

# DSC291: Machine Learning with Few Labels

## Enhancing Large Language Models: Overview

**Zhiting Hu**

Lecture 11, April 24, 2024

**UC San Diego**

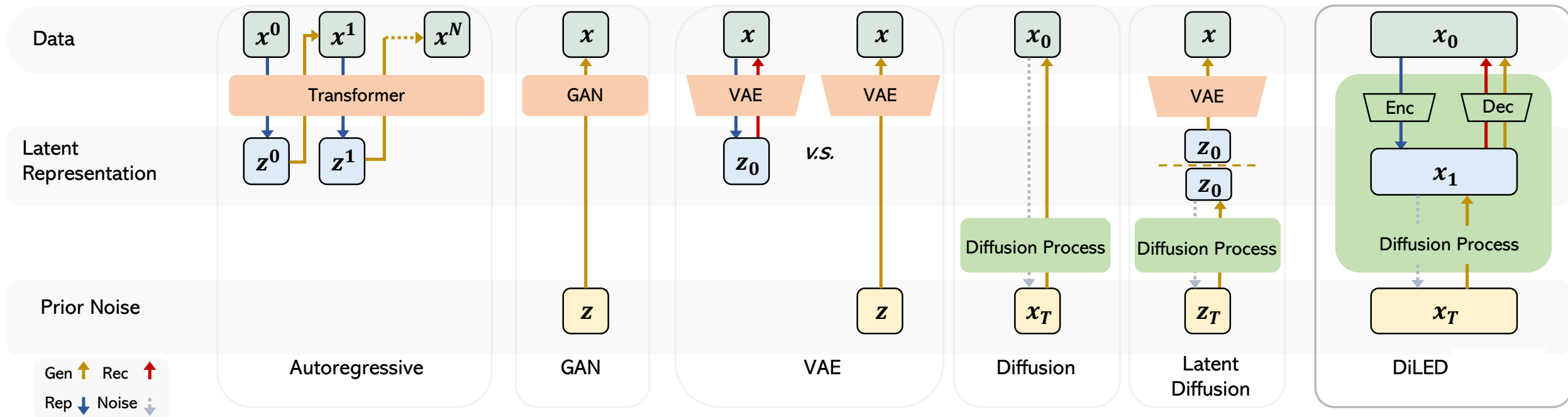
**HALICIOĞLU DATA SCIENCE INSTITUTE**

# Discussion

- **No Free Lunch (NFL) theorem** (suggested reading of Lecture#10):
  - No single machine learning algorithm is universally the best-performing algorithm for all problems
- Do generalist models (LLMs) violate this theorem?
- Does "the Bitter Lesson" contradict with this theorem?
  - (suggested reading of Lecture#6)

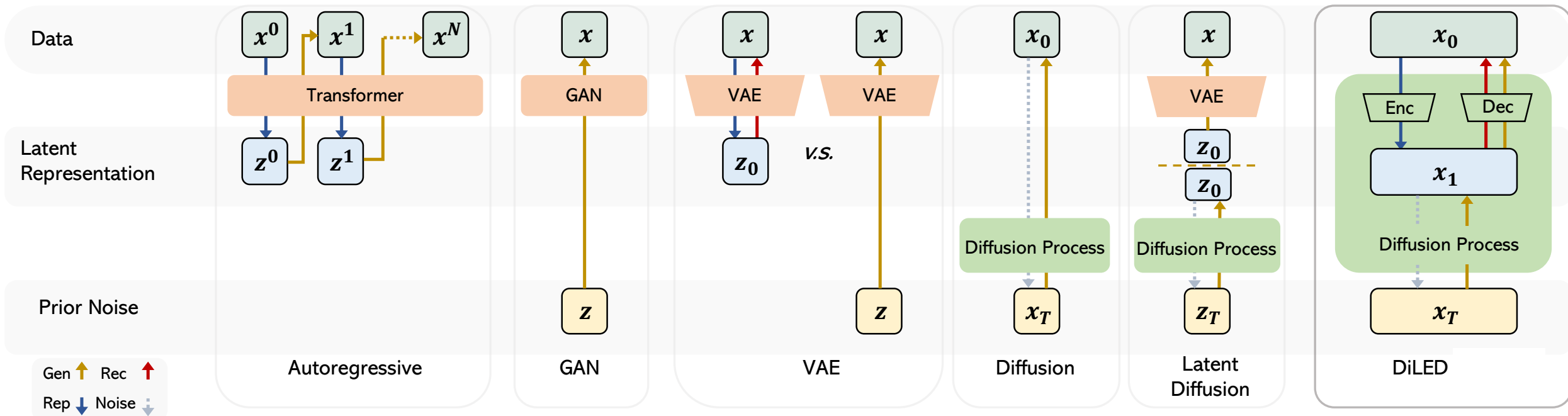
# Latent-space Reasoning (Recap)

