

# DSC291: Machine Learning with Few Labels

## Self-supervised Learning

**Zhiting Hu**

Lecture 7, April 22, 2025

# Outline

- Variational Auto-Encoders (VAEs)
- Self-Supervised Learning (SSL)

# Recap: EM and Variational Inference

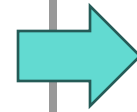
- The EM algorithm:

- E-step:  $q^{t+1} = \arg \min_q F(q, \theta^t)$

**Intractable** when  
model  $p(\mathbf{z}, \mathbf{x}|\theta)$  is  
complex

$$= p(\mathbf{z}|\mathbf{x}, \theta^t) = \frac{p(\mathbf{z}, \mathbf{x}|\theta^t)}{\sum_{\mathbf{z}} p(\mathbf{z}, \mathbf{x}|\theta^t)}$$

- M-step:  $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta)$



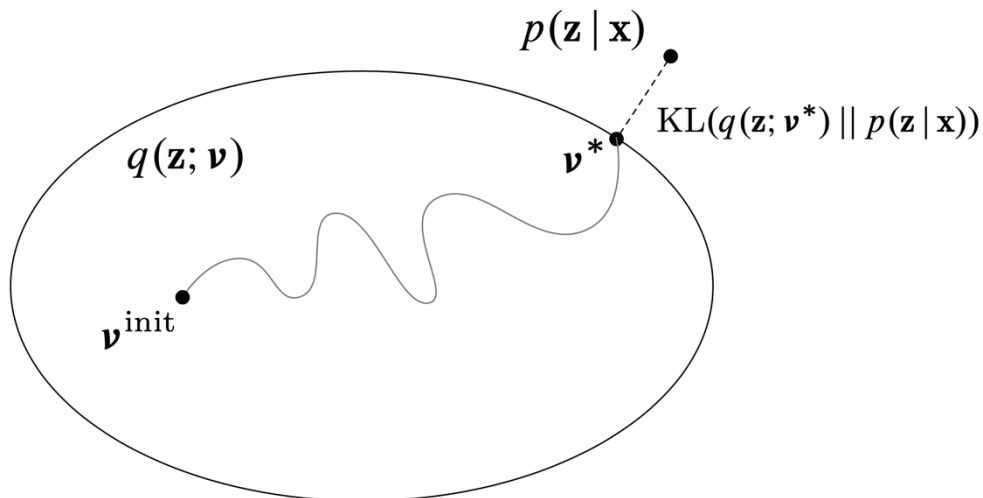
Need to approximate  $p(\mathbf{z}|\mathbf{x}, \theta^t)$   
with Variational Inference (VI)

$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

## Recap: Variational Inference

- Choose a family of distributions over the latent variables  $\mathbf{z}$  with its own set of variational parameters  $\boldsymbol{\nu}$ , i.e.  $q(\mathbf{z}|\mathbf{x}, \boldsymbol{\nu})$
- We maximize the ELBO over  $q$  to find an “optimal approximation” to  $p(\mathbf{z}|\mathbf{x})$

$$\begin{aligned} & \operatorname{argmax}_{\boldsymbol{\nu}} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \boldsymbol{\nu})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x}, \boldsymbol{\nu})} \right] \\ &= \operatorname{argmax}_{\boldsymbol{\nu}} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \boldsymbol{\nu})} [\log p(\mathbf{x}, \mathbf{z}|\theta)] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \boldsymbol{\nu})} [\log q(\mathbf{z}|\mathbf{x}, \boldsymbol{\nu})] \end{aligned}$$

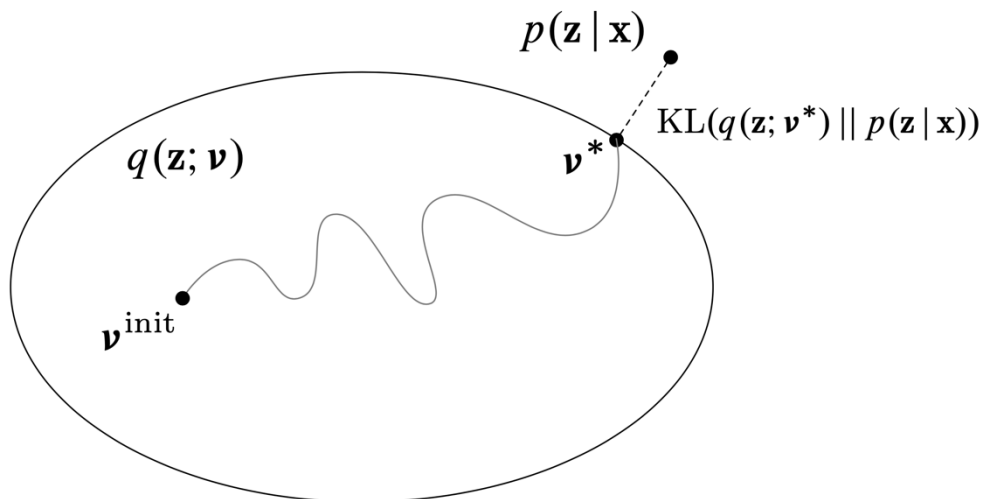


$$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$$

## Recap: Variational Inference

- Choose a family of distributions over the latent variables  $\mathbf{z}$  with its own set of variational parameters  $\mathbf{v}$ , i.e.  $q(\mathbf{z}|\mathbf{x}, \mathbf{v})$
- We maximize the ELBO over  $q$  to find an “optimal approximation” to  $p(\mathbf{z}|\mathbf{x})$

$$\begin{aligned} & \operatorname{argmax}_{\mathbf{v}} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} \right] \\ &= \operatorname{argmax}_{\mathbf{v}} \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} [\log p(\mathbf{x}, \mathbf{z}|\theta)] - \mathbb{E}_{q(\mathbf{z}|\mathbf{x}, \mathbf{v})} [\log q(\mathbf{z}|\mathbf{x}, \mathbf{v})] \end{aligned}$$



**Question:** How do we choose the variational family  $q(\mathbf{z}|\mathbf{x}, \mathbf{v})$ ?

- Factorized distribution -> mean field VI
- Mixture of Gaussian distribution -> black-box VI
- Neural-based distribution -> Variational Autoencoders (VAEs)

# Variational Auto-Encoders (VAEs)

- Model  $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})$ 
  - $p_{\theta}(\mathbf{x}|\mathbf{z})$ : a.k.a., generative model, generator, (probabilistic) decoder, ...
  - $p(\mathbf{z})$ : prior, e.g., Gaussian
- Assume variational distribution  $q_{\phi}(\mathbf{z}|\mathbf{x})$ 
  - E.g., a Gaussian distribution parameterized as **deep neural networks**
  - a.k.a, recognition model, inference network, (probabilistic) encoder, ...
- ELBO:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x}) &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z})] + H(q_{\phi}(\mathbf{z}|\mathbf{x})) \\ &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \text{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}))\end{aligned}$$

Reconstruction

Divergence from prior  
(KL divergence between two Gaussians has  
an analytic form)

# Variational Auto-Encoders (VAEs)

- ELBO:
$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \boldsymbol{x}) &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} [\log p_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{z})] + H(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})) \\ &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} [\log p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z})] - \text{KL}(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x}) || p(\boldsymbol{z}))\end{aligned}$$

$$\nabla_{\boldsymbol{\phi}} \mathcal{L} =$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L} =$$

# Variational Auto-Encoders (VAEs)

- ELBO:
$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \boldsymbol{x}) &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} [\log p_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{z})] + H(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})) \\ &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} [\log p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z})] - \text{KL}(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x}) || p(\boldsymbol{z}))\end{aligned}$$

$$\nabla_{\boldsymbol{\phi}} \mathcal{L} =$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L} = \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{z})]$$



# Variational Auto-Encoders (VAEs)

- ELBO:
$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \boldsymbol{x}) &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} [\log p_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{z})] + H(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})) \\ &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} [\log p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z})] - \text{KL}(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x}) || p(\boldsymbol{z}))\end{aligned}$$

$$\nabla_{\boldsymbol{\phi}} \mathcal{L} =$$

- Reparameterization:
  - $[\boldsymbol{\mu}; \boldsymbol{\sigma}] = f_{\boldsymbol{\phi}}(\boldsymbol{x})$  (a neural network)
  - $\boldsymbol{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{1})$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L} = \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{z})]$$

# Variational Auto-Encoders (VAEs)

- ELBO:

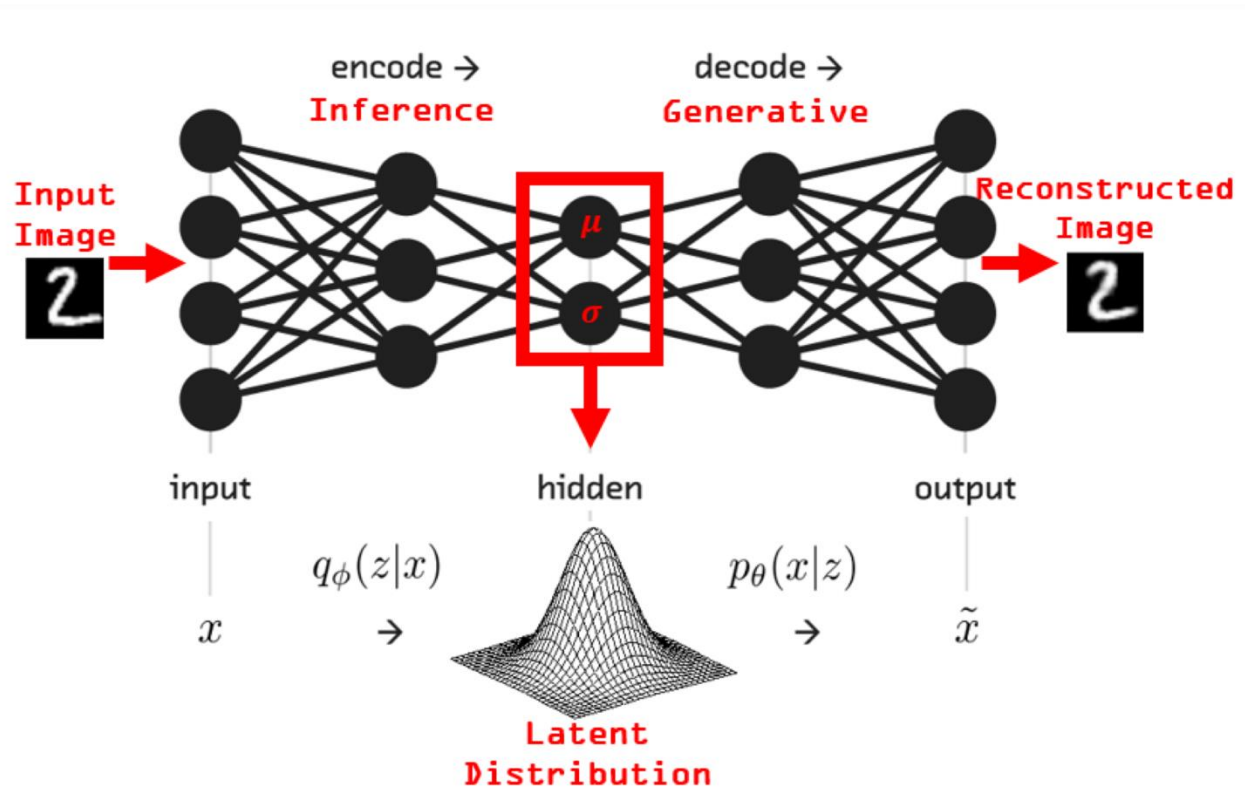
$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \boldsymbol{x}) &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} [\log p_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{z})] + H(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})) \\ &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} [\log p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z})] - \text{KL}(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x}) || p(\boldsymbol{z}))\end{aligned}$$

$$\nabla_{\boldsymbol{\phi}} \mathcal{L} = \mathbb{E}_{\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{1})} [\nabla_{\boldsymbol{z}} [\log p_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{z}) - \log q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})] \nabla_{\boldsymbol{\phi}} \boldsymbol{z}(\boldsymbol{\epsilon}, \boldsymbol{\phi})]$$

