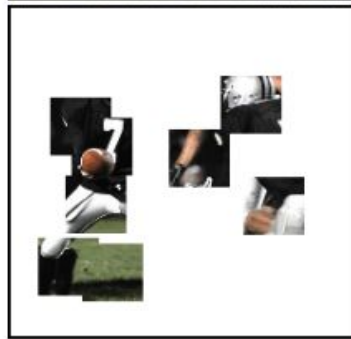
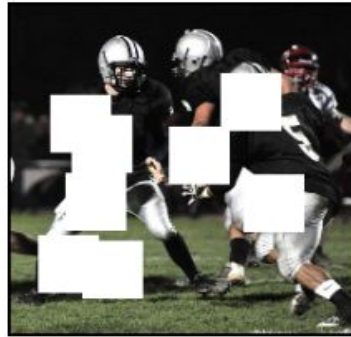


Pathak et al. 2016
Slide: CS294-158

Context Encoder



(a) Central region



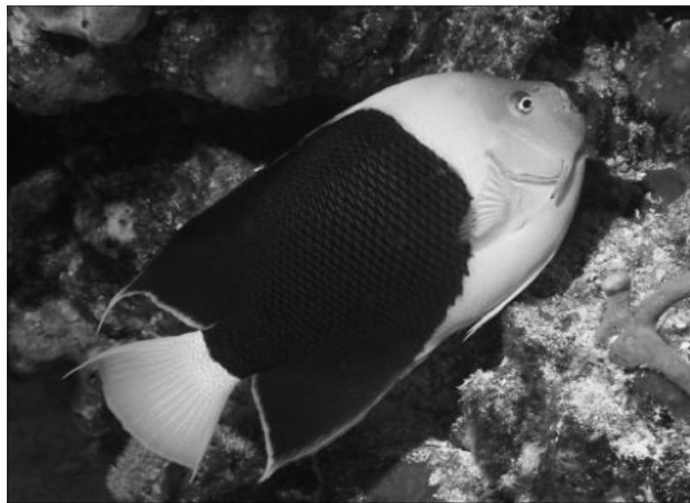
(b) Random block



(c) Random region

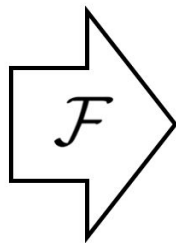
Pathak et al 2016
Slide: CS294-158

Predicting one view from another



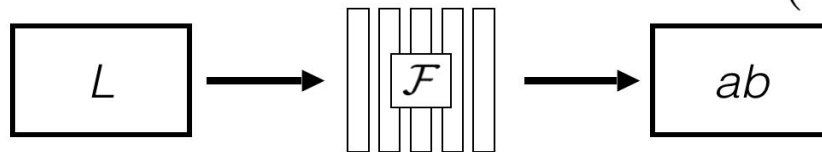
Grayscale image: L channel

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

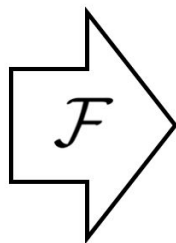
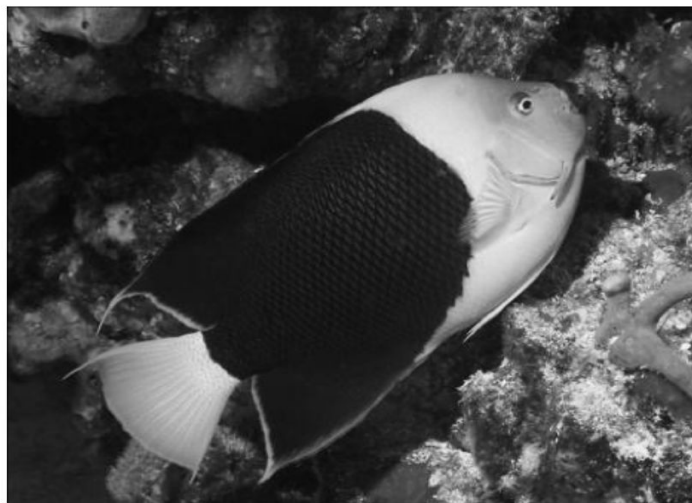