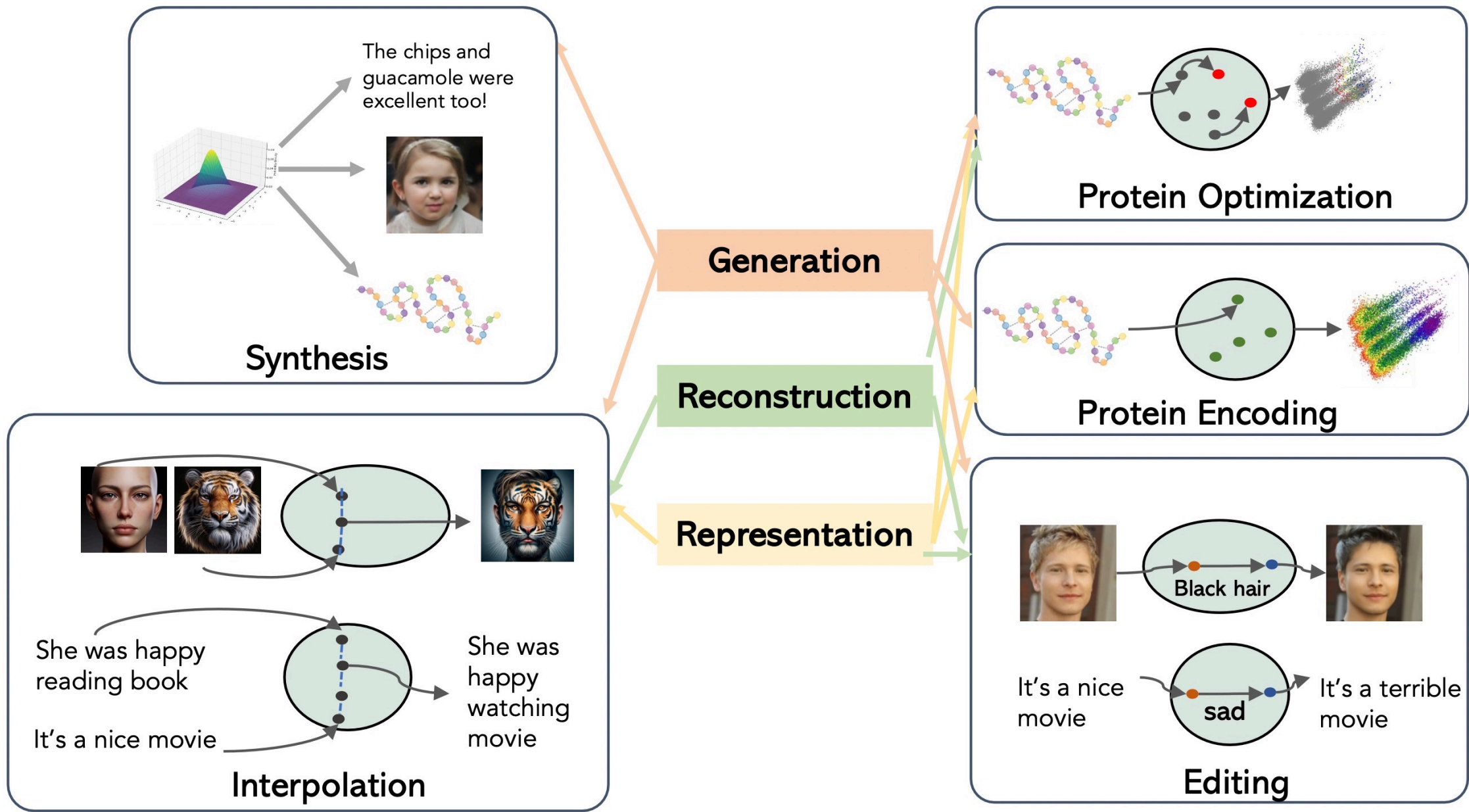
- But how to learn a good latent space in the first place?