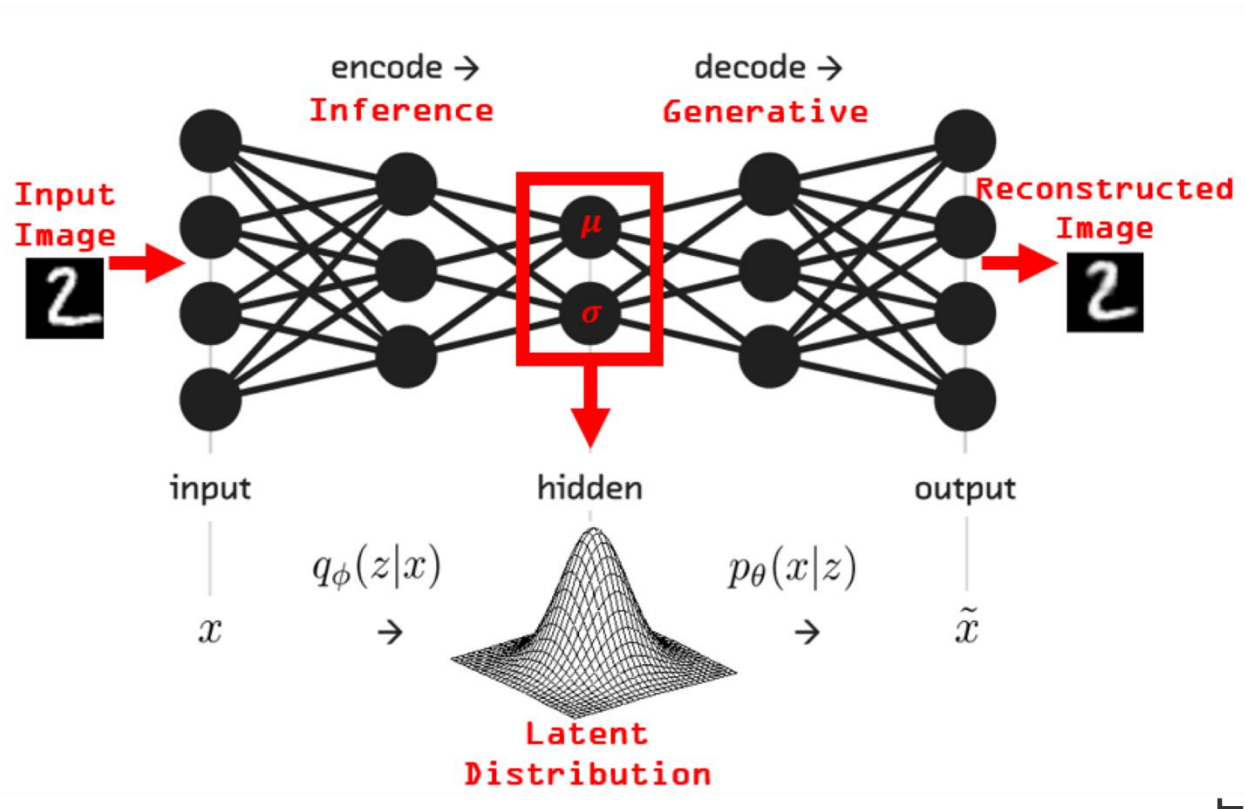
- Reparameterization:
  - $[\boldsymbol{\mu}; \boldsymbol{\sigma}] = f_{\boldsymbol{\phi}}(\boldsymbol{x})$  (a neural network)
  - $\boldsymbol{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{1})$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L} = \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{z})]$$

# Example: VAEs for images



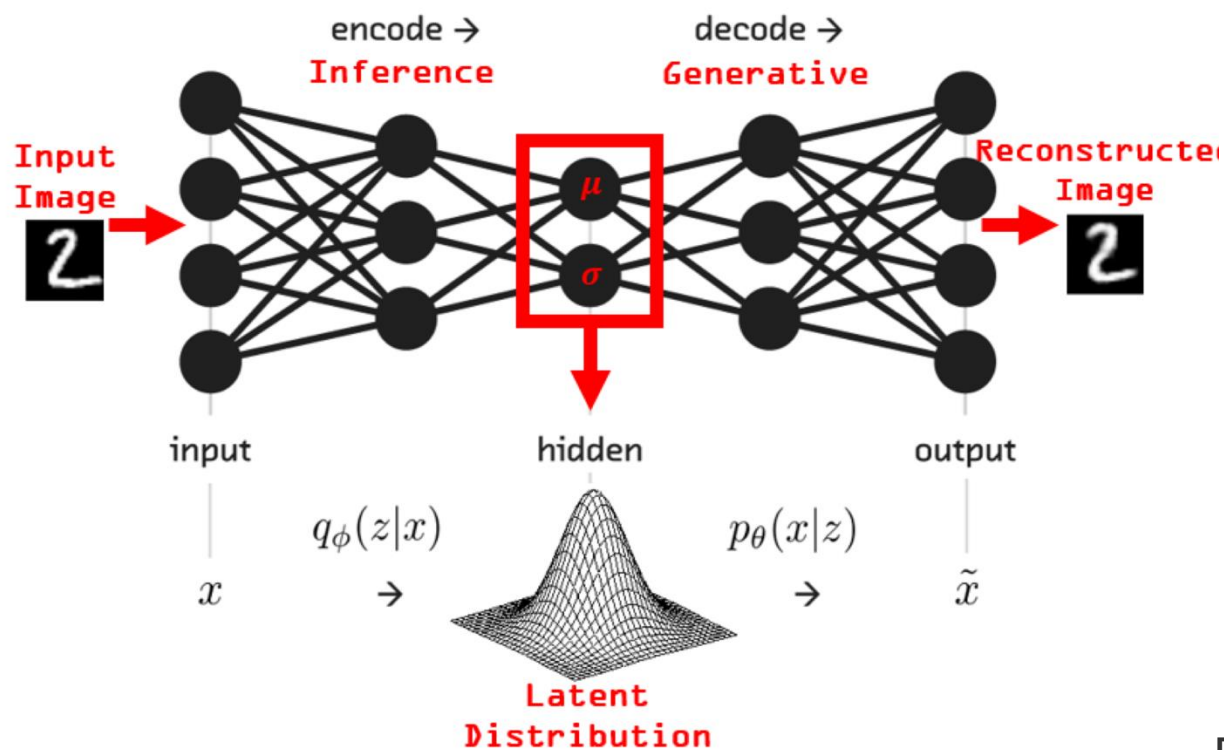
# Example: VAEs for images



Input Data

$\tilde{x}$

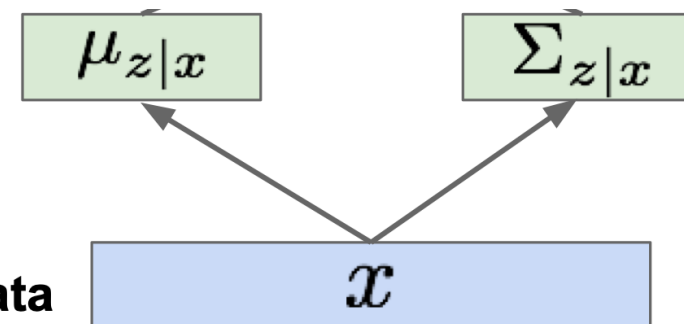
# Example: VAEs for images



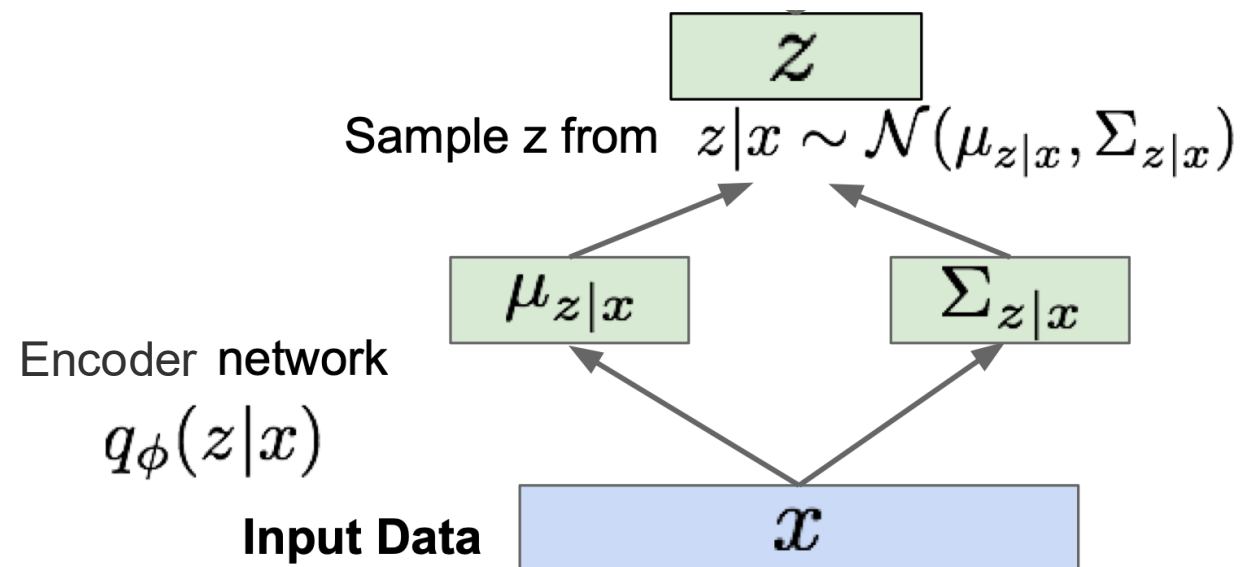
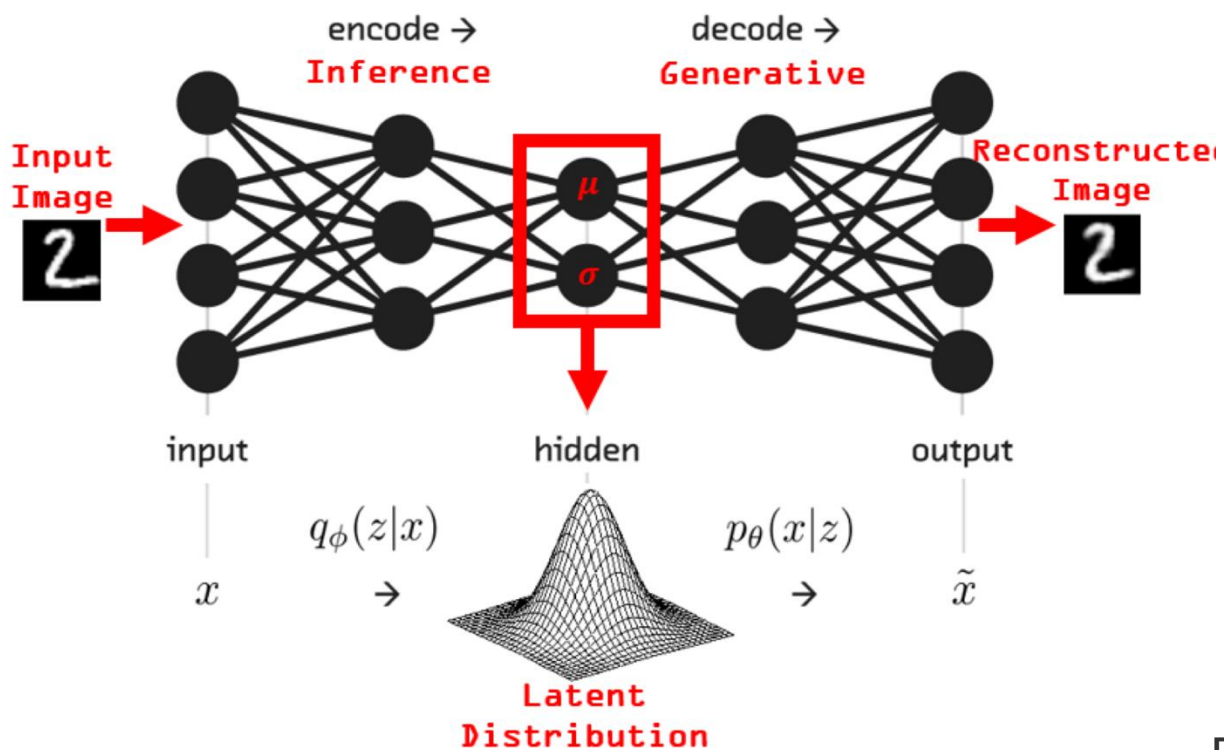
Encoder network

