# Generative model intro & Autoregressive model

Instructor: Seunghoon Hong

# Course overview

- Image classification
- Object detection
- Semantic segmentation
- Visualization
- Style transfer
- Adversarial attacks

- Text modeling
- Machine translation
- Image captioning
- Visual question answering
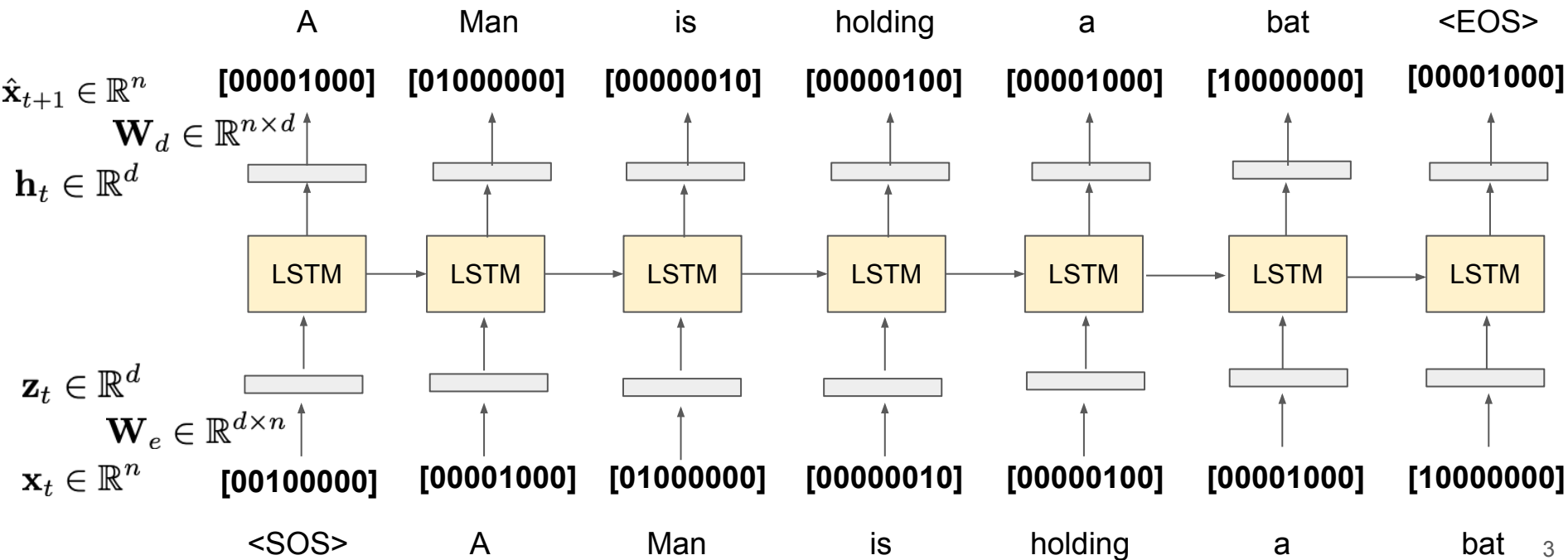
- Image generation
- Text generation
- Img-to-img translation

- Attention and versatile networks
- Self- and Semi-supervised learning
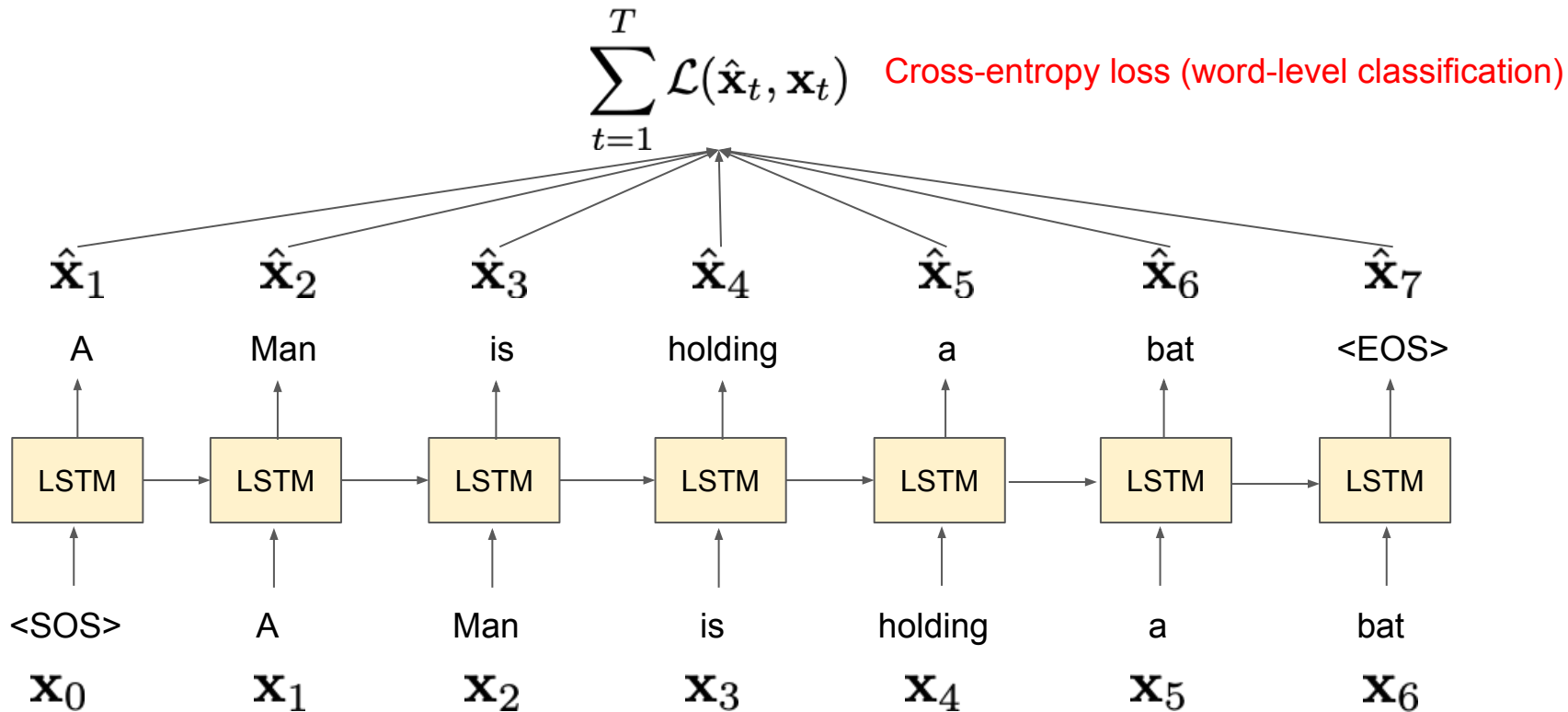- Multi-modal learning
- Graph neural networks