# Latent-space Reasoning

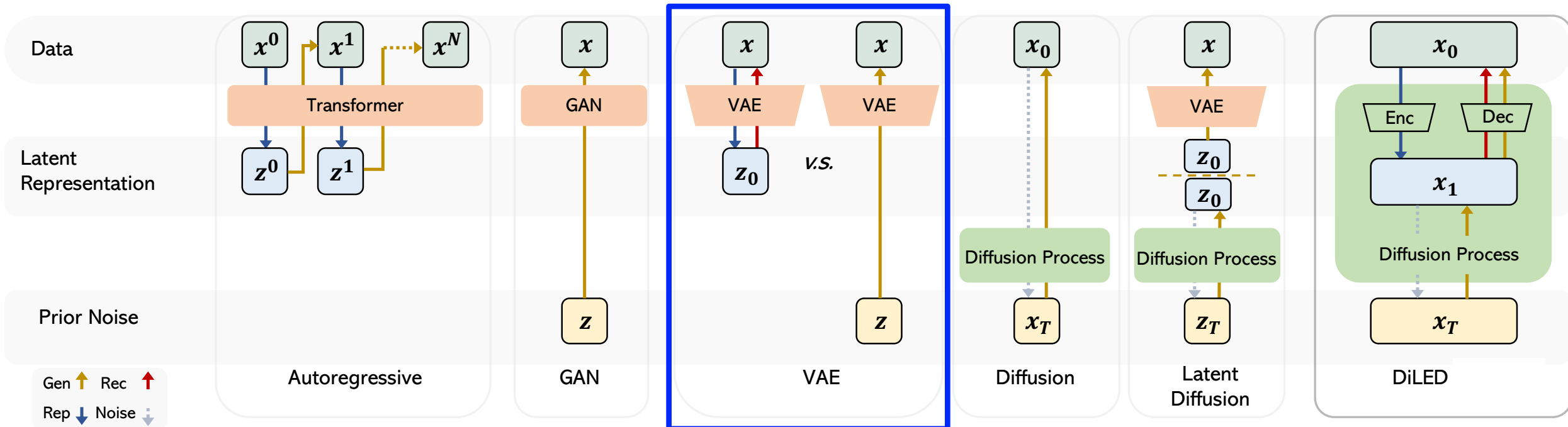
- But how to learn a good latent space in the first place?
  - Compact and well-structured **representation** of the world, enabling realistic **generation** and consistent **reconstruction**





# Latent-space Reasoning

- But how to learn a good latent space in the first place?
  - Compact and well-structured **representation** of the world, enabling realistic **generation** and consistent **reconstruction**



## Variational Autoencoders (VAEs)

# Supervised Learning, Unsupervised Learning

# Probabilistic Models: Why Probability?

- The world is a very uncertain place
  - “What will the weather be like today?”
  - “Will I like this movie?”
- We often can’t prove something is true, but we can still ask how likely different outcomes are or ask for the most likely explanations
- Predictions need to have associated confidence
  - Confidence -> probability
- Not all machine learning models are probabilistic
  - ... but most of them have probabilistic interpretations





# Notations

- A random variable  $\mathbf{x}$  represents outcomes or states of the world.
  - We write  $p(\mathbf{x}_0)$  to mean Probability( $\mathbf{x} = \mathbf{x}_0$ )
- Sample space: the space of all possible outcomes (may be discrete, continuous, or mixed)
- $p(\mathbf{x})$  is the probability mass (density) function
  - Assigns a number to each point in sample space
  - Non-negative, sums (integrates) to 1
  - Intuitively: how often does  $\mathbf{x}$  occur, how much do we believe in  $\mathbf{x}$ .