$$q_\phi(z|x)$$

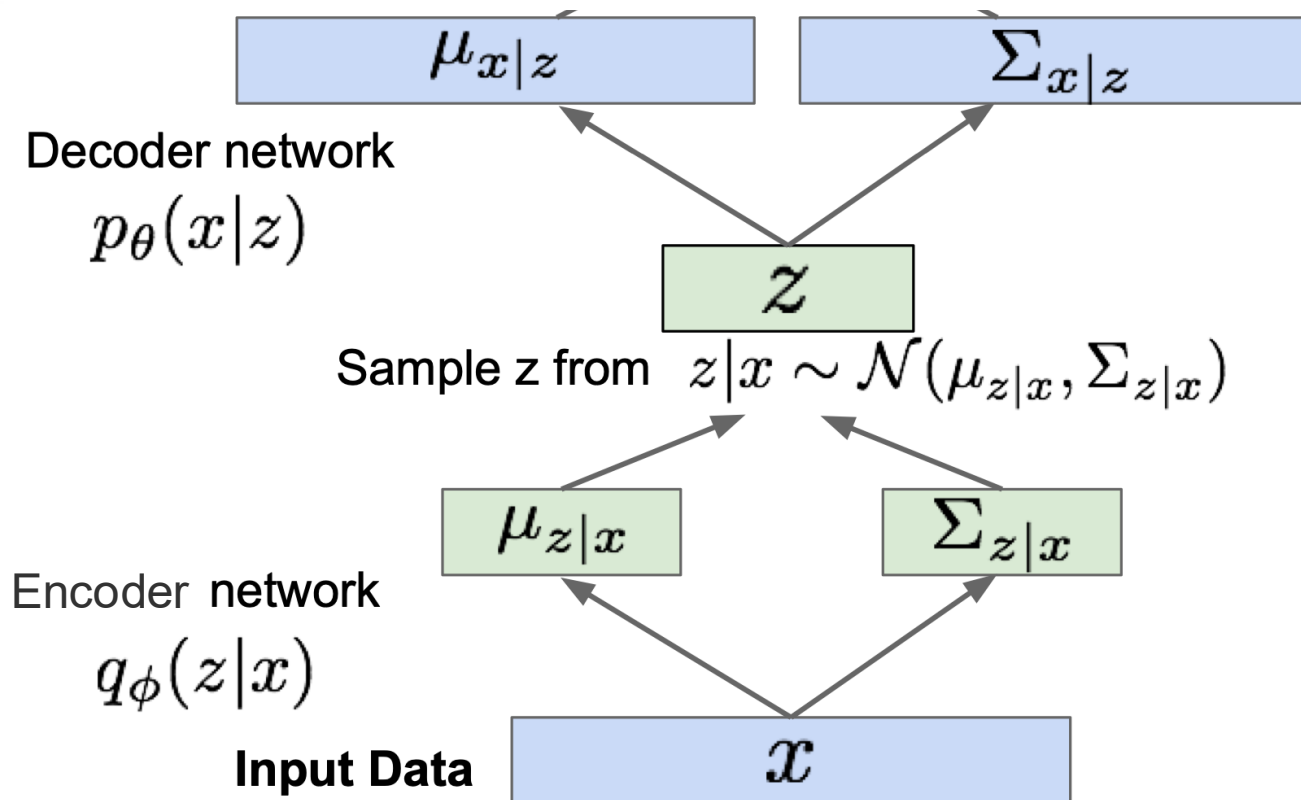
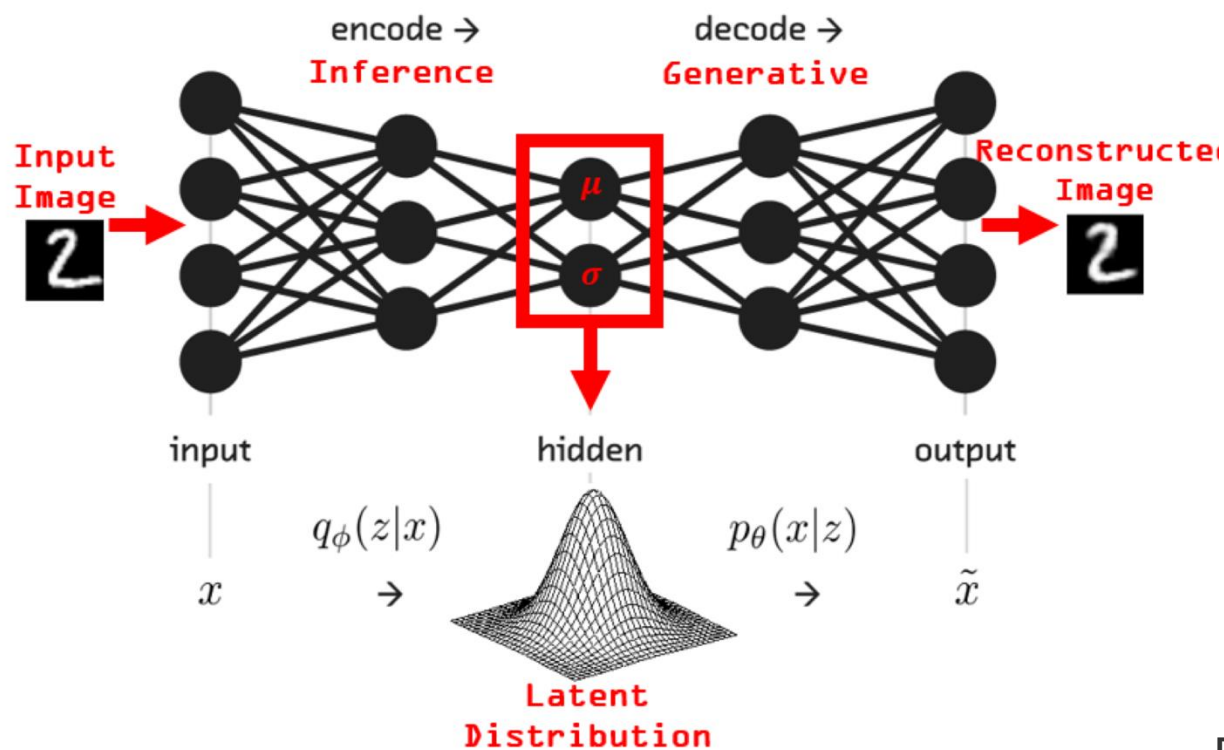
Input Data



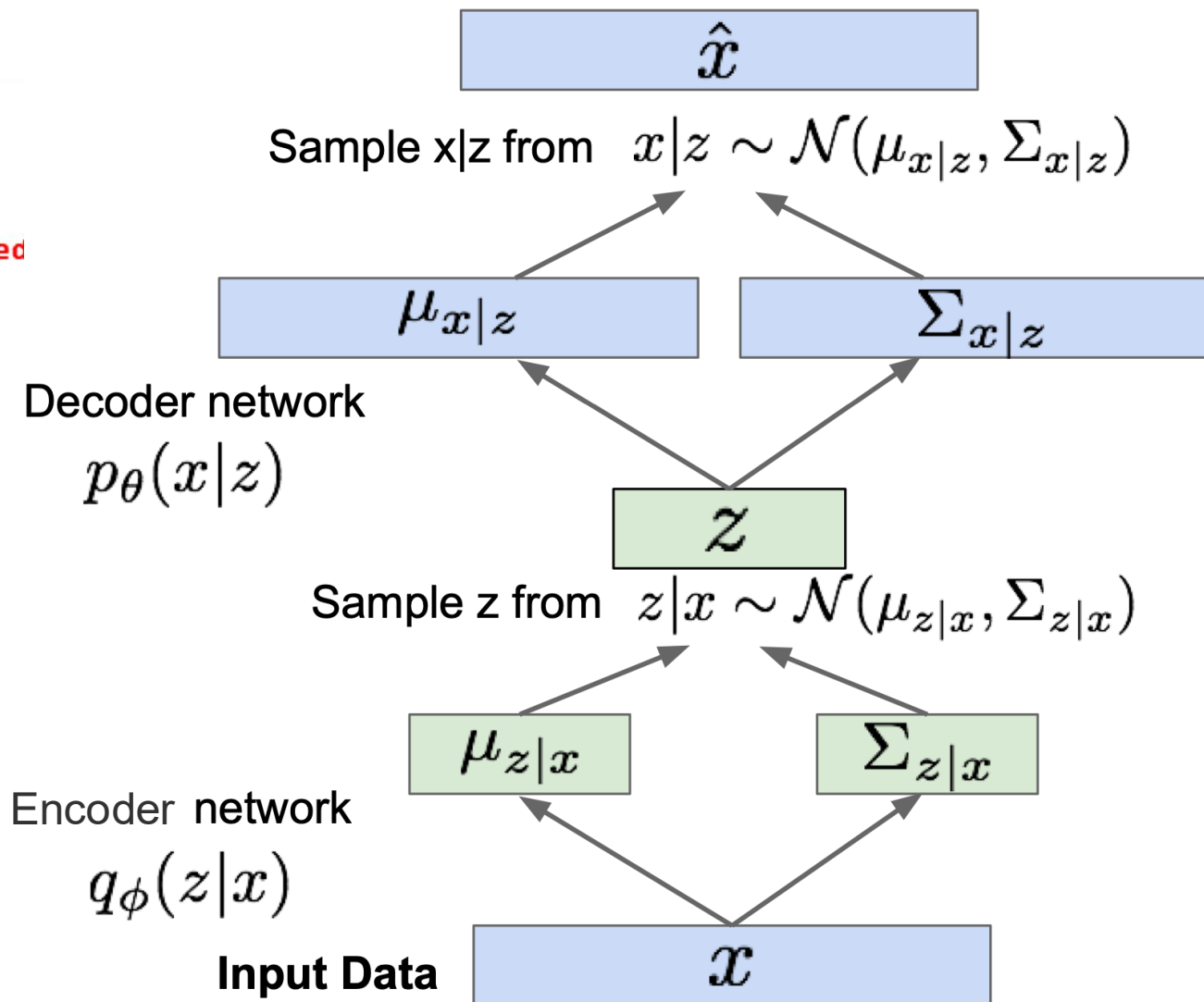
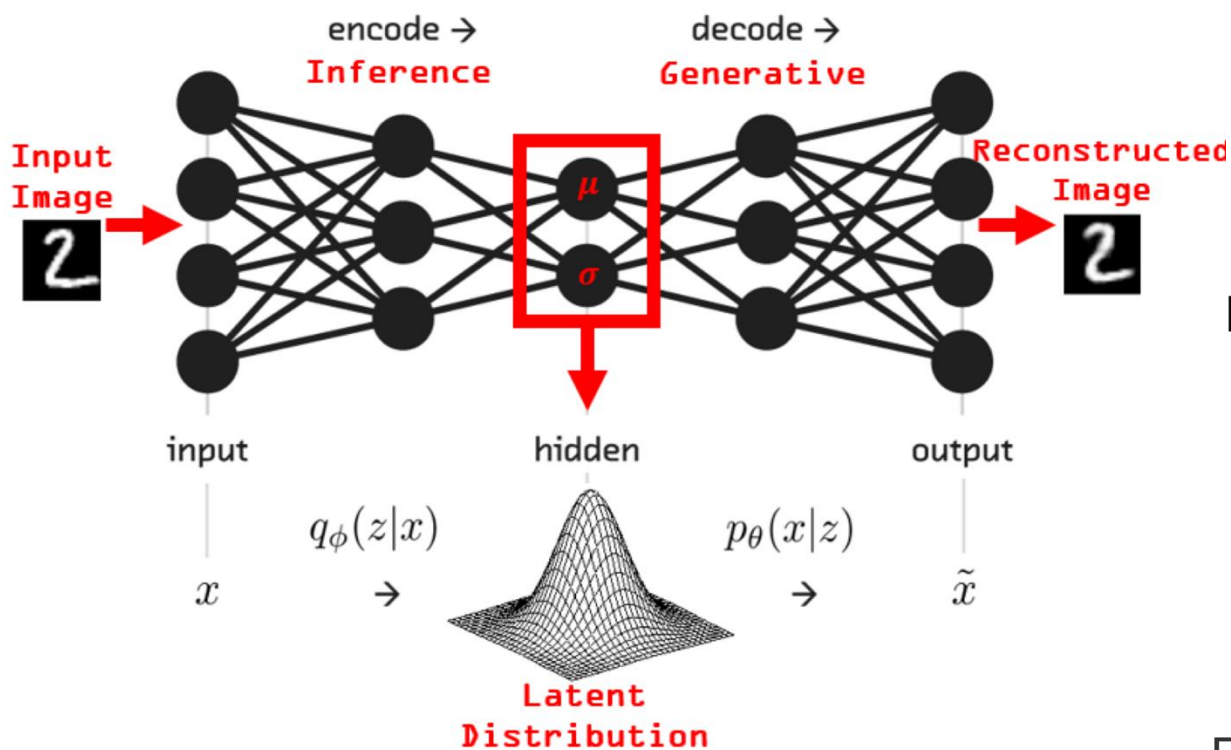
# Example: VAEs for images