Concatenate (L, ab) channels

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



Slide: Richard Zhang

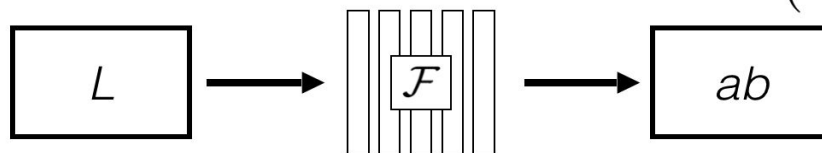


Grayscale image: L channel

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

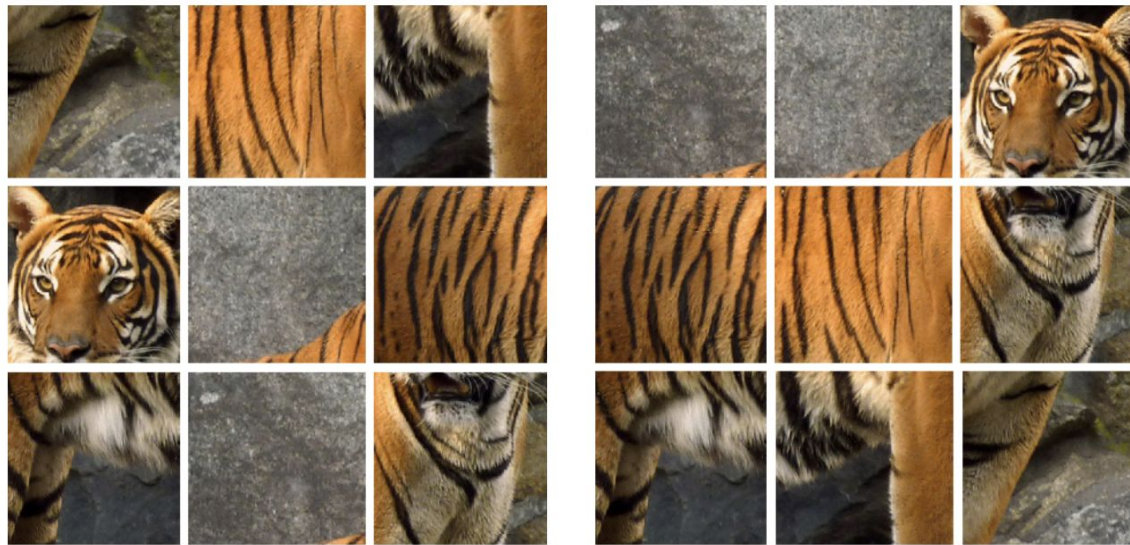
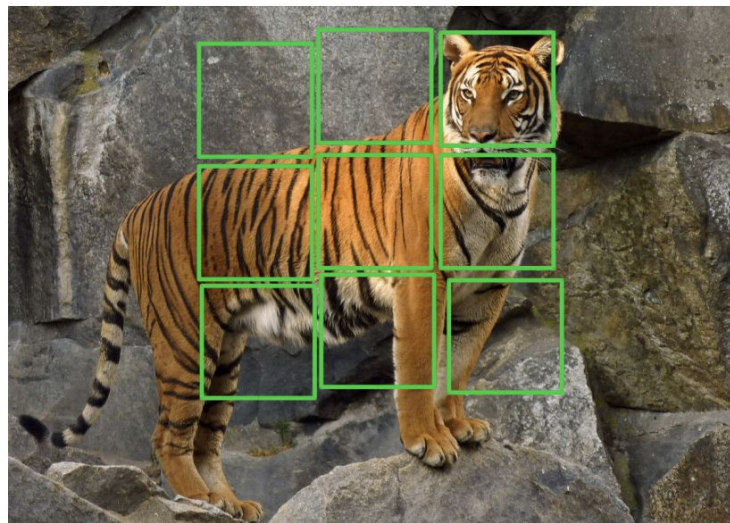
Concatenate (L, ab) channels