# Notations

- Joint distribution  $p(\mathbf{x}, \mathbf{y})$
- Conditional distribution  $p(\mathbf{y}|\mathbf{x})$

- $p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{y})}{p(\mathbf{x})}$

- Expectation:

$$\mathbb{E}[f(\mathbf{x})] = \sum_{\mathbf{x}} f(\mathbf{x}) p(\mathbf{x})$$

or

$$\mathbb{E}[f(\mathbf{x})] = \int_{\mathbf{x}} f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

# Rules of Probability

- Sum rule

$$p(x) = \sum_y p(x, y) \quad (\text{Marginalize out } y)$$

$$p(x_1) = \sum_{x_2} \sum_{x_3} \dots \sum_{x_N} p(x_1, x_2, \dots, x_N)$$

- Product/chain rule

$$p(x, y) = p(y | x) p(x)$$

$$p(x_1, \dots, x_N) = p(x_1) p(x_2 | x_1) \dots p(x_N | x_1, \dots, x_{N-1})$$

# Bayes' Rule

$$p(\mathbf{y}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{y})p(\mathbf{y})}{p(\mathbf{x})}$$

- This gives us a way of “reversing” conditional probabilities
- We call  $p(\mathbf{y})$  the “prior”, and  $p(\mathbf{y}|\mathbf{x})$  the “posterior”
- Ex: Bayes' Rule in machine learning:
  - $\mathcal{D}$ : data (evidence)
  - $\theta$ : unknown quantities, such as model parameters, predictions

**Posterior** belief on the unknown quantities you see data  $\mathcal{D}$

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$

**Likelihood**: How likely is the observed data under the particular unknown quantities  $\theta$

**Prior** belief on the unknown quantities **before** you see data  $\mathcal{D}$

# Independence

- Two random variables are said to be **independent** iff their joint distribution factors

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$$

- Two random variables are **conditionally independent** given a third if they are independent after conditioning on the third

$$p(\mathbf{x}, \mathbf{y}|\mathbf{z}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{y}|\mathbf{z})$$

# Entropy

- Shannon entropy  $H(p) = - \sum_x p(x) \log p(x)$ 
  - The average level of "information", "surprise", or "uncertainty" inherent to the variable  $x$ 's possible outcomes

# KL Divergence

- Kullback-Leibler (KL) divergence: measures the closeness of two distributions  $p(\mathbf{x})$  and  $q(\mathbf{x})$

$$\text{KL}(q(\mathbf{x}) \parallel p(\mathbf{x})) = \sum_{\mathbf{x}} q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})}$$

- a.k.a. Relative entropy
- $\text{KL} \geq 0$  (Jensen's inequality)
- Intuitively:
  - If  $q$  is high and  $p$  is high, then we are happy (i.e. low KL divergence)
  - If  $q$  is high and  $p$  is low then we pay a price (i.e. high KL divergence).
  - If  $q$  is low then we don't care (i.e. also low KL divergence, regardless of  $p$ )
- not a true "distance":
  - not commutative (symmetric)  $\text{KL}(p \parallel q) \neq \text{KL}(q \parallel p)$
  - doesn't satisfy triangle inequality

# Supervised Learning

- Model to be learned  $p_{\theta}(\mathbf{x})$
- Observe **full** data  $\mathcal{D} = \{ \mathbf{x}_i \}_{i=1}^N$ 
  - e.g.,  $\mathbf{x}_i$  includes both input (e.g., image) and output (e.g., image label)
  - $\mathcal{D}$  defines an empirical data distribution  $\tilde{p}(\mathbf{x})$ 
    - $\mathbf{x} \sim \mathcal{D} \Leftrightarrow \mathbf{x} \sim \tilde{p}(\mathbf{x})$