# Example: VAEs for images



# Example: VAEs for images

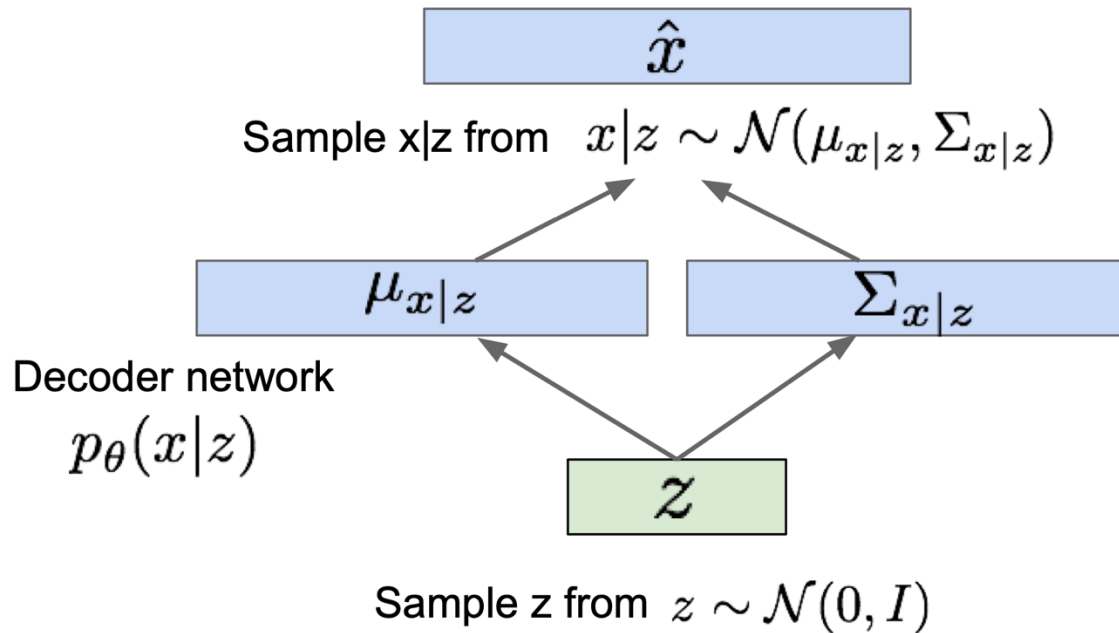




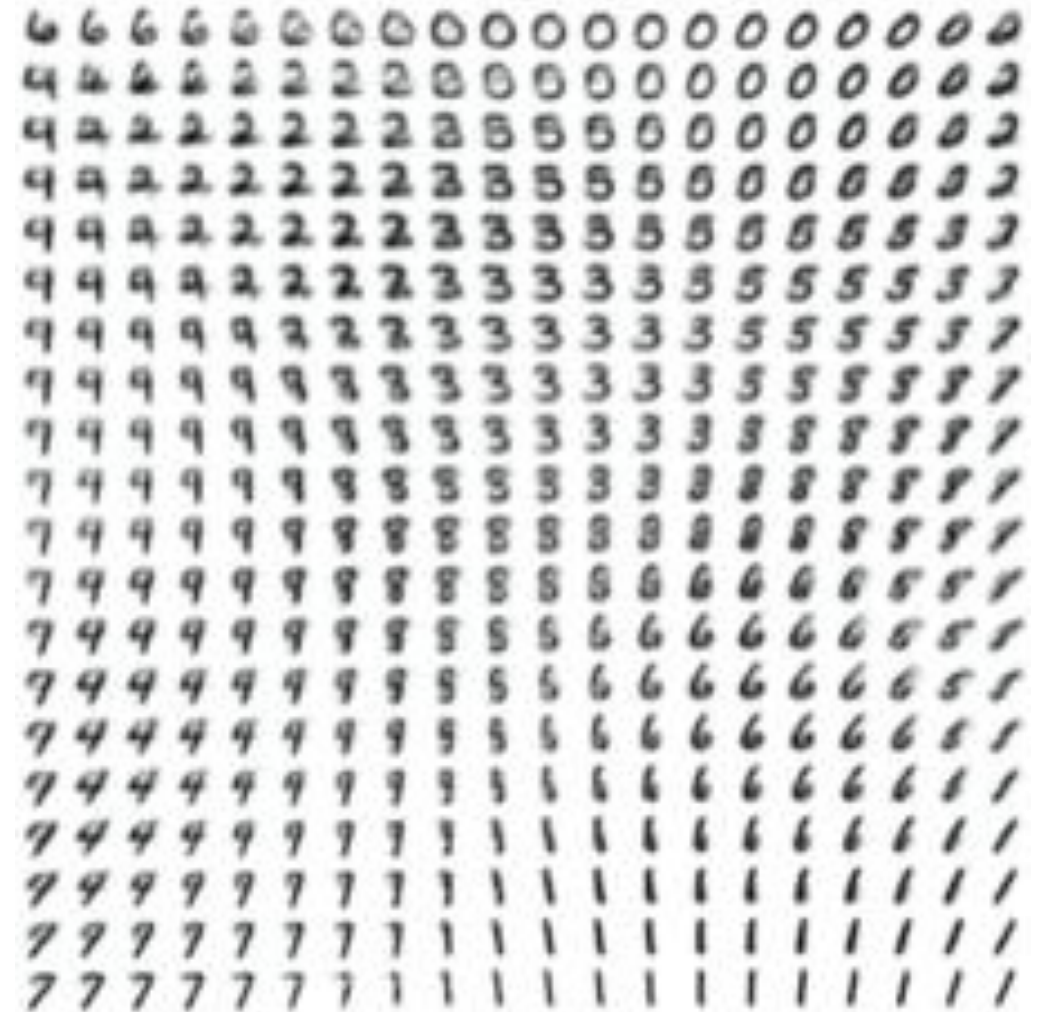
# Example: VAEs for images

Generating samples:

- Use decoder network. Now sample  $z$  from prior!



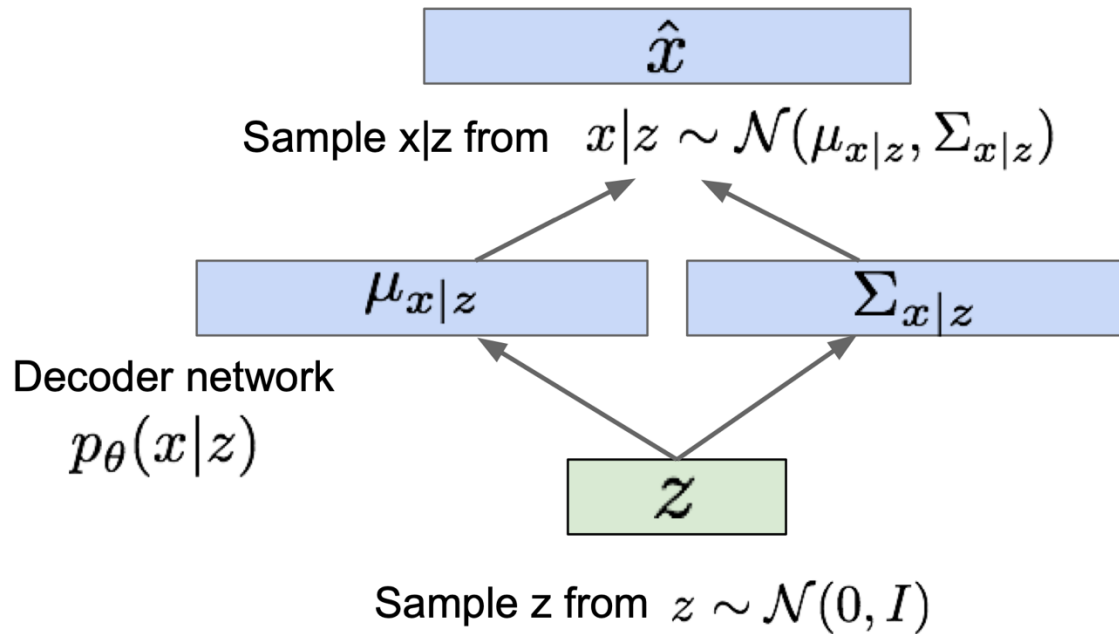
Data manifold for 2-d  $z$



# Example: VAEs for images

Generating samples:

- Use decoder network. Now sample  $z$  from prior!



Data manifold for 2-d  $z$

Vary  $z_1$   
(Degree of smile)



Vary  $z_2$  (head pose)

## Example: VAEs for text

- Latent code interpolation and sentences generation from VAEs [Bowman et al., 2015].

---

**“ i want to talk to you . ”**

*“i want to be with you . ”*

*“i do n’t want to be with you . ”*

*i do n’t want to be with you .*

**she did n’t want to be with him .**

---

# Variational Auto-encoders: Summary

- A combination of the following ideas:
  - Variational Inference: ELBO
  - Variational distribution parametrized as neural networks
  - Reparameterization trick

$$\mathcal{L}(\theta, \phi; x) = [\log p_{\theta}(x|z)] - \text{KL}(q_{\phi}(z|x) || p(z))$$

Reconstruction

Divergence from prior

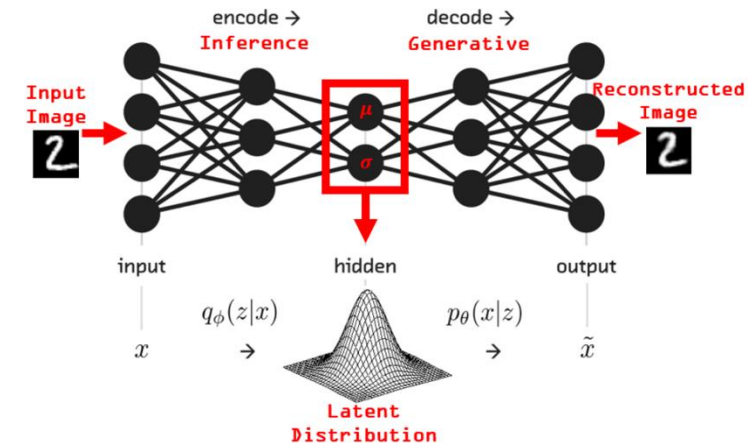
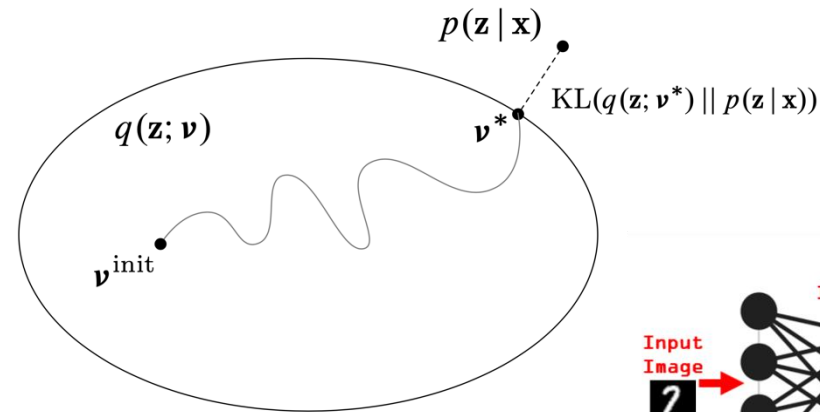


(Razavi et al., 2019)

- Pros:
  - Principled approach to generative models
  - Allows inference of  $q(z|x)$ , can be useful feature representation for other tasks
- Cons:
  - Samples blurrier and lower quality compared to GANs
  - Tend to collapse on text data

# Summary: Supervised / Unsupervised Learning

- Supervised Learning
  - Maximum likelihood estimation (MLE)
- Unsupervised learning
  - Maximum likelihood estimation (MLE) with latent variables
    - Marginal log-likelihood
  - EM algorithm for MLE
    - ELBO / Variational free energy
  - Variational Inference
    - ELBO / Variational free energy
    - Variational distributions
      - Factorized (mean-field VI)
      - Mixture of Gaussians (Black-box VI)
      - Neural-based (VAEs)



# Self-Supervised Learning

# “X”-supervised learning

- Supervised learning
- Unsupervised learning
- Self-supervised learning
- Weakly-/distantly-supervised learning
- Semi-supervised learning
- ...

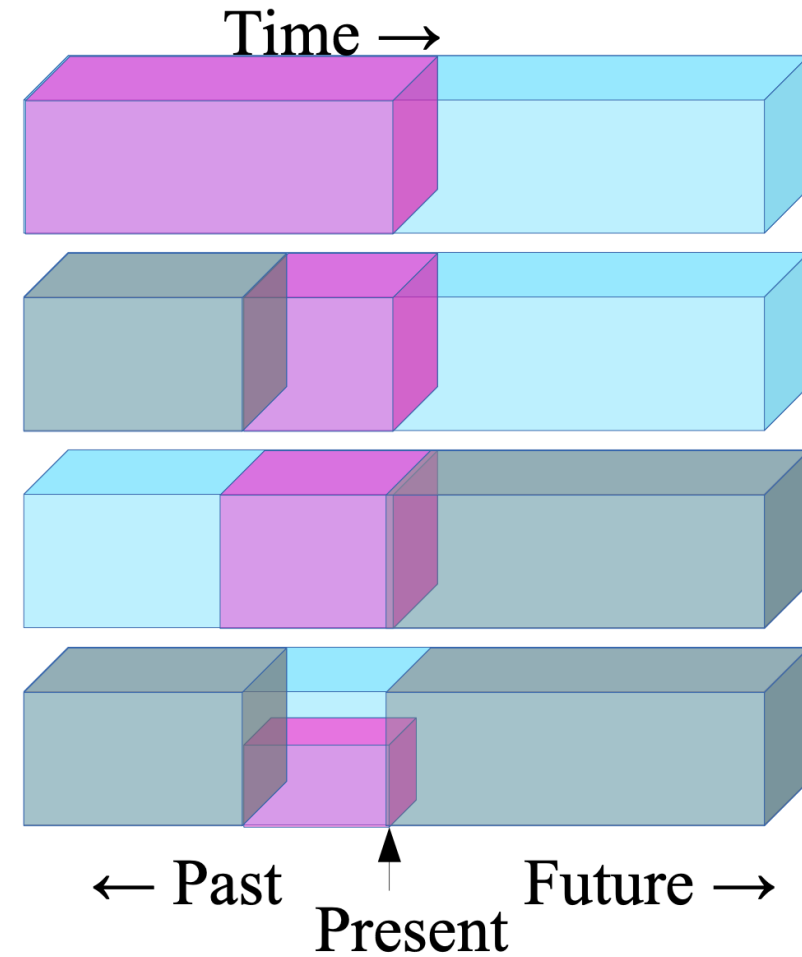
# Self-Supervised Learning

- Given an observed data instance  $t$
- One could derive various supervision signals based on the structure of the data
- By applying a “split” function that artificially partition  $t$  into two parts
  - $(x, y) = \text{split}(t)$
  - sometimes split in a stochastic way
- Treat  $x$  as the input and  $y$  as the output
- Train a model  $p_{\theta}(y|x)$