We are here!

Convolutional
Neural Networks (CNN)

Recurrent
Neural Networks (RNN)

Deep generative models

Advanced topics

# Recap: RNN as a language model

- Sentence generation = predicting a next token



$\hat{\mathbf{x}}_{t+1} \in \mathbb{R}^n$

$\mathbf{W}_d \in \mathbb{R}^{n \times d}$

$\mathbf{h}_t \in \mathbb{R}^d$

$\mathbf{z}_t \in \mathbb{R}^d$

$\mathbf{W}_e \in \mathbb{R}^{d \times n}$

$\mathbf{x}_t \in \mathbb{R}^n$

| A | Man | is | holding | a | bat | <EOS> |
|---|---|---|---|---|---|---|
| [00001000] | [01000000] | [00000010] | [00000100] | [00001000] | [10000000] | [00001000] |

LSTM → LSTM → LSTM → LSTM → LSTM → LSTM → LSTM

| [00100000] | [00001000] | [01000000] | [00000010] | [00000100] | [00001000] | [10000000] |
|---|---|---|---|---|---|---|
| <SOS> | A | Man | is | holding | a | bat |

3

# Recap: Training: RNN-based language model

$$\sum_{t=1}^{T} \mathcal{L}(\hat{\mathbf{x}}_t, \mathbf{x}_t)$$

Cross-entropy loss (word-level classification)

$\hat{\mathbf{x}}_1$    $\hat{\mathbf{x}}_2$    $\hat{\mathbf{x}}_3$    $\hat{\mathbf{x}}_4$    $\hat{\mathbf{x}}_5$    $\hat{\mathbf{x}}_6$    $\hat{\mathbf{x}}_7$

A    Man    is    holding    a    bat    <EOS>

| LSTM | → | LSTM | → | LSTM | → | LSTM | → | LSTM | → | LSTM | → | LSTM |

<SOS>    A    Man    is    holding    a    bat

$\mathbf{x}_0$    $\mathbf{x}_1$    $\mathbf{x}_2$    $\mathbf{x}_3$    $\mathbf{x}_4$    $\mathbf{x}_5$    $\mathbf{x}_6$

We feed ground-truth words as inputs (also known as **teacher forcing**)

# Recap: Machine translation

- Translate a sentence in one language to another

Er liebte zu essen .

Softmax

Encoder → Decoder

Embed

NULL Er liebte zu essen

He loved to eat .

# Recap: Bayes' Theorem

- Conditional probability

$$P(A|B) = \frac{P(A,B)}{P(B)}, \quad P(B|A) = \frac{P(A,B)}{P(A)}$$

$$P(A,B) = P(A|B)P(B) = P(B|A)P(A)$$

- Bayes' Theorem

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$posterior \propto likelihood \times prior$$

# Today's agenda

- Introduction to generative models
- Autoregressive models

# Introduction to generative models

# Machine Learning for Understanding Data

- Learning to **perceive** and **reason** from a data



Concepts?
**Person, elephant, field, sky, fence**

Relationship between concepts?
**One person is holding another**
**Two people are standing next to fence**
**An elephant is standing on a grass**

Context?
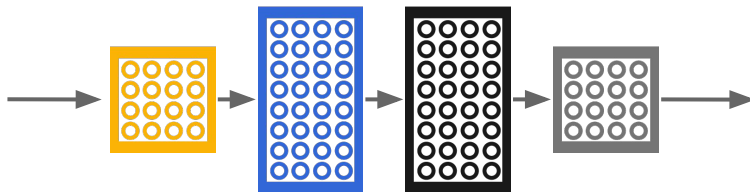**A father went to zoo with his son watching an elephant**

# Understanding via Categorization

- Learning to associate input to pre-defined, task-specific labels
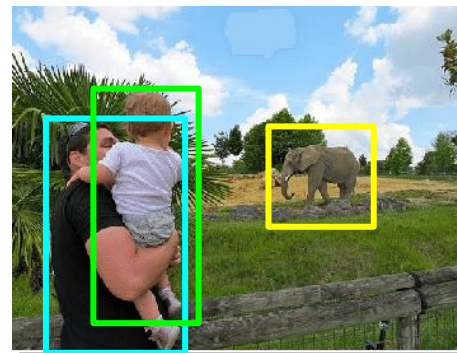- Examples: **classification** (concept)

input x

model

$$f_\theta(x)$$

output y

"dog"

**"person"**

"apple"

**"elephant"**

⋮

**"field"**

# Understanding via Categorization

- Learning to associate input to pre-defined, task-specific labels
- Examples: **classification** (relationship)