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



Slide: Richard Zhang

Solving Jigsaw Puzzles

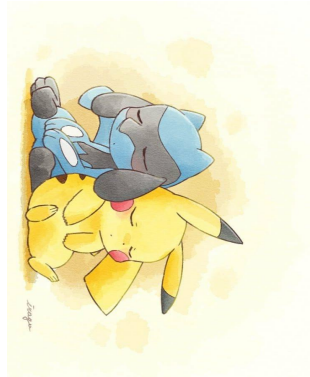


Slide: CS294-158

Rotation



90° rotation



270° rotation

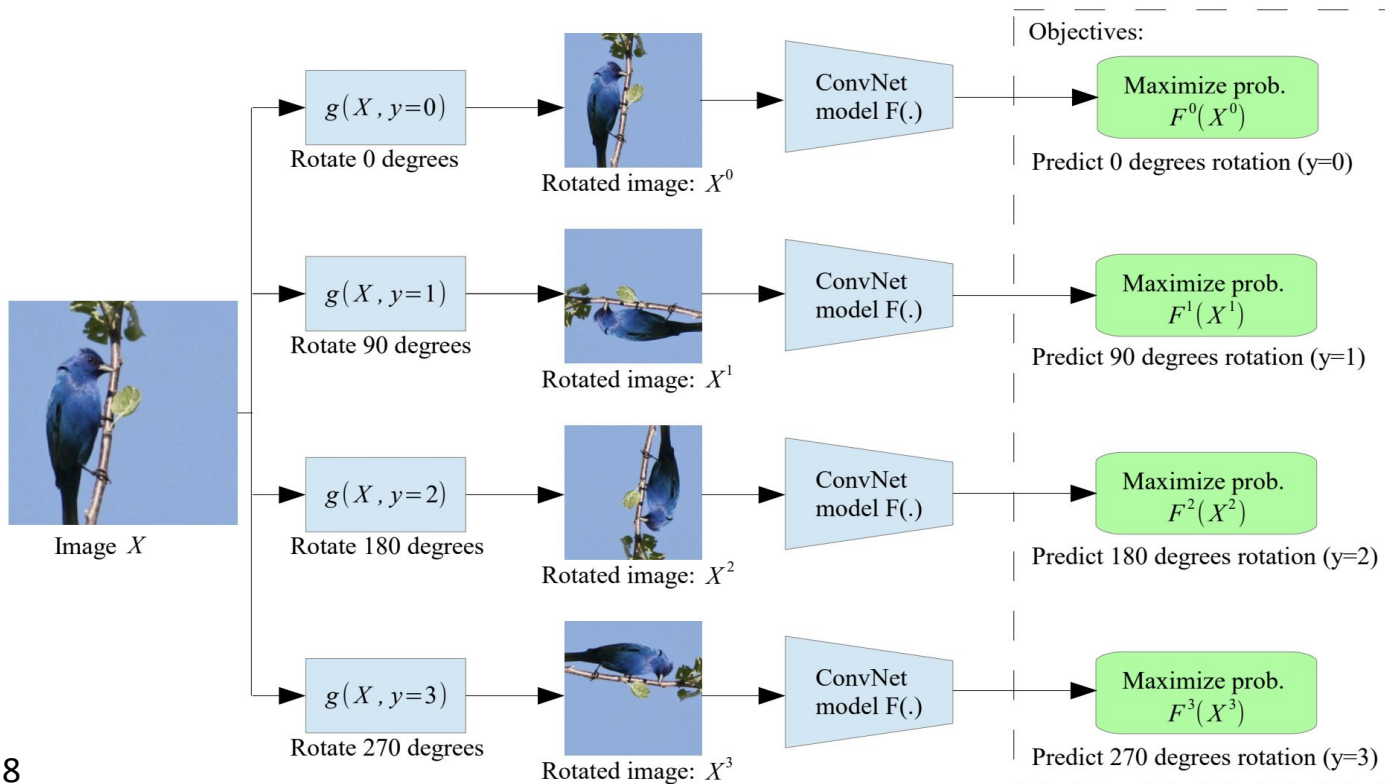


180° rotation



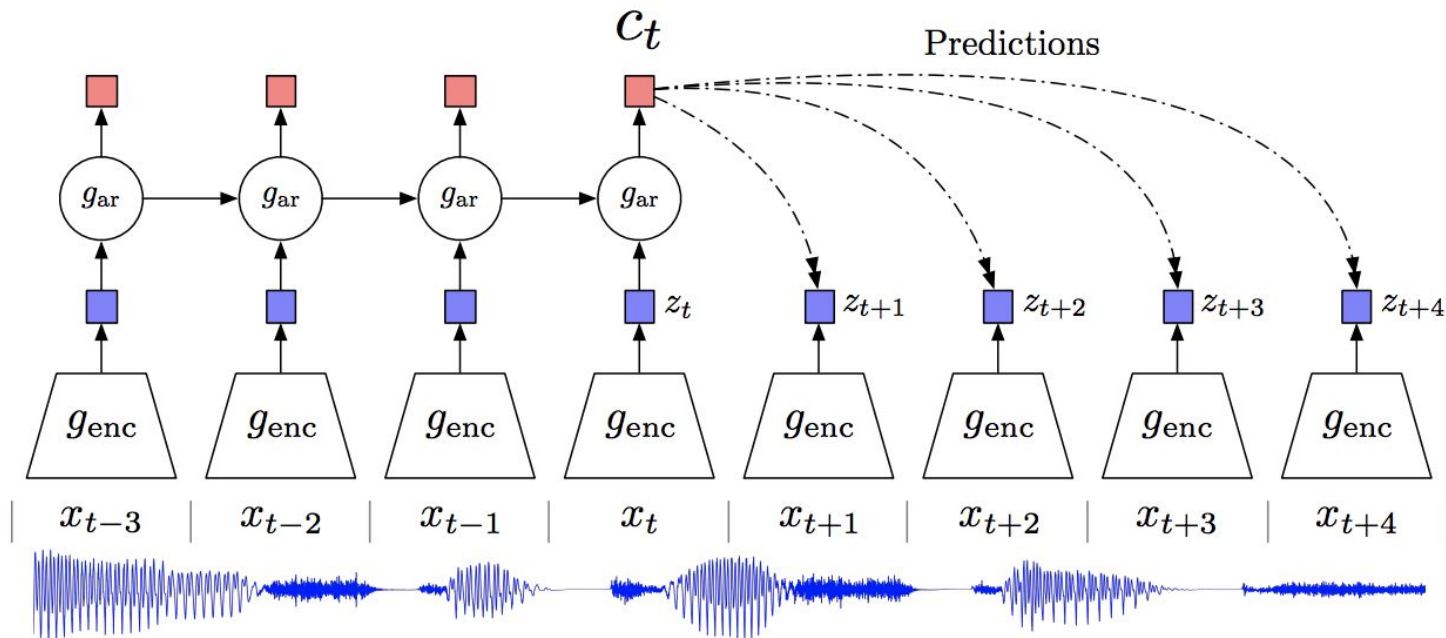
0° rotation

Rotation



Slide: CS294-158

Contrastive Predictive Coding

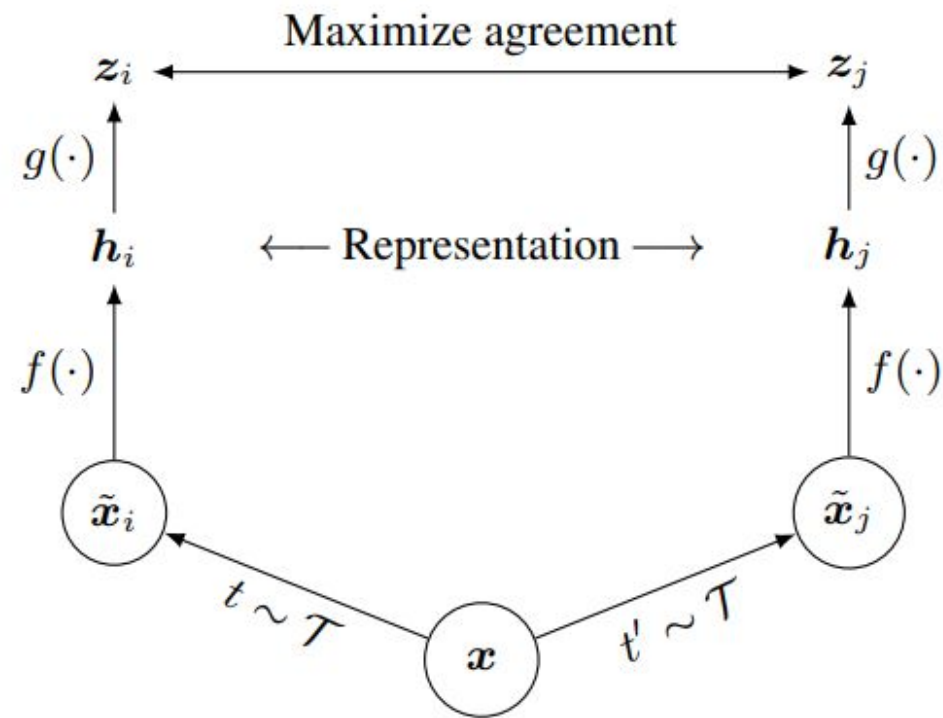


Remember Word2Vec? They are almost the same idea.

van den Oord et al. 2020

$$\mathcal{L}_N = -\mathbb{E}_X \left[\log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)} \right]$$

SimCLR



Chen et al. 2020

SimCLR



(a) Original



(b) Crop and resize



(c) Crop, resize (and flip)



(d) Color distort. (drop)



(e) Color distort. (jitter)



(f) Rotate $\{90^\circ, 180^\circ, 270^\circ\}$



(g) Cutout



(h) Gaussian noise



(i) Gaussian blur



(j) Sobel filtering

Reference

- CS294-158 Deep Unsupervised Learning Lecture 7
- AAIL 2020 Keynotes Turing Award Winners Event
- Learning From Text - OpenAI
- Learning from Unlabeled Data - Thang Luong