- Maximum Likelihood Estimation (MLE)
  - The most classical learning algorithm

$$\min_{\theta} - \mathbb{E}_{\mathbf{x} \sim \tilde{p}(\mathbf{x})} \left[ \log p_{\theta}(\mathbf{x}) \right]$$

- MLE is minimizing the KL divergence between the empirical data distribution and the model distribution

$$\text{KL}(\tilde{p}(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) = -\mathbb{E}_{\tilde{p}(\mathbf{x})} [\log p_{\theta}(\mathbf{x})] + H(\tilde{p}(\mathbf{x}))$$

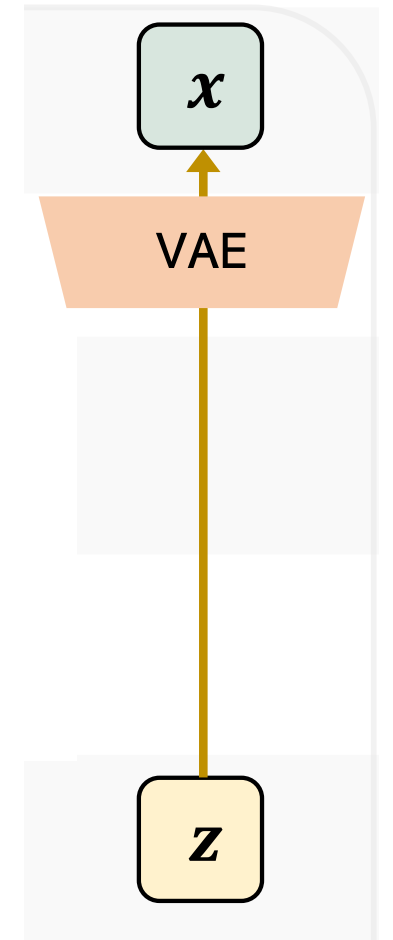


Cross entropy



# Unsupervised Learning

- Each data instance is partitioned into two parts:
  - observed variables  $\mathbf{x}$
  - latent (unobserved) variables  $\mathbf{z}$
- Want to learn a model  $p_{\theta}(\mathbf{x}, \mathbf{z})$



# Latent (unobserved) variables

- A variable can be unobserved (latent) because:
  - imaginary quantity: meant to provide some simplified and abstractive view of the data generation process
    - e.g., speech recognition models, mixture models, ...