# Understanding via Categorization

- Learning to associate input to pre-defined, task-specific labels
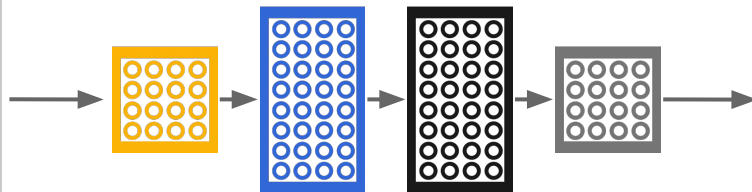- Examples: classification, **detection**

input x

$$f_\theta(x)$$

output y

# Understanding via Categorization

- Learning to associate input to pre-defined, task-specific labels
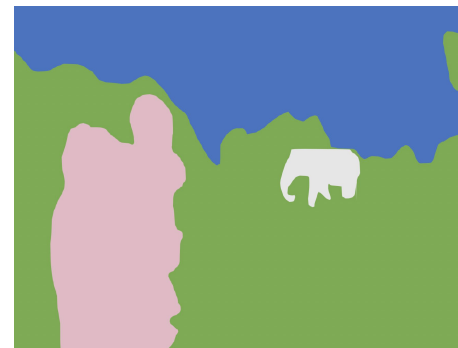- Examples: classification, detection, **segmentation**, …

input x
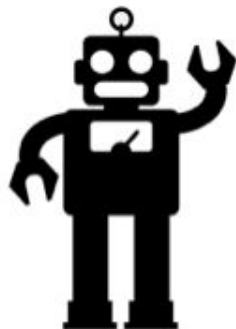
$$f_\theta(x)$$

output y

# Understanding via Categorization

- Limitations
  - Requires labels (human annotations) for training
  - Learns a biased knowledge to solve the specific task

# Understanding via Generation

- Learning to synthesize the data itself
- Why do we care about generation?
  - Generation requires implicit understanding of underlying structure of data
  - No need for labels → <mark>unsupervised learning</mark>
  - Can learn something useful for downstream tasks
  - Generated data itself can be useful



generate

data x

# How do we define the task of *generation*?

- What is the task of generative modeling?
- What is the objective function for learning generative models?

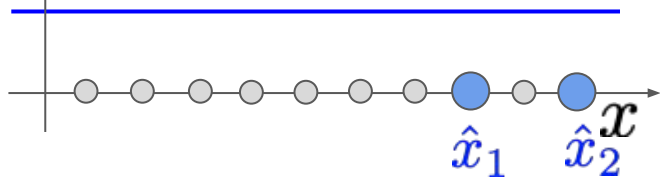# Q1: what is it like building a generative model?

- Consider that our data is composed of 1d points
- Example 1: 1d data is distributed according to Uniform distribution

$p(x)$

Training data의 distribution을 알 수 있다면...
그 distribution을 기반으로 sampling 가능

$$p(x) = \text{Uniform}(0, a)$$

once we know the underlying function that produces data, we can create new samples by sampling from (or to follow) this function

$\hat{x}_1$  $\hat{x}_2$  $x$

# Q1: what is it like building a generative model?

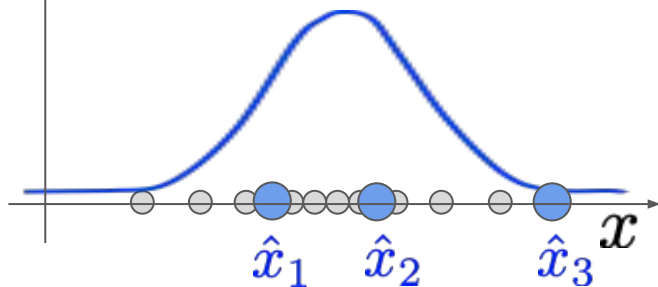- Consider that our data is composed of 1d points
- Example 2: data is distributed according to 1d Gaussian distribution

**Generative model**
**- Learn the probability distribution from training data!**

$$p(x) = c \cdot \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Again, once we know the underlying function that produces data, we can create new samples by sampling from (or to follow) this function

# Q1: what is it like building a generative model?

- If we know the ground-truth distribution of data, we can generate new one by sampling from the distribution (or to follow the distribution)

$$\hat{x} \sim P(X)$$

즉, we assume that training data follows some distribution

- Since the ground-truth distribution is unknown, we use a neural network to approximate the distribution

$$G_\theta \approx P(X)$$

# Deep Generative Models

Assumption:
many many training data will
follow some distribution

- Learning a model that its outputs follow the true data distribution

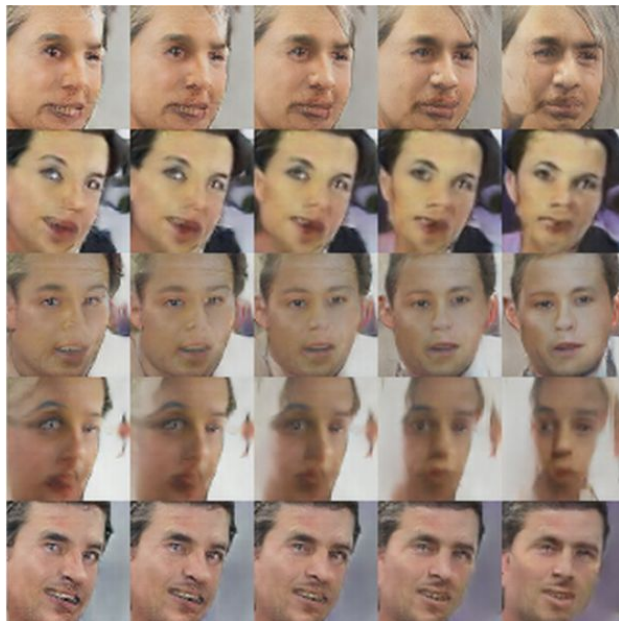$$G_\theta \sim P(X)$$

Generated images from G



True Images X

# Recent applications of generative models

# Generative models have improved enormously



2014

Goodfellow et al.

2016

Radford et al.