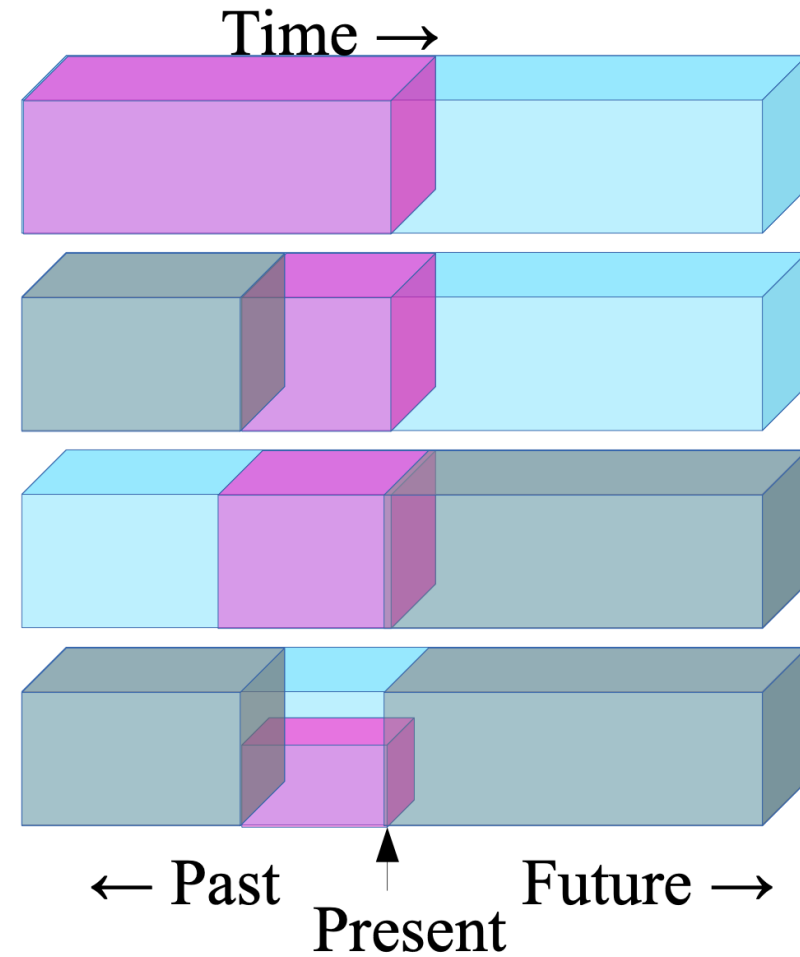
# Self-Supervised Learning: Examples

- ▶ 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**.



# Self-Supervised Learning: Examples

- ▶ 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.**



# Self-Supervised Learning: Motivation (I)

- ▶ Our brains do this all the time
- ▶ Filling in the visual field at the retinal blind spot
- ▶ Filling in occluded images, missing segments in speech
- ▶ Predicting the state of the world from partial (textual) descriptions
- ▶ Predicting the consequences of our actions
- ▶ Predicting the sequence of actions leading to a result
- ▶ **Predicting any part of the past, present or future percepts from whatever information is available.**



# Self-Supervised Learning: Motivation (I)

- Successfully learning to predict everything from everything else would result in **the accumulation of lots of background knowledge about how the world works**
- The model is forced to learn what we really care about, e.g. a semantic representation, in order to solve the prediction problem

[Courtesy: Lecun “Self-supervised Learning”]

[Courtesy: Zisserman “Self-supervised Learning”]

# Self-Supervised Learning: Motivation (II)

- The machine predicts any part of its input from any observed part
  - **A lot of** supervision signals in each data instance
- Untapped/availability of vast numbers of unlabeled text/images/videos..
  - Facebook: one billion images uploaded per day
  - 300 hours of video are uploaded to YouTube every minute

# Self-Supervised Learning (SSL): Examples

- SSL from text
- SSL from images
- SSL from videos

# Self-Supervised Learning from Text

Examples:

- Language models
- Learning contextual text representations

# Language Models

- Calculates the probability of a sentence:
  - Sentence:

$$\mathbf{y} = (y_1, y_2, \dots, y_T)$$

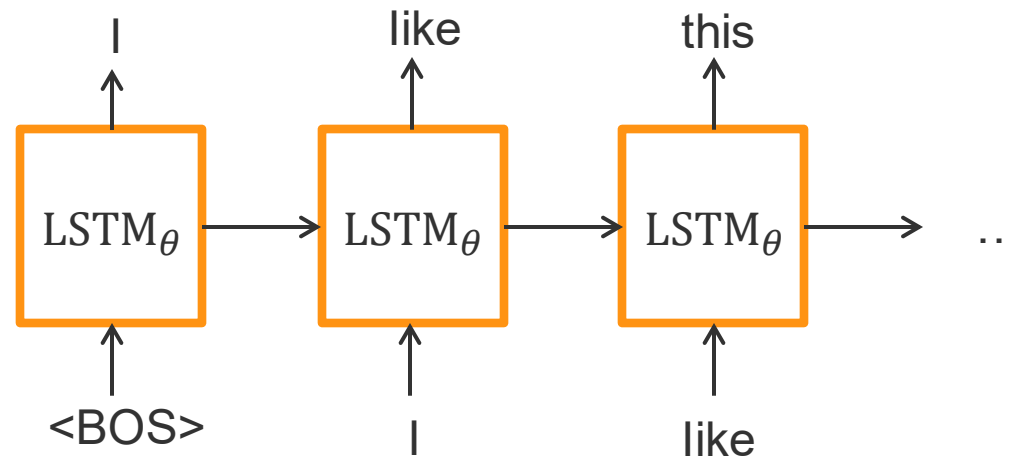
$$p_{\theta}(\mathbf{y}) = \prod_{t=1}^T p_{\theta}(y_t \mid \mathbf{y}_{1:t-1})$$

Example:

*(I, like, this, ...)*

$\dots p_{\theta}(\text{like} \mid I) p_{\theta}(\text{this} \mid I, \text{like}) \dots$

Model: LSTM RNN





# Language Models

- Calculates the probability of a sentence:
  - Sentence:

$$\mathbf{y} = (y_1, y_2, \dots, y_T)$$

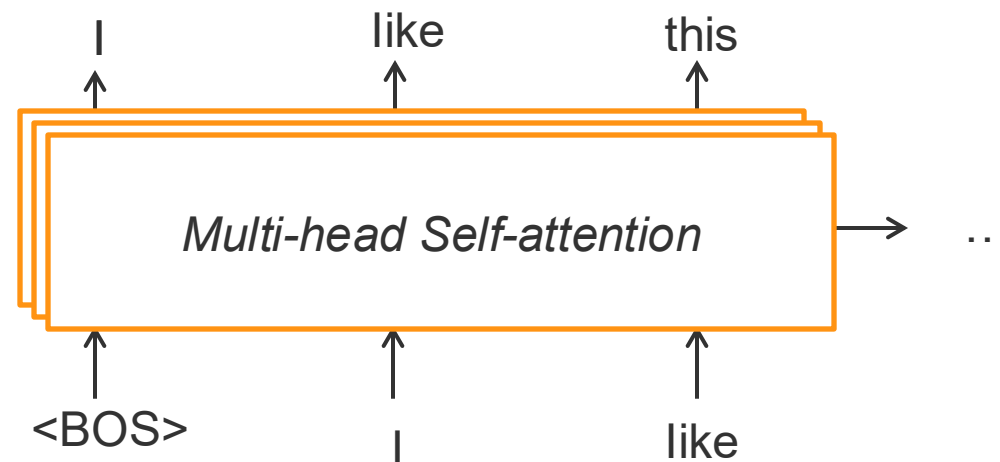
$$p_{\theta}(\mathbf{y}) = \prod_{t=1}^T p_{\theta}(y_t \mid \mathbf{y}_{1:t-1})$$

Example:

*(I, like, this, ...)*

$\dots p_{\theta}(\text{like} \mid I) p_{\theta}(\text{this} \mid I, \text{like}) \dots$

Model: Transformer



# Language Models: Training

- Given data example  $\mathbf{y}^*$
- Minimizes negative log-likelihood of the data

$$\min_{\theta} \mathcal{L}(\theta) = -\log p_{\theta}(\mathbf{y}^*) = -\prod_{t=1}^T p_{\theta}(y_t^* \mid \mathbf{y}_{1:t-1}^*)$$

- Next word prediction

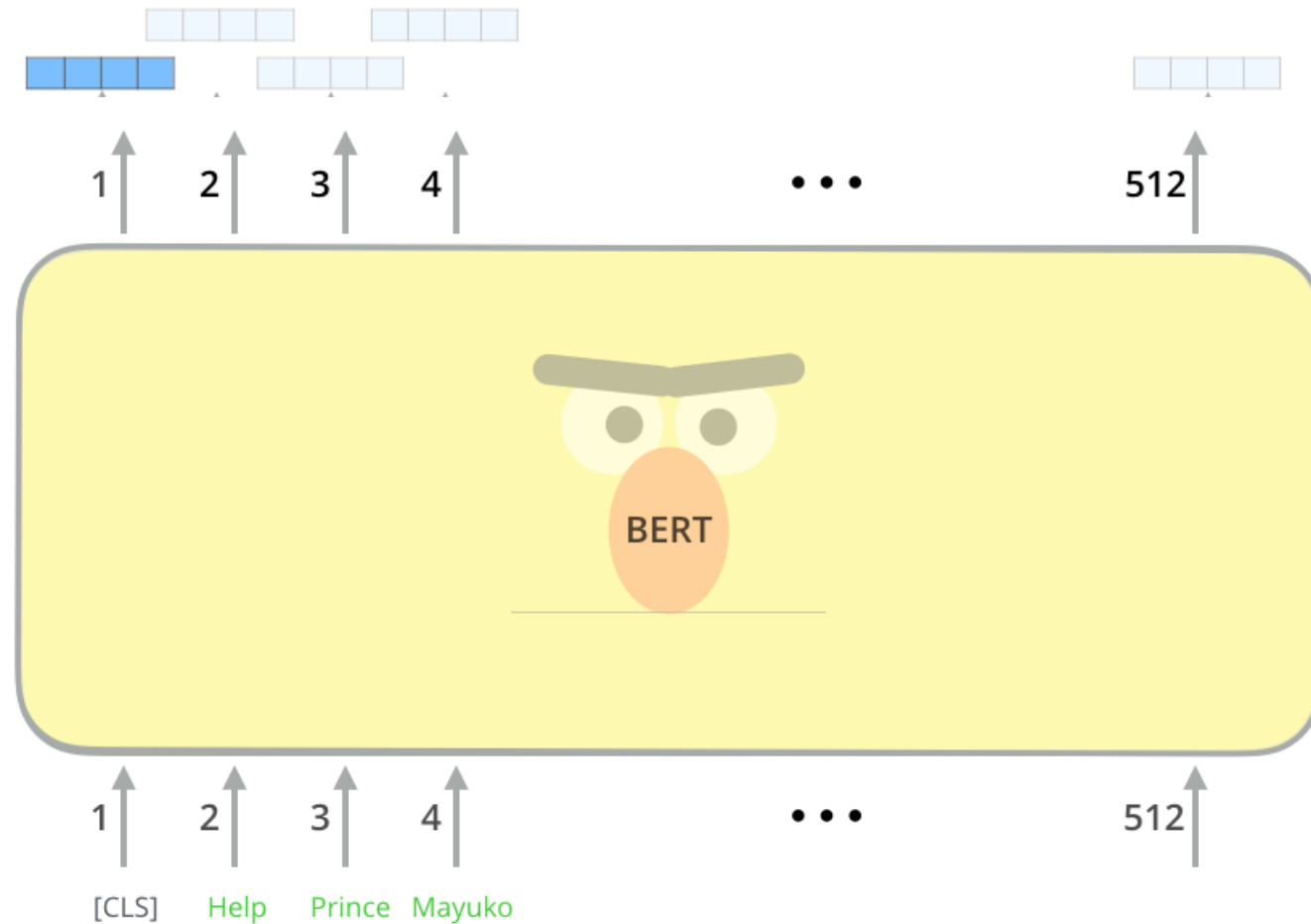
# Self-Supervised Learning from Text

Examples:

- Language models
- Learning contextual text representations

# BERT

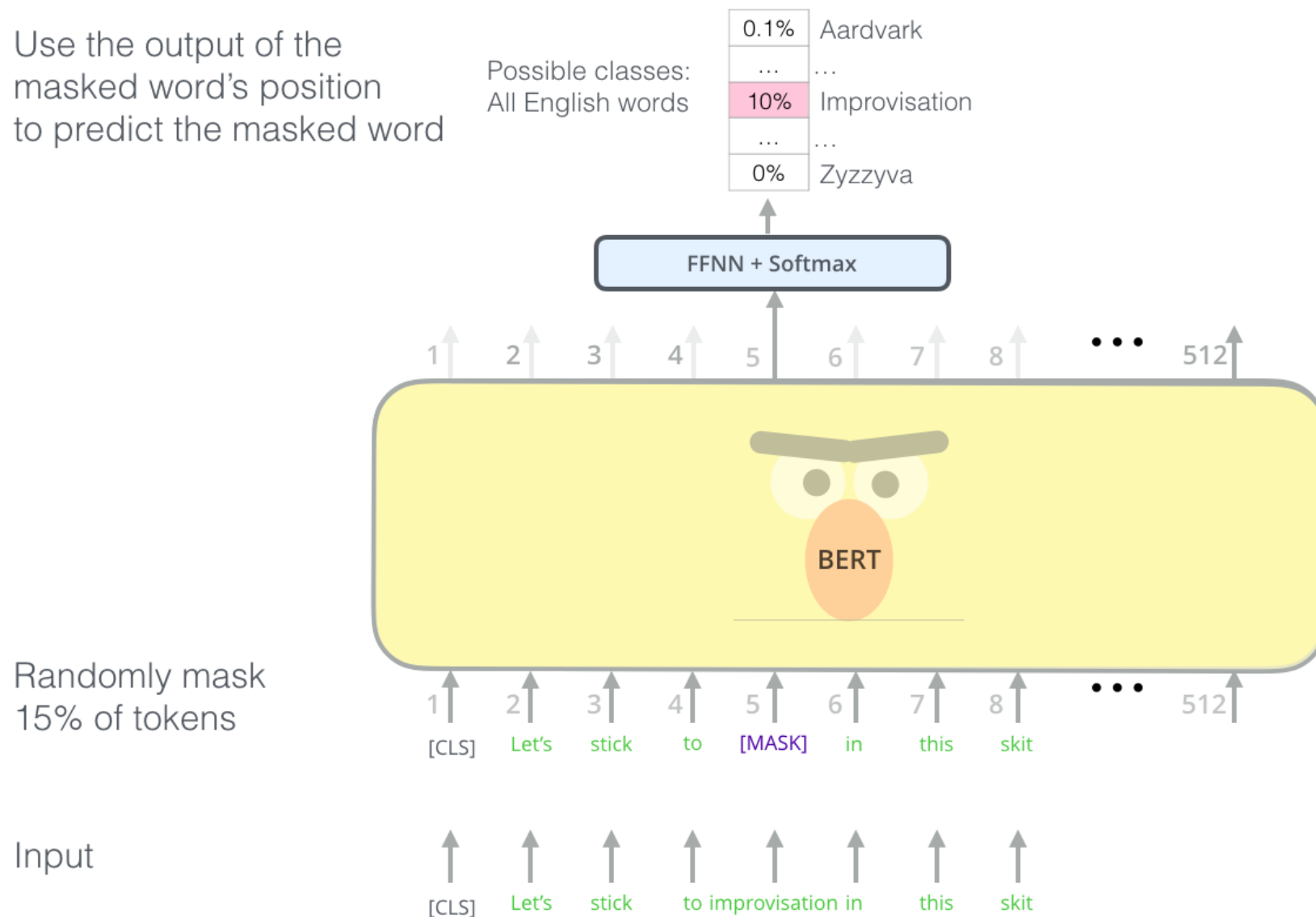
- BERT: A bidirectional model to extract contextual word embedding



# BERT: Pre-training with Self-supervised Learning

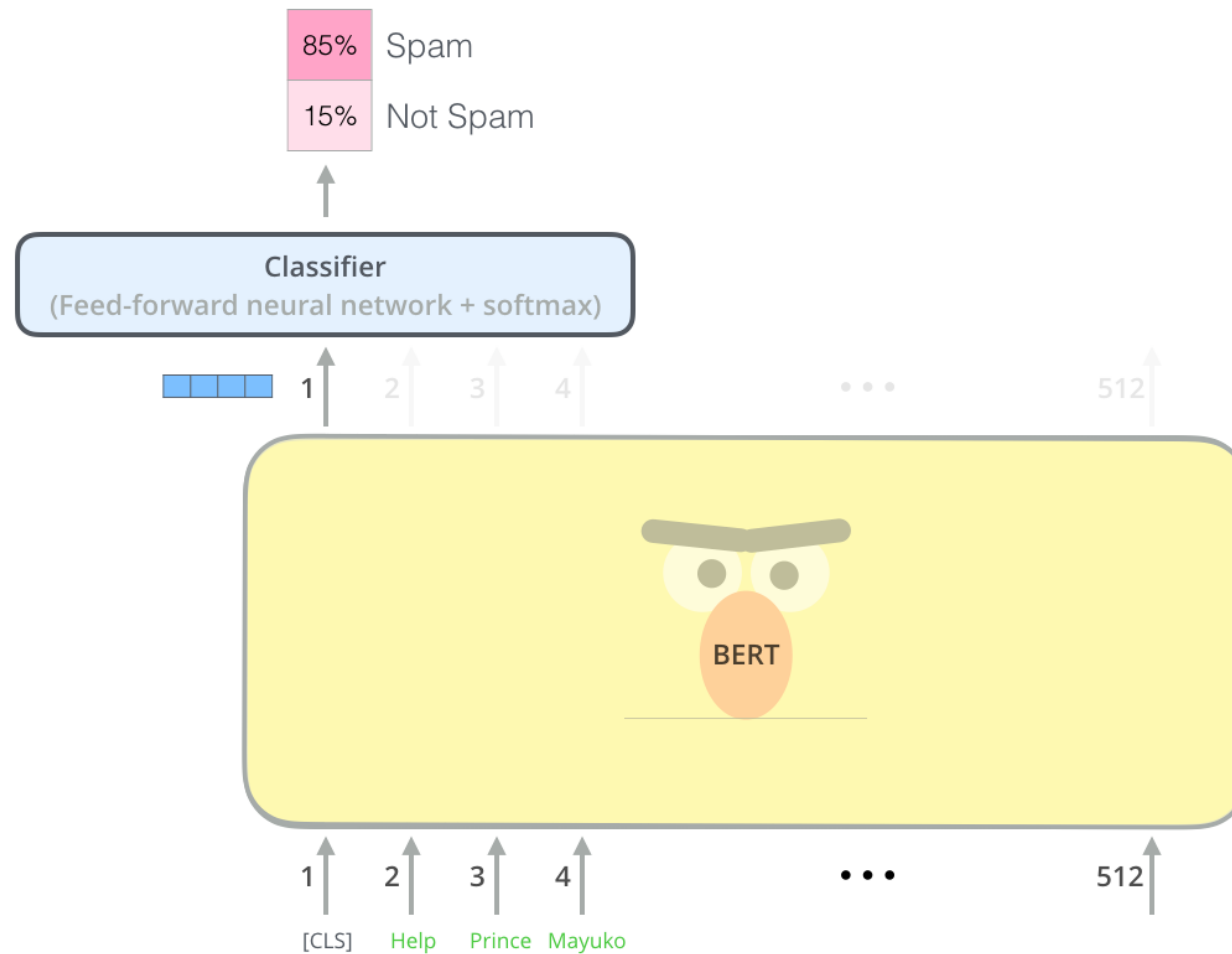
- Masked LM

Use the output of the masked word's position to predict the masked word



# BERT: Downstream Fine-tuning

- Use BERT for sentence classification



# BERT Results

- Huge improvements over SOTA on 12 NLP task

System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT <sub>BASE</sub>	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT <sub>LARGE</sub>	<b>86.7/85.9</b>	<b>72.1</b>	<b>91.1</b>	<b>94.9</b>	<b>60.5</b>	<b>86.5</b>	<b>89.3</b>	<b>70.1</b>	<b>81.9</b>

Table 1: GLUE Test results, scored by the GLUE evaluation server. The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set. OpenAI GPT = (L=12, H=768, A=12); BERT<sub>BASE</sub> = (L=12, H=768, A=12); BERT<sub>LARGE</sub> = (L=24, H=1024, A=16). BERT and OpenAI GPT are single-model, single task. All results obtained from <https://gluebenchmark.com/leaderboard> and <https://blog.openai.com/language-unsupervised/>.

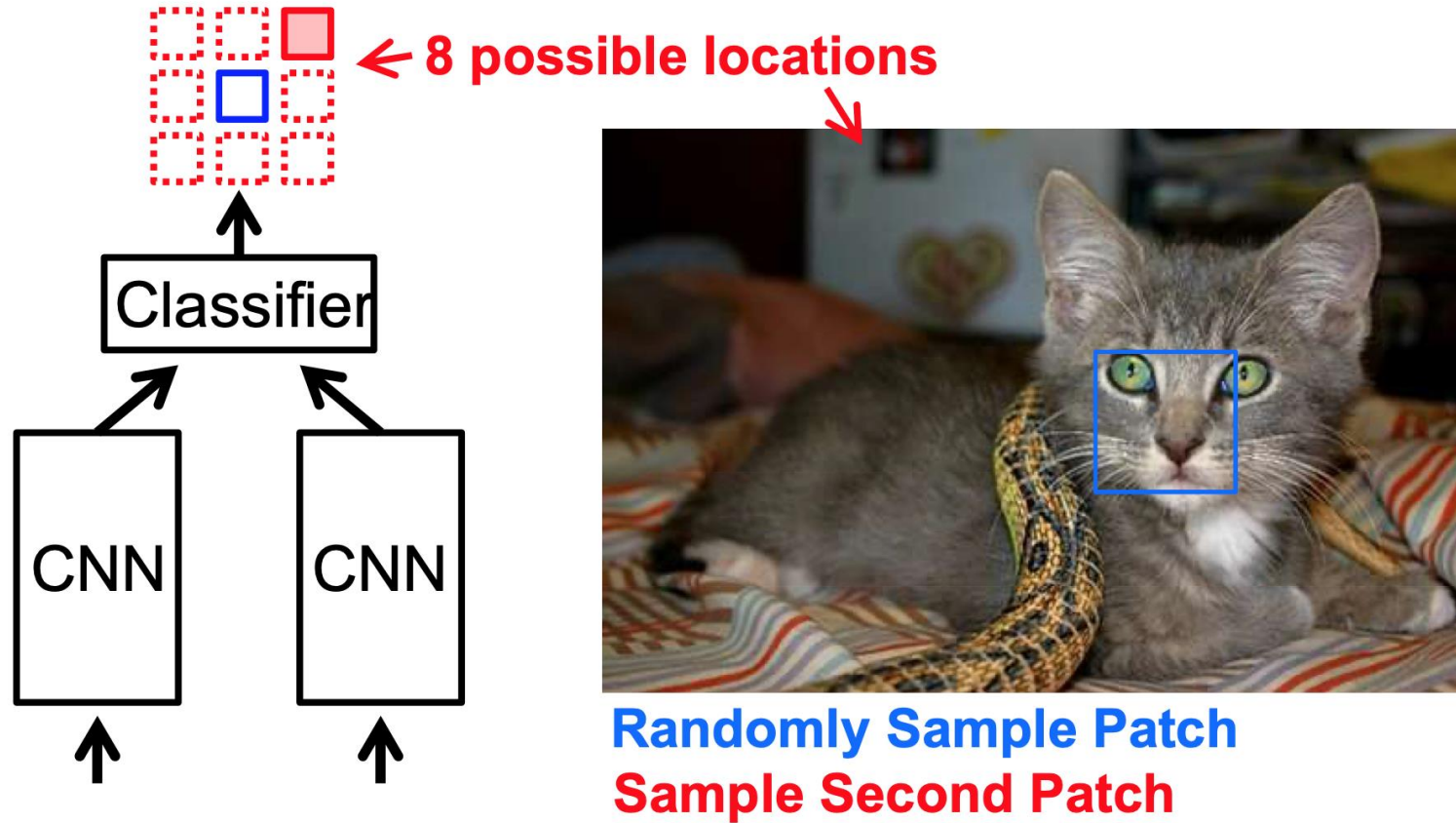
# Self-Supervised Learning (SSL): Examples

- SSL from text
- SSL from images
- SSL from videos



# SSL from Images, EX (I): relative positioning

Train network to predict relative position of two regions in the same image



Unsupervised visual representation learning by context prediction,  
Carl Doersch, Abhinav Gupta, Alexei A. Efros, ICCV 2015

## SSL from Images, EX (I): relative positioning



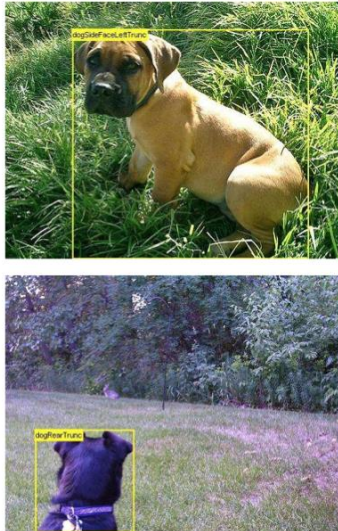
Unsupervised visual representation learning by context prediction,  
Carl Doersch, Abhinav Gupta, Alexei A. Efros, ICCV 2015

# SSL from Images, EX (I): relative positioning

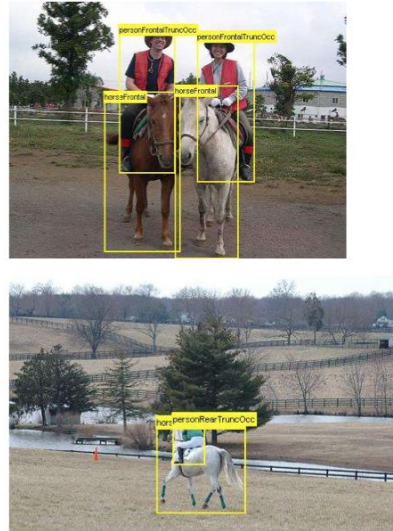
## Evaluation: PASCAL VOC Detection

- 20 object classes (car, bicycle, person, horse ...)
- Predict the bounding boxes of all objects of a given class in an image (if any)