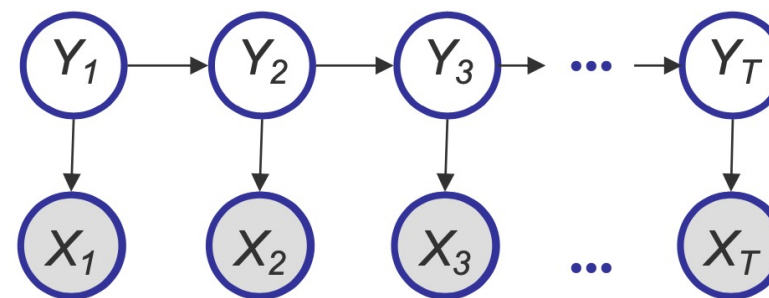
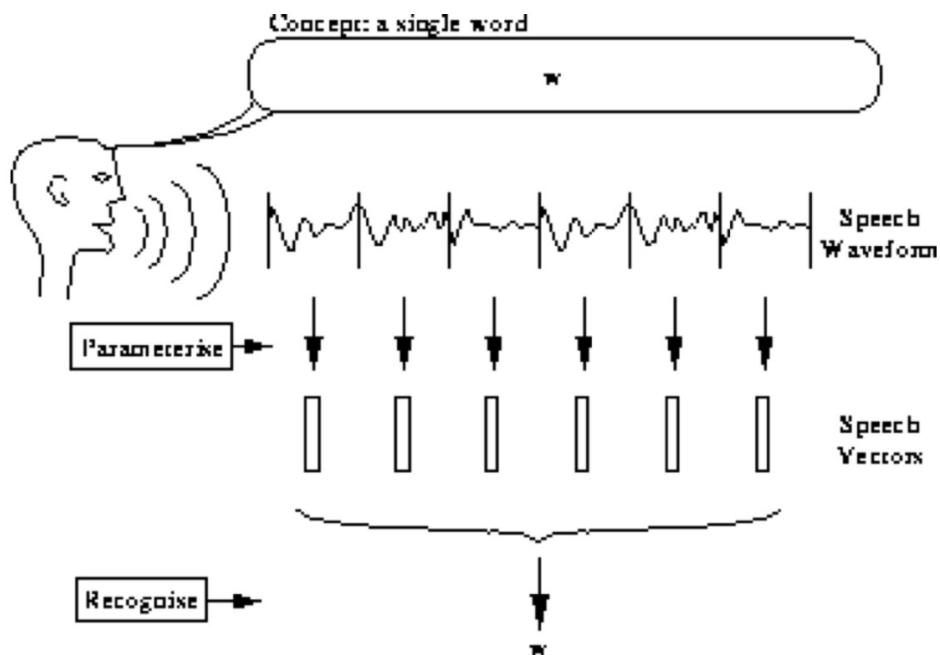
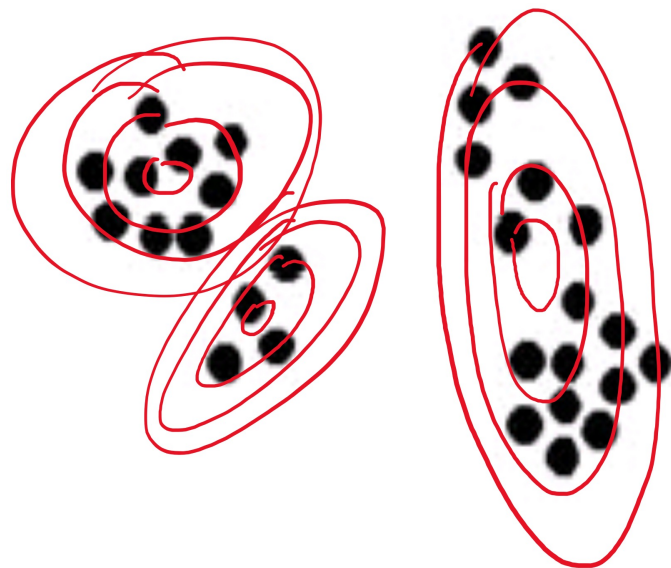


Fig. 1.2 Isolated Word Problem

# Latent (unobserved) variables

- A variable can be unobserved (latent) because:
  - imaginary quantity: meant to provide some simplified and abstractive view of the data generation process
    - e.g., speech recognition models, mixture models, ...



# Latent (unobserved) variables

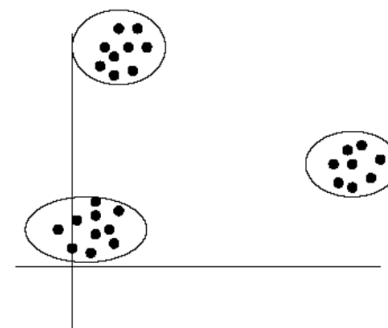
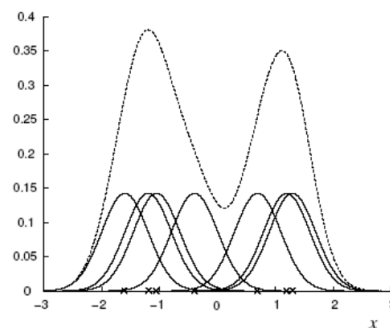
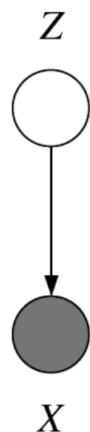
- A variable can be unobserved (latent) because:
  - imaginary quantity: meant to provide some simplified and abstractive view of the data generation process
    - e.g., speech recognition models, mixture models, ...
  - a real-world object (and/or phenomena), but difficult or impossible to measure
    - e.g., the temperature of a star, causes of a disease, evolutionary ancestors ...
  - a real-world object (and/or phenomena), but sometimes wasn't measured, because of faulty sensors, etc.
- Discrete latent variables can be used to partition/cluster data into sub-groups
- Continuous latent variables (factors) can be used for dimensionality reduction (e.g., factor analysis, etc.)

# Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:

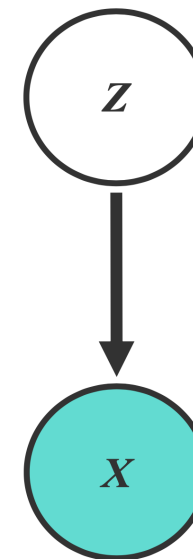
$$p(x_n | \mu, \Sigma) = \sum_k \pi_k N(x, | \mu_k, \Sigma_k)$$

mixture proportion      mixture component



- This model can be used for unsupervised clustering.
  - This model (fit by AutoClass) has been used to discover new kinds of stars in astronomical data, etc.

# Example: Gaussian Mixture Models (GMMs)



- Consider a mixture of  $K$  Gaussian components:
  - $Z$  is a latent class indicator vector:

$$p(z_n) = \text{multi}(z_n : \pi) = \prod_k (\pi_k)^{z_n^k}$$

- $X$  is a conditional Gaussian variable with a class-specific mean/covariance

$$p(x_n | z_n^k = \mathbf{1}, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$

Parameters to be learned:

- The likelihood of a sample:

$$\begin{aligned} p(x_n | \mu, \Sigma) &= \sum_k p(z^k = \mathbf{1} | \pi) p(x_n | z^k = \mathbf{1}, \mu, \Sigma) \\ &= \sum_{z_n} \prod_k \left( (\pi_k)^{z_n^k} N(x_n : \mu_k, \Sigma_k)^{z_n^k} \right) = \sum_k \pi_k N(x_n | \mu_k, \Sigma_k) \end{aligned}$$

mixture proportion

mixture component

# Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:  $p(x_n | \mu, \Sigma) = \sum_k \pi_k N(x_n | \mu_k, \Sigma_k)$
- Recall MLE for completely observed data

- Data log-likelihood: 
$$\ell(\theta; D) = \log \prod_n p(z_n, x_n) = \log \prod_n p(z_n | \pi) p(x_n | z_n, \mu, \sigma)$$

$$= \sum_n \log \prod_k \pi_k^{z_n^k} + \sum_n \log \prod_k N(x_n; \mu_k, \sigma)^{z_n^k}$$

$$= \sum_n \sum_k z_n^k \log \pi_k - \sum_n \sum_k z_n^k \frac{1}{2\sigma^2} (x_n - \mu_k)^2 + C$$

- MLE:

$$\hat{\pi}_{k,MLE} = \arg \max_{\pi} \ell(\theta; D),$$

$$\hat{\mu}_{k,MLE} = \arg \max_{\mu} \ell(\theta; D)$$

$$\hat{\sigma}_{k,MLE} = \arg \max_{\sigma} \ell(\theta; D)$$

$$\Rightarrow \hat{\mu}_{k,MLE} = \frac{\sum_n z_n^k x_n}{\sum_n z_n^k}$$

- What if we do not know  $z_n$ ?

# Why is Learning Harder?

- **Complete log likelihood:** if both  $\mathbf{x}$  and  $\mathbf{z}$  can be observed, then

$$\ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z}|\theta) = \log p(\mathbf{z}|\theta_z) + \log p(\mathbf{x}|\mathbf{z}, \theta_x)$$

- Decomposes into a sum of factors, the parameter for each factor can be estimated separately
- But given that  $\mathbf{z}$  is not observed,  $\ell_c(\theta; \mathbf{x}, \mathbf{z})$  is a random quantity, cannot be maximized directly
- **Incomplete (or marginal) log likelihood:** with  $\mathbf{z}$  unobserved, our objective becomes the log of a marginal probability:

$$\ell(\theta; \mathbf{x}) = \log p(\mathbf{x}|\theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z}|\theta)$$

- All parameters become coupled together
- In other models when  $\mathbf{z}$  is complex (continuous) variables (as we'll see later), marginalization over  $\mathbf{z}$  is intractable.



Questions?