2018

Karras et al.
[slide credit: Tim Saliman]

# Generative models have improved enormously



a pitcher is about to throw the ball to the batter.

a picture of a very clean living room.

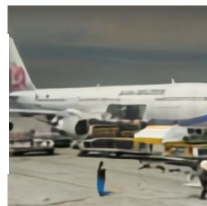a sheep standing in a open grass field.

a very cute cat laying by a big bike.

china airlines plain on the ground at an airport with baggage cars nearby.

a table that has a train model on it with other cars and things

A cute corgi lives in a house made out of sushi.

**2016**

Reed et al.

**2021**

Ramesh et al.

**2022**

Saharia et al.

# Generative models are started to become useful

- Image upsampling / image compression

Ledig et al. (2016)



| original | bicubic (21.59dB/0.6423) | SRResNet (23.44dB/0.7777) | SRGAN (20.34dB/0.6562) |

# Generative models are started to become useful

- Video synthesis



https://tcwang0509.github.io/vid2vid/

# Generative models are started to become useful

- Text-to-image synthesis



A wall in a royal castle. There are two paintings on the wall. The one on the left a detailed oil painting of the royal raccoon king. The one on the right a detailed oil painting of the royal raccoon queen.

A group of teddy bears in suit in a corporate office celebrating the birthday of their friend. There is a pizza cake on the desk.

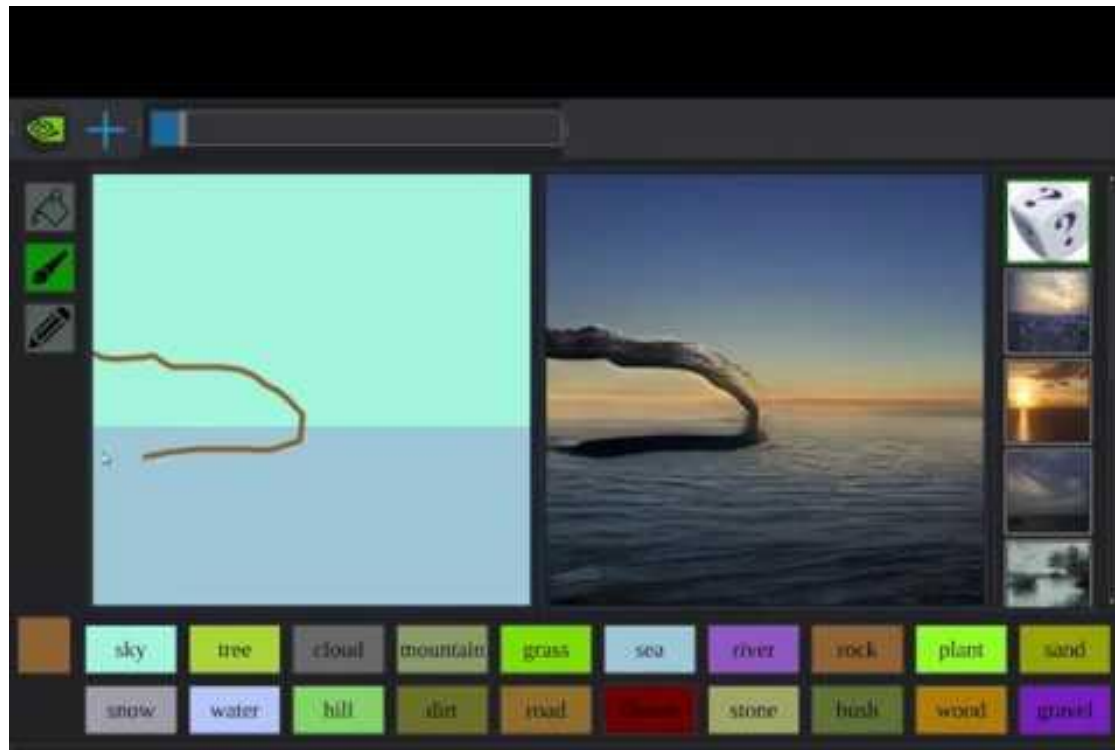A chrome-plated duck with a golden beak arguing with an angry turtle in a forest.