Dog



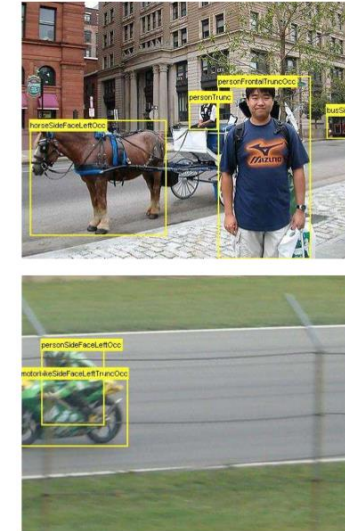
Horse



Motorbike



Person



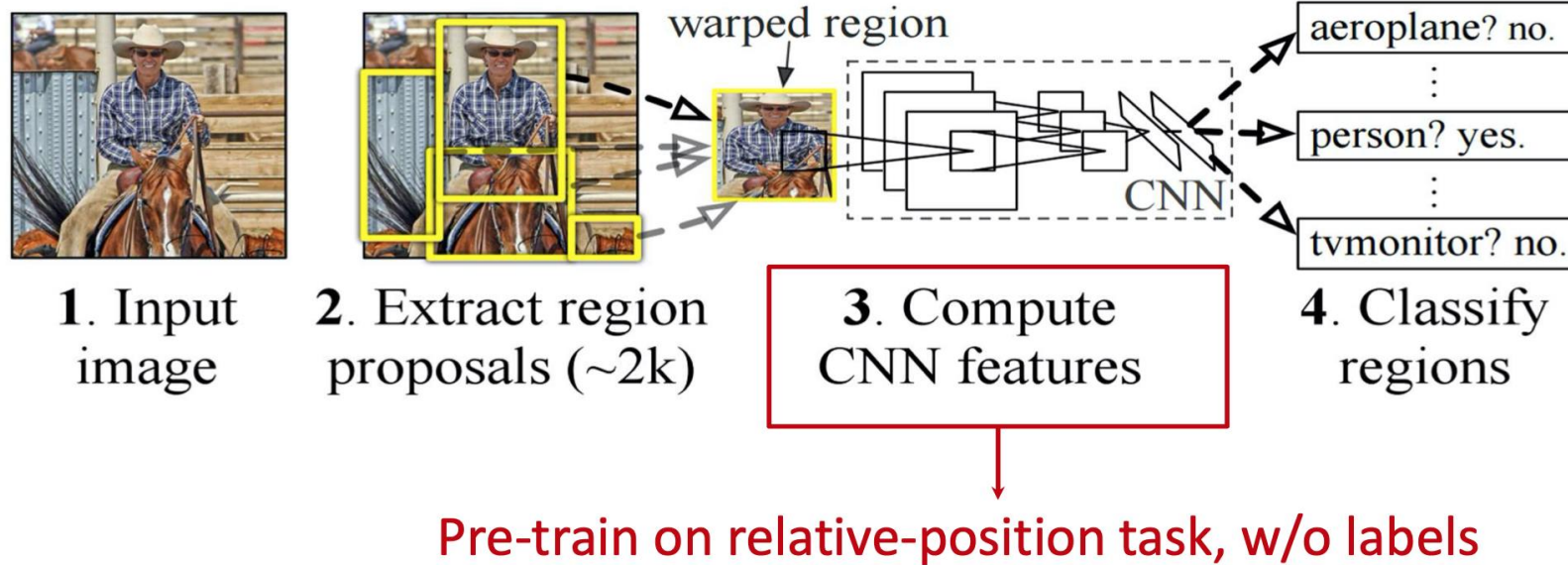


# SSL from Images, EX (I): relative positioning

## Evaluation: PASCAL VOC Detection

- Pre-train CNN using self-supervision (no labels)
- Train CNN for detection in R-CNN object category detection pipeline

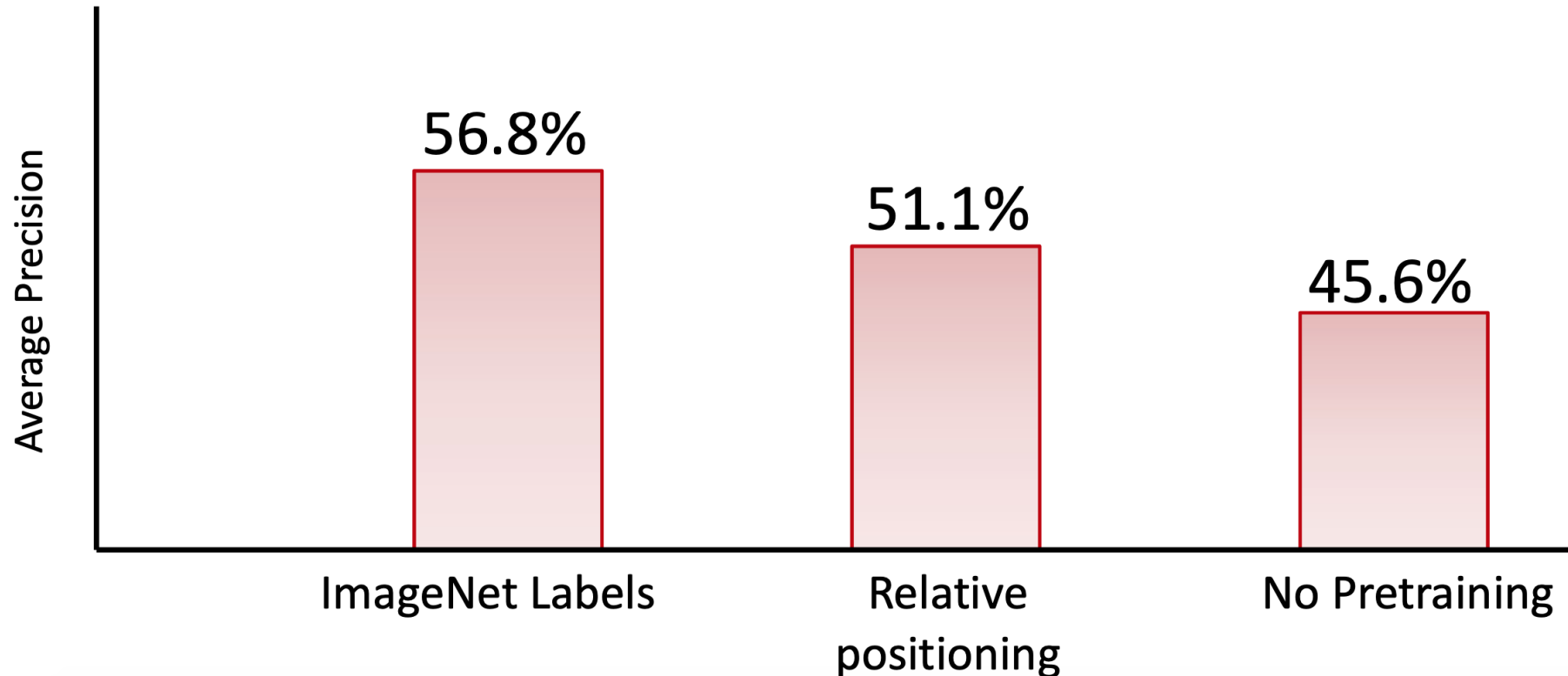
R-CNN



[Girshick et al. 2014]

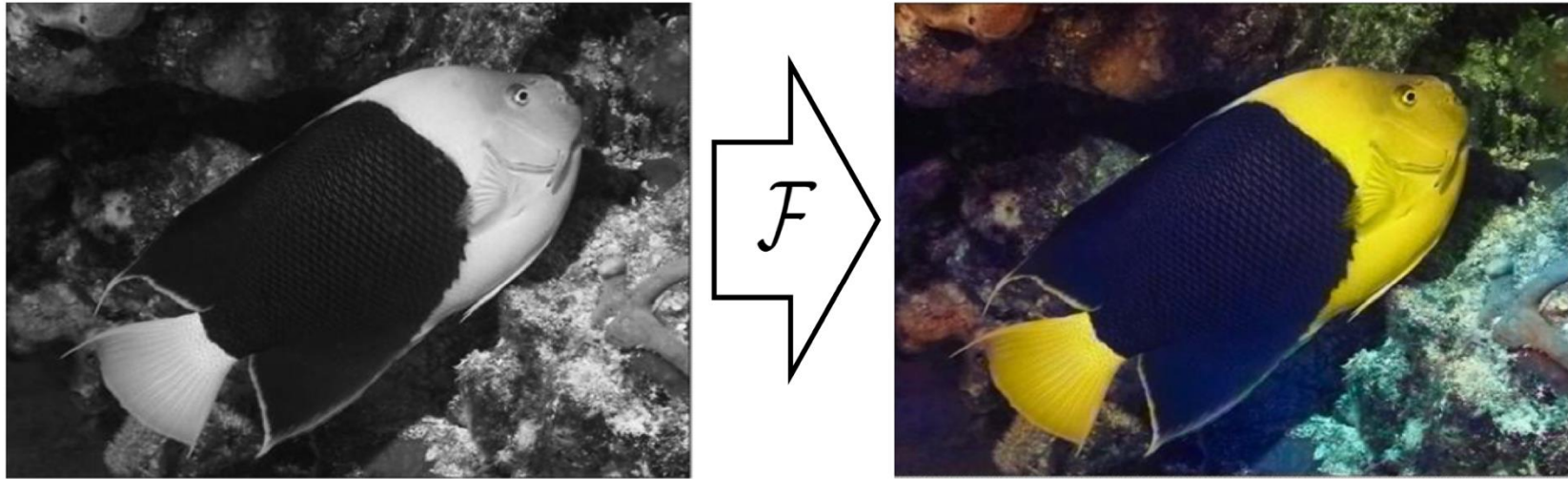
# SSL from Images, EX (I): relative positioning

## Evaluation: PASCAL VOC Detection



# SSL from Images, EX (II): colorization

Train network to predict pixel colour from a monochrome input

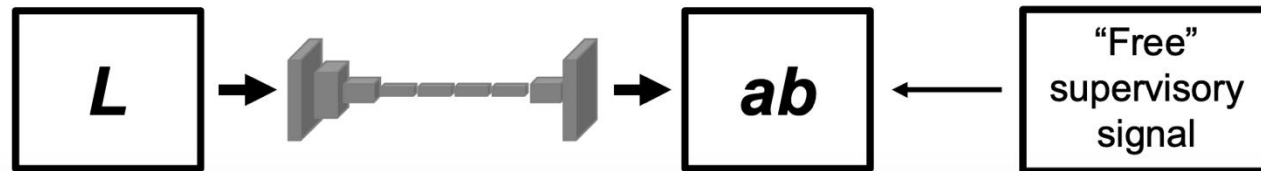


Grayscale image:  $L$  channel

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

Concatenate ( $L, ab$ )

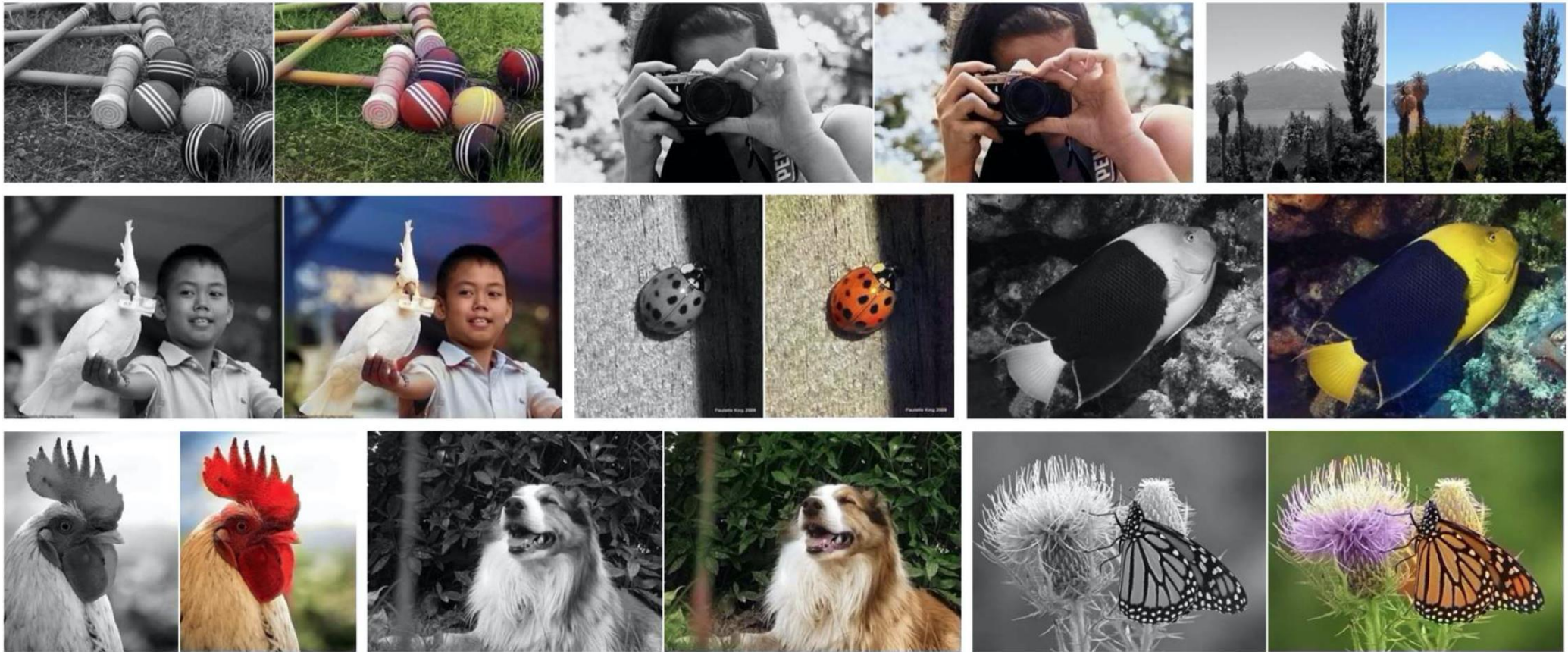
$$(\mathbf{X}, \hat{\mathbf{Y}})$$





# SSL from Images, EX (II): colorization

Train network to predict pixel colour from a monochrome input



# SSL from Images, EX (III): exemplar networks

- Exemplar Networks (Dosovitskiy et al., 2014)
- Perturb/distort image patches, e.g. by cropping and affine transformations
- Train to classify these exemplars as same class