Saharia et al., Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding

# Generative models are started to become useful

- Text-to-image synthesis

# Generative models are started to become useful

- Interactive drawing
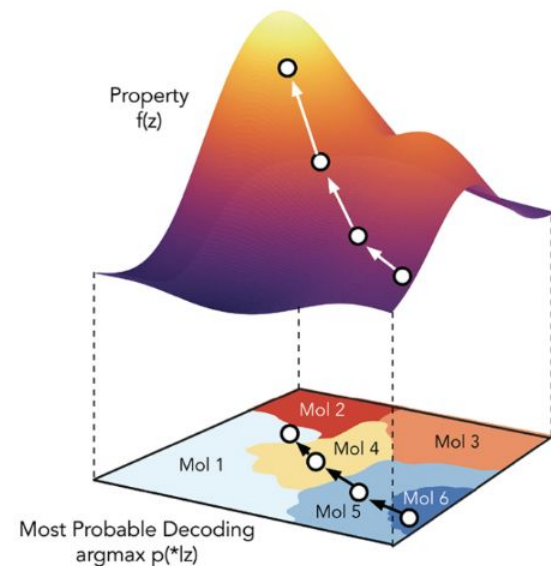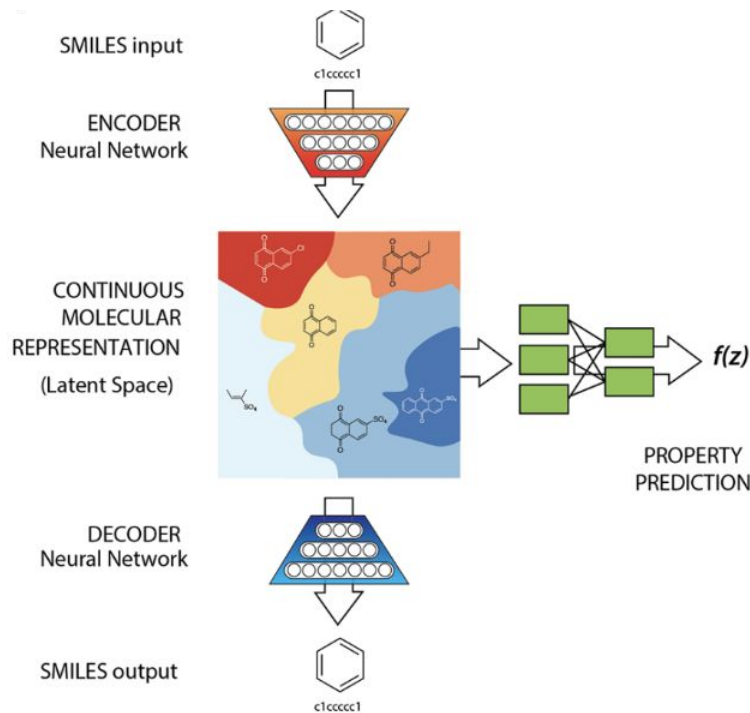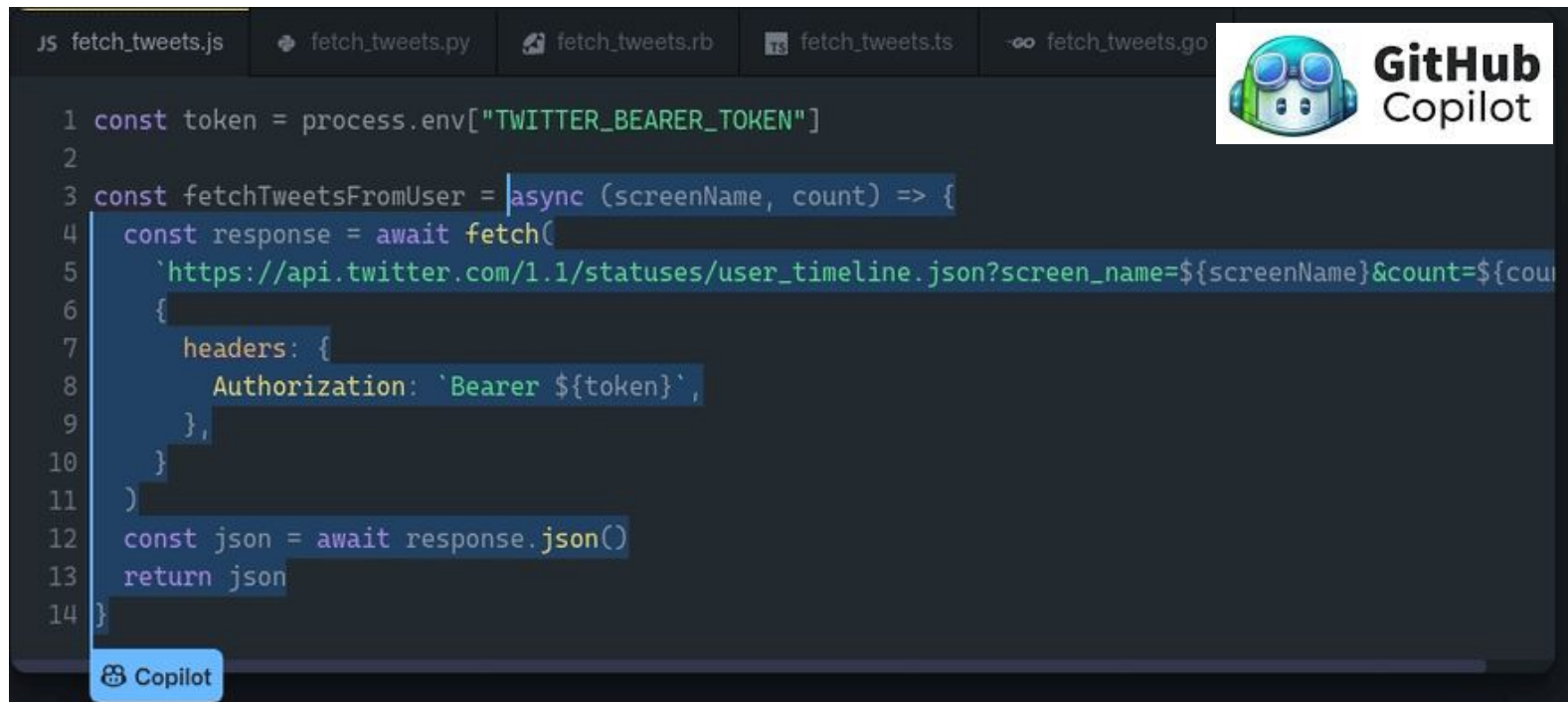
# Generative models are started to become useful

- Drug discovery

# Generative models are started to become useful

- Code generation

# Generative models are started to become useful

- Text-to-speech or speech-to-text synthesis

hand-coded

Wavenet generative model

# Autoregressive models

# What is the loss for generative model?

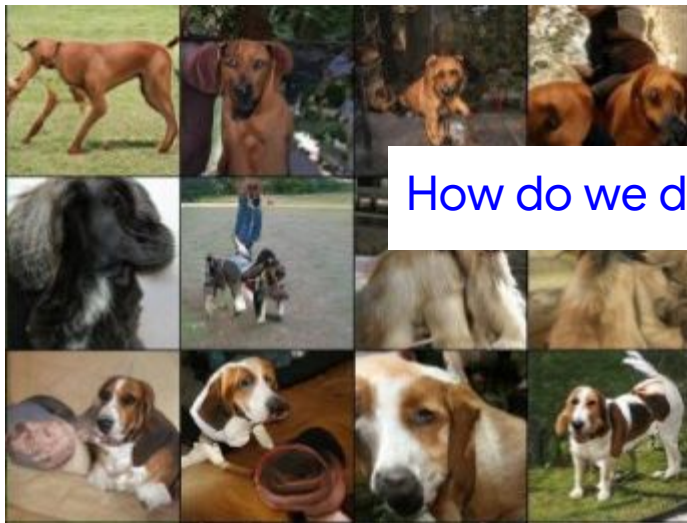- Learning objective:

  train a generator such that its outputs are distributed according to the target distribution

  $$G_\theta \approx P(X)$$

Generated images from G

True Images X



How do we define a loss for this?

# What is the loss for generative model?

- Let's consider that our generator (e.g. neural network) produces a probability measure for any input data x

- Generator output for input x: $p_\theta(\mathbf{x})$

- **Assume** that we know the true probability: $p^*(\mathbf{x})$

- Then, we will define a loss as a discrepancy b/w two probability distributions



generated distribution

true data distribution

$p_\theta(x)$

$p(x)$

loss

image space

image space

# What is the loss for generative model?

- There are various distance measures for two probability distributions
- Here, we consider **KL divergence**

$$\min_\theta D_{KL}[p^*(x)||p_\theta(x)] = \mathbb{E}_{p^*(x)}\left[\log\frac{p^*(x)}{p_\theta(x)}\right]$$

$$= \int p^*(x)\log\frac{p^*(x)}{p_\theta(x)}dx$$

Cross-entropy loss!

In the sense of minimizing cross-entropy loss, it is similar to supervised learning

$$= \int p^*(x)\log p^*(x)dx \boxed{- \int p^*(x)\log p_\theta(x)dx}$$

$\mathsf{L}$ But we don't know actually the true dist. $p^*(x)$.

This term is irrelevant to parameters

$$\leftrightarrow \min_\theta \boxed{\mathbb{E}_{p^*(x)}[-\log p_\theta(x)]}$$

# What is the loss for generative model?

$$\min_{\theta} \mathbb{E}_{\boxed{p^*(x)}}[-\log p_\theta(x)]$$

- Remaining issues:
    - We do not know the true distribution
    - We do not have an access to infinite amount of data for expectation!
      Even if we have one, it is computationally intractable.

(위의 두 문제를 동시에 해결!)

Solution using samples:    ∵ Sample이 true dist. 로부터 나왔으므로

    - We assume that the training data approximate the true distribution
    - Then we can optimize the following    $x_i \sim p^*(x)$.    $N$개 : weights 는 Sample에 의해
                                                                    자연스럽게 반영됨. (Think (Taylor 10k)

$$\min_{\theta} \boxed{-\frac{1}{N}\sum_{I=1}^{N}\log p_\theta(x_i)} = \mathbb{E}_{p^*(x)}[-\log p_\theta(x)] \quad \text{if } N \to \infty$$

# What is the loss for generative model?

- Maximum Likelihood Estimation (MLE)

Intuition!
increasing prob of dog images
<-> decreasing prob. of other images
즉, 다른 경우를 surpress!

$$\hat{\theta} = \arg\max_{\theta \in \Theta} \sum_{x_i \in \mathcal{X}} \boxed{\log p_\theta(x_i)}$$

Log Likelihood

$$\leftrightarrow \arg\max_{\theta \in \Theta} \prod_{x_i \in \mathcal{X}} \boxed{p_\theta(x_i)}$$

Likelihood

Find model parameters that maximize the probability of sampling training data (likelihood)

$$\leftrightarrow \arg\min_{\theta \in \Theta} \sum_{x_i \in \mathcal{X}} -\log p_\theta(x_i)$$

In practice, we minimize the negative log likelihood for gradient descent

# Challenges in evaluating likelihood

$$\hat{\theta} = \arg\min_{\theta \in \Theta} \sum_{x_i \in \mathcal{X}} - \log \boxed{p_\theta(x_i)}$$

- For high-dimensional data, it is difficult to optimize the joint distribution at once

$$x_i = [x_i^1, x_i^2, x_i^3, \ldots, x_1^d] \in \mathbb{R}^d$$
$$p(x_i) = p(x_i^1, x_i^2, x_i^3, \ldots, x_1^d)$$

- Examples of high-dimensional data
  - Image (d = number of pixels)
  - Sentence (d = length of sentence)

# Auto-Regressive Model (AR)

- Factorizing the likelihood via **chain rule**

$$p(a, b) = p(a|b)p(b)$$

# Auto-Regressive Model (AR)

- Factorizing the likelihood via chain rule

(즉, $x$는 T-dim vector 가정)

$$p_\theta(x) = p_\theta(x_1, x_2, x_3, ..., x_T)$$

$$= p_\theta(x_T | x_1, x_2, ..., x_{T-1}) \boxed{p_\theta(x_1, x_2, ..., x_{T-1})}$$

**apply recursively**
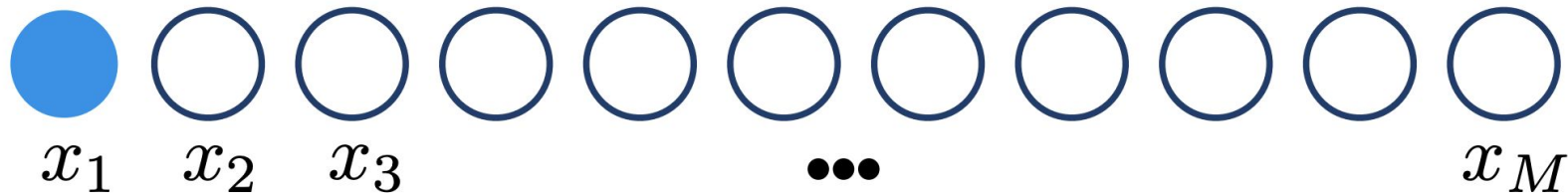
$$= \prod_{t=1}^{T} p_\theta(x_t | x_1, ..., x_{t-1})$$

each conditional has **lower dimension** 을.