# SSL from Images, EX (IV): masked autoencoder (MAE)

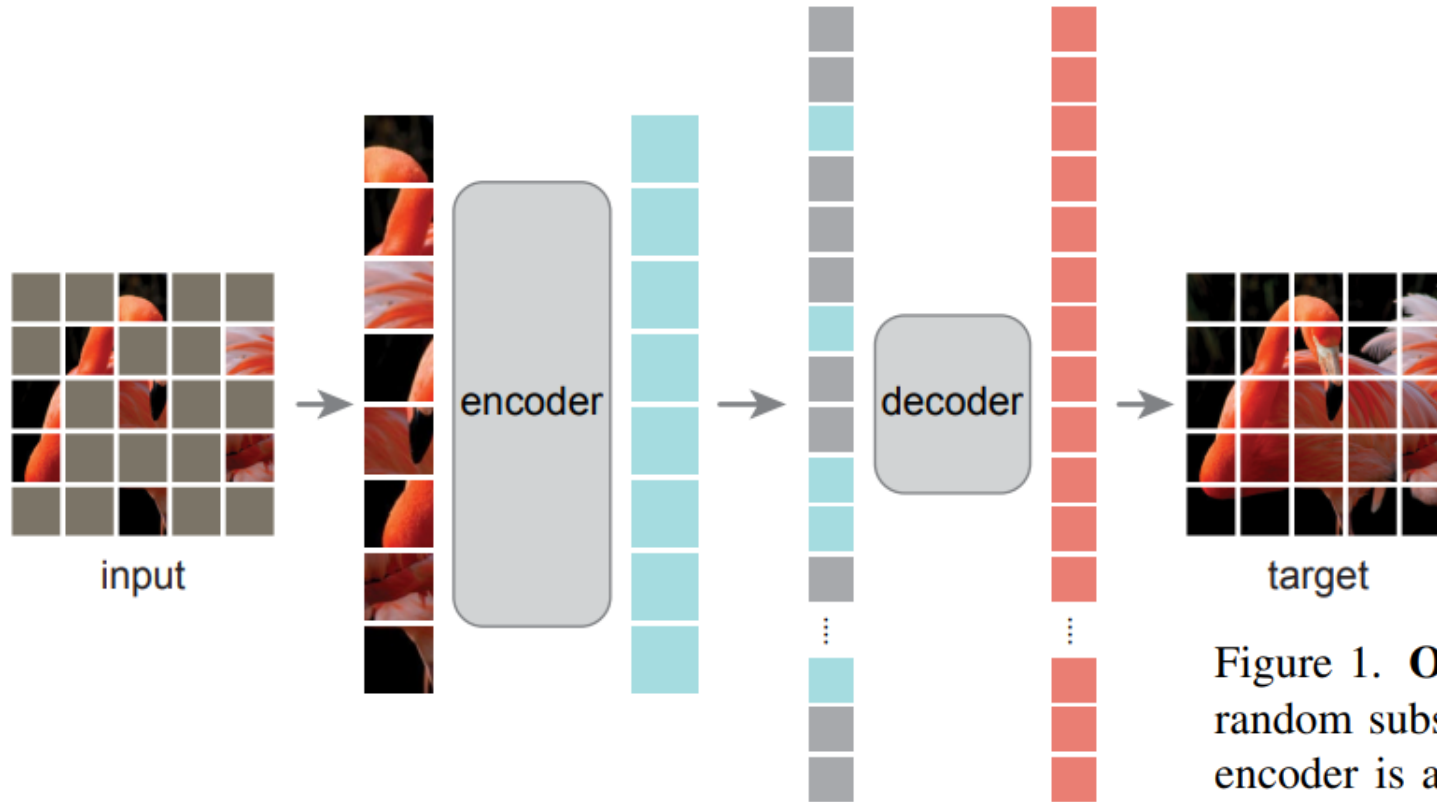
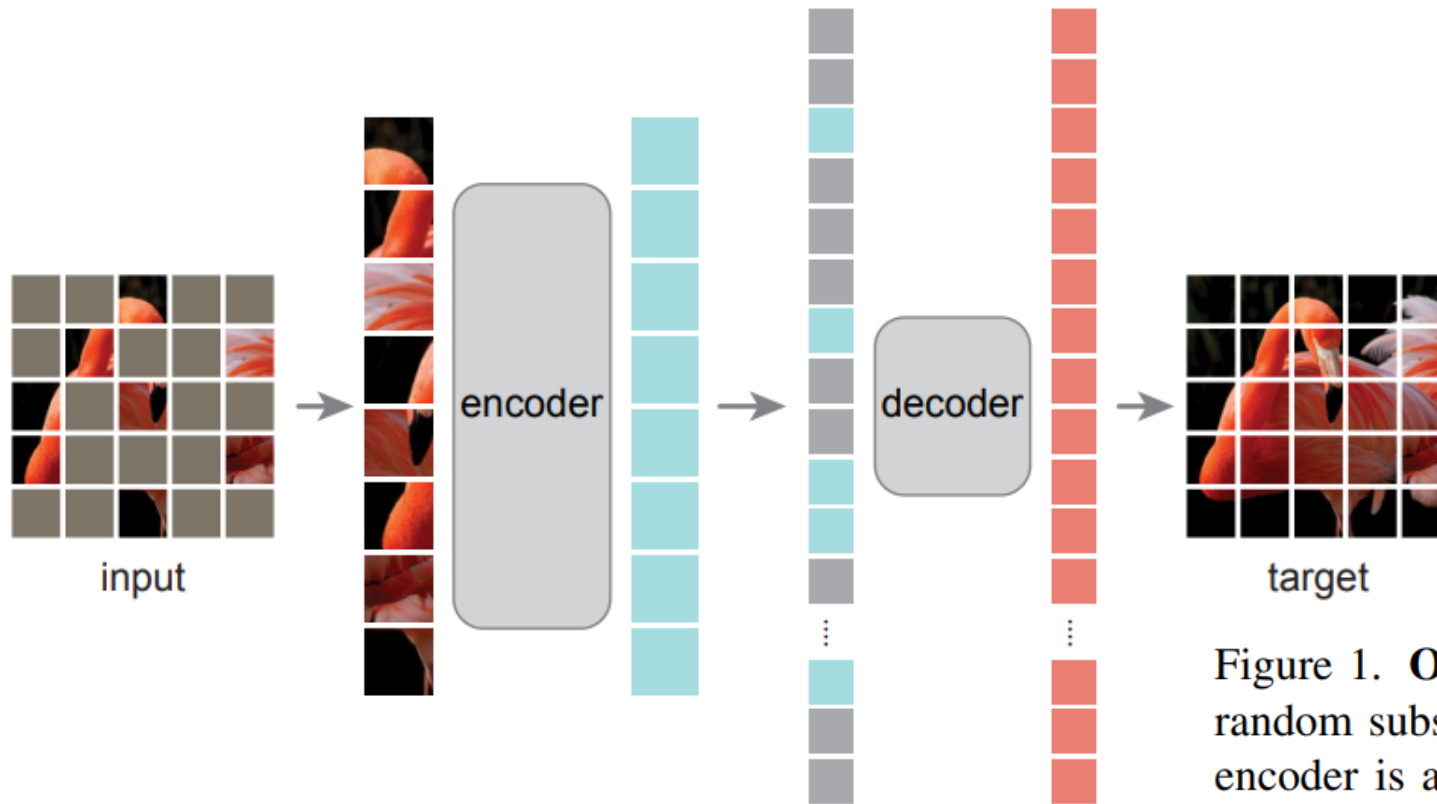


Figure 1. **Our MAE architecture.** During pre-training, a large random subset of image patches (*e.g.*, 75%) is masked out. The encoder is applied to the small subset of *visible patches*. Mask tokens are introduced *after* the encoder, and the full set of encoded patches and mask tokens is processed by a small decoder that reconstructs the original image in pixels. After pre-training, the decoder is discarded and the encoder is applied to uncorrupted images (full sets of patches) for recognition tasks.

# SSL from Images, EX (IV): masked autoencoder (MAE)

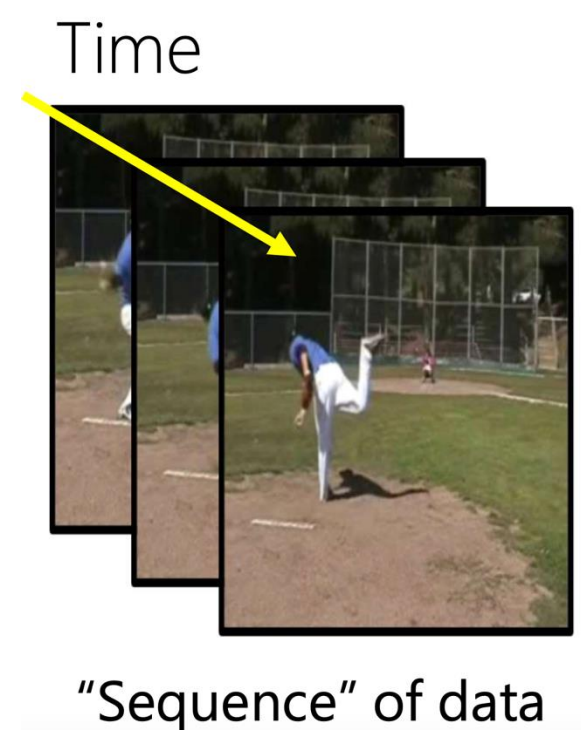


**Question:** Why is this (75%) much larger than the mask rate in BERT (15%)?

Figure 1. **Our MAE architecture.** During pre-training, a large random subset of image patches (*e.g.*, 75%) is masked out. The encoder is applied to the small subset of *visible patches*. Mask tokens are introduced *after* the encoder, and the full set of encoded patches and mask tokens is processed by a small decoder that reconstructs the original image in pixels. After pre-training, the decoder is discarded and the encoder is applied to uncorrupted images (full sets of patches) for recognition tasks.

# SSL from Videos

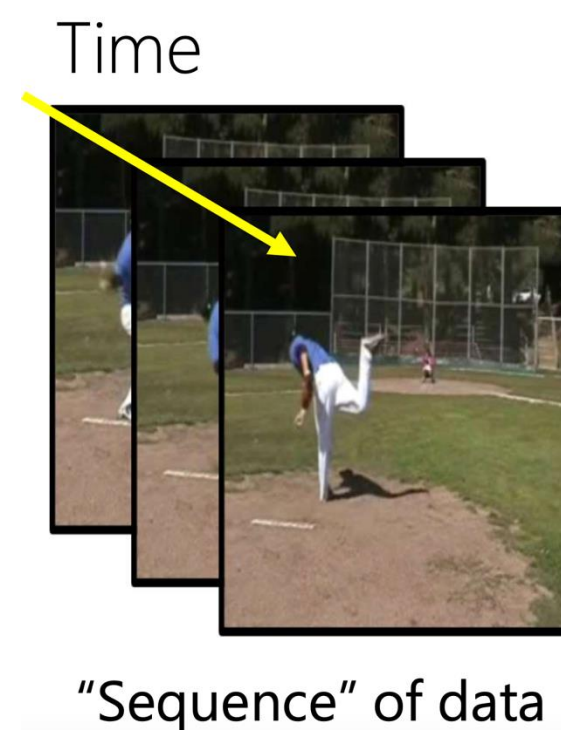
**Question:** What're your ideas of SSL from videos?



# SSL from Videos

Four example tasks:

- Video sequence order
  - Sequential Verification: Is this a valid sequence?



# SSL from Videos

Four example tasks:

- Video sequence order
  - Sequential Verification: Is this a valid sequence?
- Video direction
  - Predict if video playing forwards or backwards

# SSL from Videos

Four example tasks:

- Video sequence order
  - Sequential Verification: Is this a valid sequence?
- Video direction
  - Predict if video playing forwards or backwards
- Video tracking
  - Given a color video, colorize all frames of a gray scale version using a reference frame



# SSL from Videos

Four example tasks:

- Video sequence order
  - Sequential Verification: Is this a valid sequence?
- Video direction
  - Predict if video playing forwards or backwards
- Video tracking
  - Given a color video, colorize all frames of a gray scale version using a reference frame
- Video next frame prediction

# Key Takeaways

- Self supervision learning
  - Predicting any part of the observations given any available information
  - The prediction task forces models to learn semantic representations
  - Massive/unlimited data supervisions
- SSL for text:
  - Language models: next word prediction
  - BERT text representations: masked language model (MLM)
- SSL for images/videos:
  - Various ways of defining the prediction task



**Questions?**