— predicting one value at a time.

# Auto-Regressive Model (AR)

$$p_\theta(x) = \prod_{t=1}^{T} p_\theta(x_t | x_1, ..., x_{t-1})$$

$p(x_1)$



$x_1 \quad x_2 \quad x_3 \qquad\qquad \bullet\bullet\bullet \qquad\qquad x_M$

# Auto-Regressive Model (AR)

$$p_\theta(x) = \prod_{t=1}^{T} p_\theta(x_t | x_1, ..., x_{t-1})$$
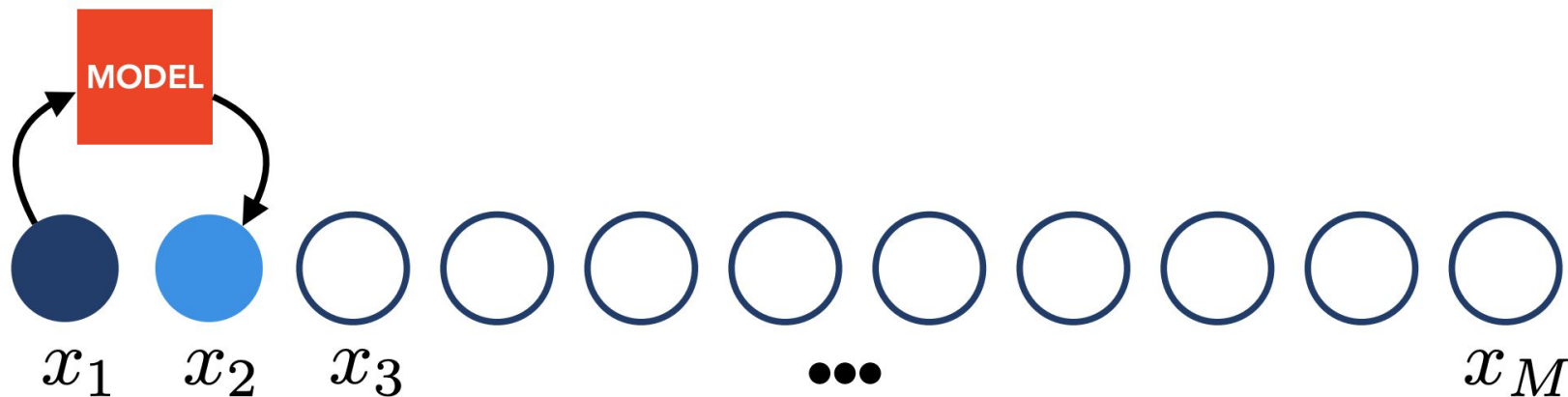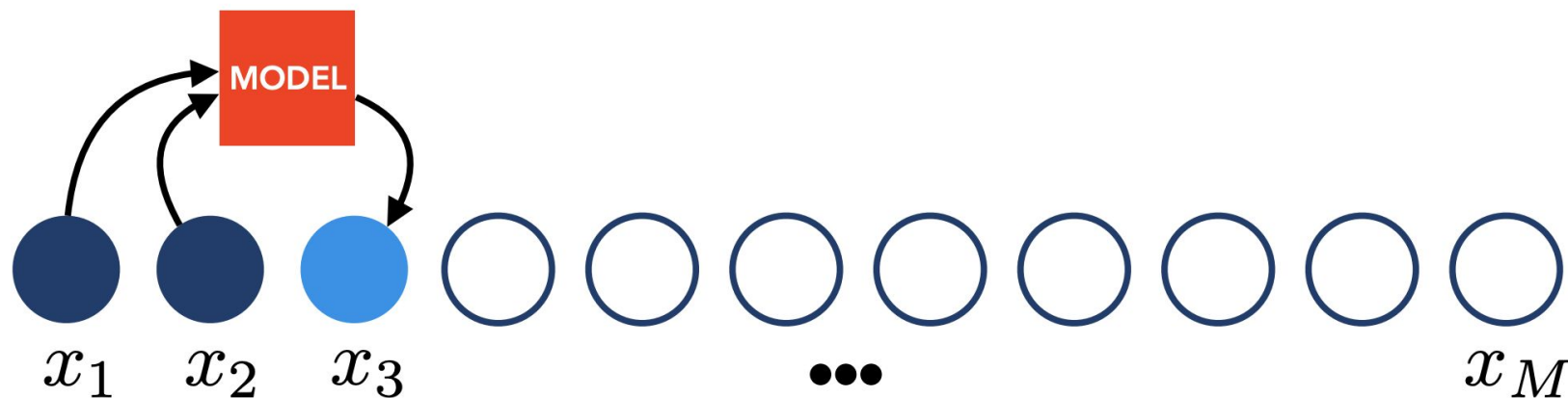
$p(x_2 | x_1)$

# Auto-Regressive Model (AR)

$$p_\theta(x) = \prod_{t=1}^{T} p_\theta(x_t | x_1, ..., x_{t-1})$$

$$p(x_3 | x_2, x_1)$$

# Auto-Regressive Model (AR)

$$p_\theta(x) = \prod_{t=1}^{T} p_\theta(x_t | x_1, ..., x_{t-1})$$

$$p(x_4 | x_3, x_2, x_1)$$



$x_1 \quad x_2 \quad x_3 \qquad\qquad \bullet\bullet\bullet \qquad\qquad x_M$

# Auto-Regressive Model (AR)

$$p_\theta(x) = \prod_{t=1}^{T} p_\theta(x_t | x_1, ..., x_{t-1})$$

$$p(x_M | x_{M-1}, \ldots, x_1)$$



**MODEL**

$x_1 \quad x_2 \quad x_3 \qquad \bullet\bullet\bullet \qquad x_M$

# Autoregressive models

- Factorized objective function

$$\hat{\theta} = \arg\min_{\theta \in \Theta} \sum_{x_i \in \mathcal{X}} -\log p_\theta(x_i)$$

$$= \arg\min_{\theta \in \Theta} \sum_{x_i \in \mathcal{X}} \sum_{1 \leq t \leq d} -\log p_\theta(x_i^t | x_i^1, \ldots, x_i^{t-1})$$

# RNN revisited

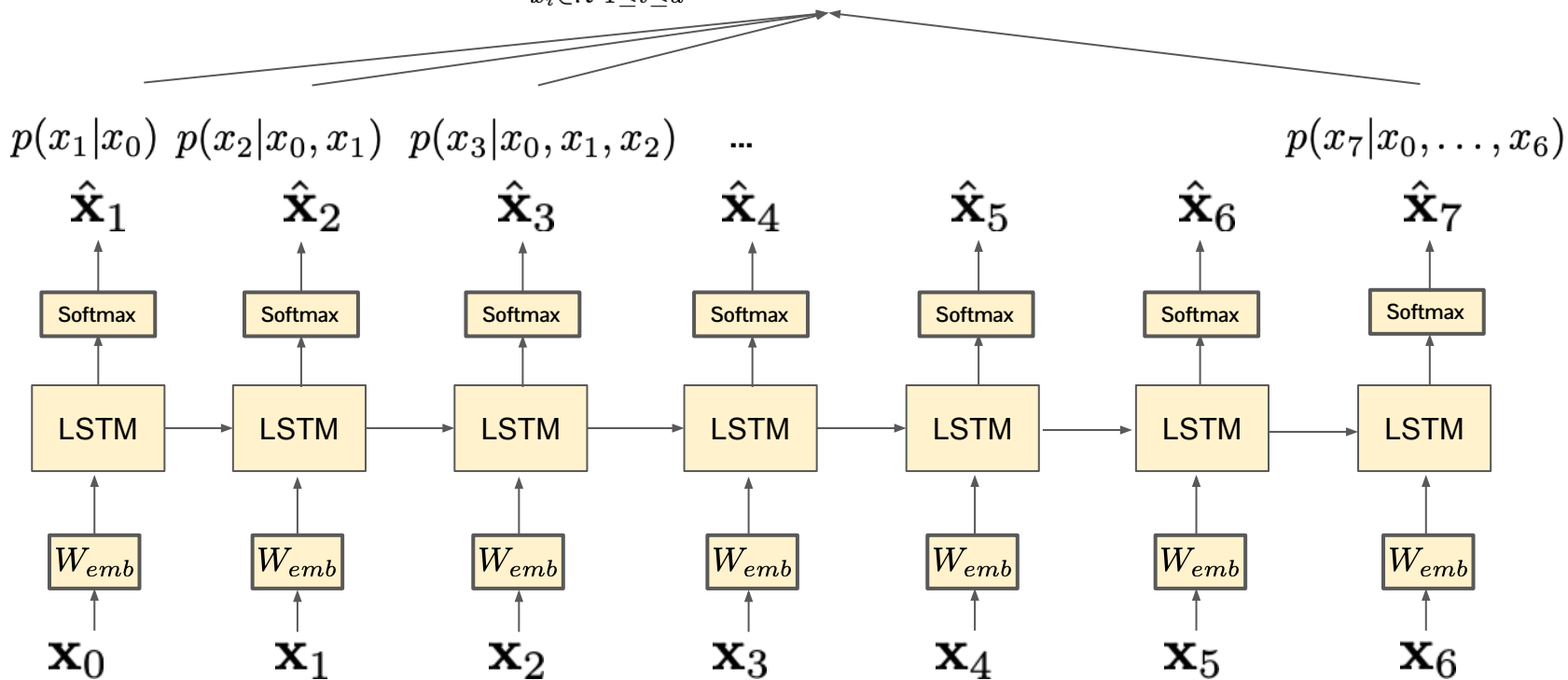즉, RNN은 text generalization에서 generative model!

$$\sum_{x_i \in \mathcal{X}} \sum_{1 \le t \le d} - \log p_\theta(\hat{x}_i^t = x_i^t | x_i^1, \ldots, x_i^{t-1})$$

$p(x_1|x_0)$  $p(x_2|x_0, x_1)$  $p(x_3|x_0, x_1, x_2)$  ...  $p(x_7|x_0, \ldots, x_6)$

$\hat{\mathbf{x}}_1$  $\hat{\mathbf{x}}_2$  $\hat{\mathbf{x}}_3$  $\hat{\mathbf{x}}_4$  $\hat{\mathbf{x}}_5$  $\hat{\mathbf{x}}_6$  $\hat{\mathbf{x}}_7$

| Softmax | Softmax | Softmax | Softmax | Softmax | Softmax | Softmax |

| LSTM | LSTM | LSTM | LSTM | LSTM | LSTM | LSTM |

| $W_{emb}$ | $W_{emb}$ | $W_{emb}$ | $W_{emb}$ | $W_{emb}$ | $W_{emb}$ | $W_{emb}$ |

$\mathbf{x}_0$  $\mathbf{x}_1$  $\mathbf{x}_2$  $\mathbf{x}_3$  $\mathbf{x}_4$  $\mathbf{x}_5$  $\mathbf{x}_6$

# Auto-Regressive Model: A Summary

- Maximizing factorized likelihood
    - Generate data one-by-one conditioned on previous outputs
- Appropriate to handle sequential data
    - Text, audio, video
- Fully-observable model
    - No latent representation of data

# Challenges

- Modeling long-term dependency
- Serial processing → difficult for parallelization

# Next

- Case study: autoregressive models
  - AR with attention for modeling long-term dependency
  - Task: machine translation