

# CS395T: Structured Models for NLP

## Lecture 21: Deep Generative Models I



Greg Durrett



# Administrivia

---

- ▶ Final project proposals due today
- ▶ Project 3 grades back tonight



# Project 3 Results

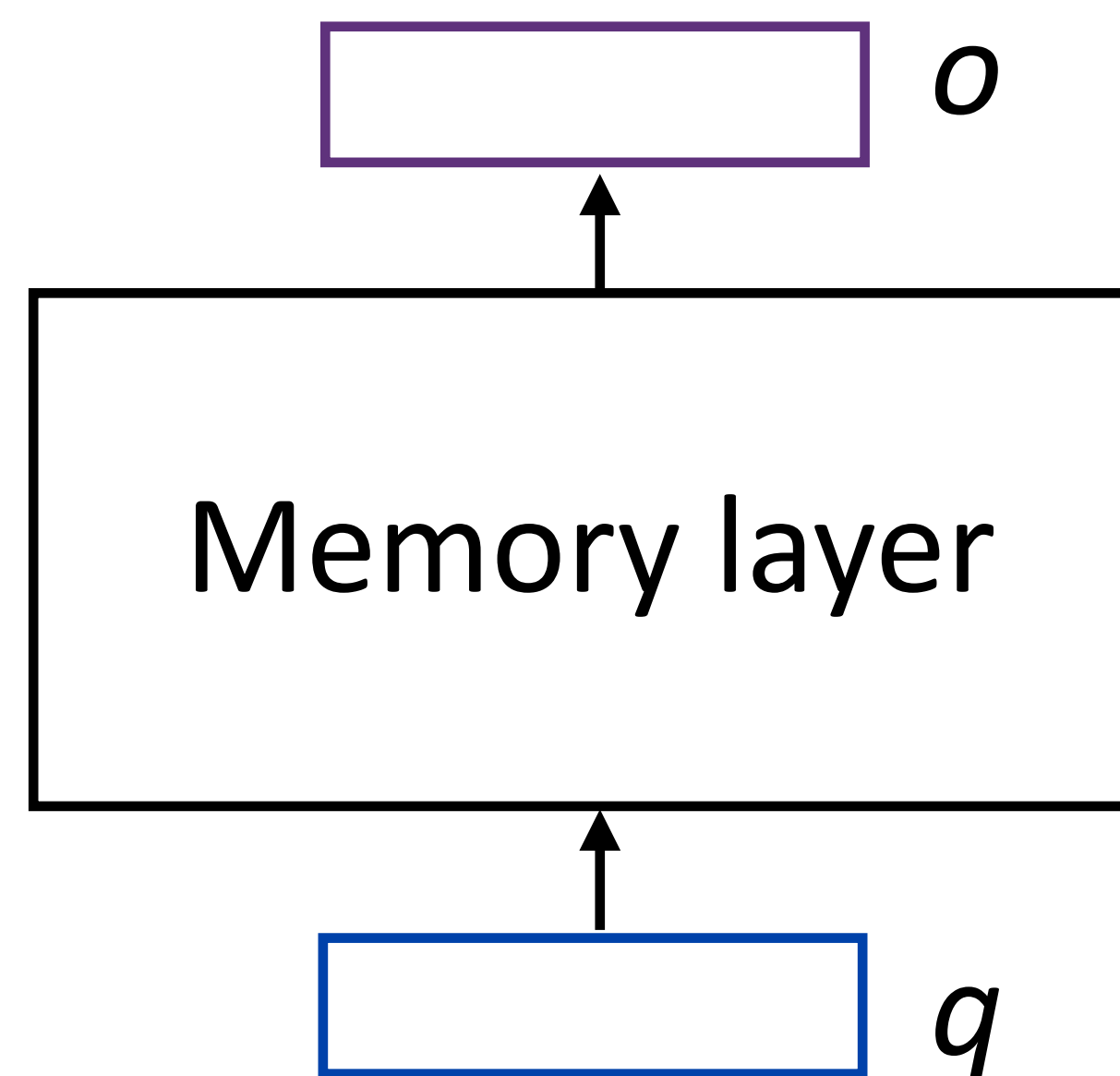
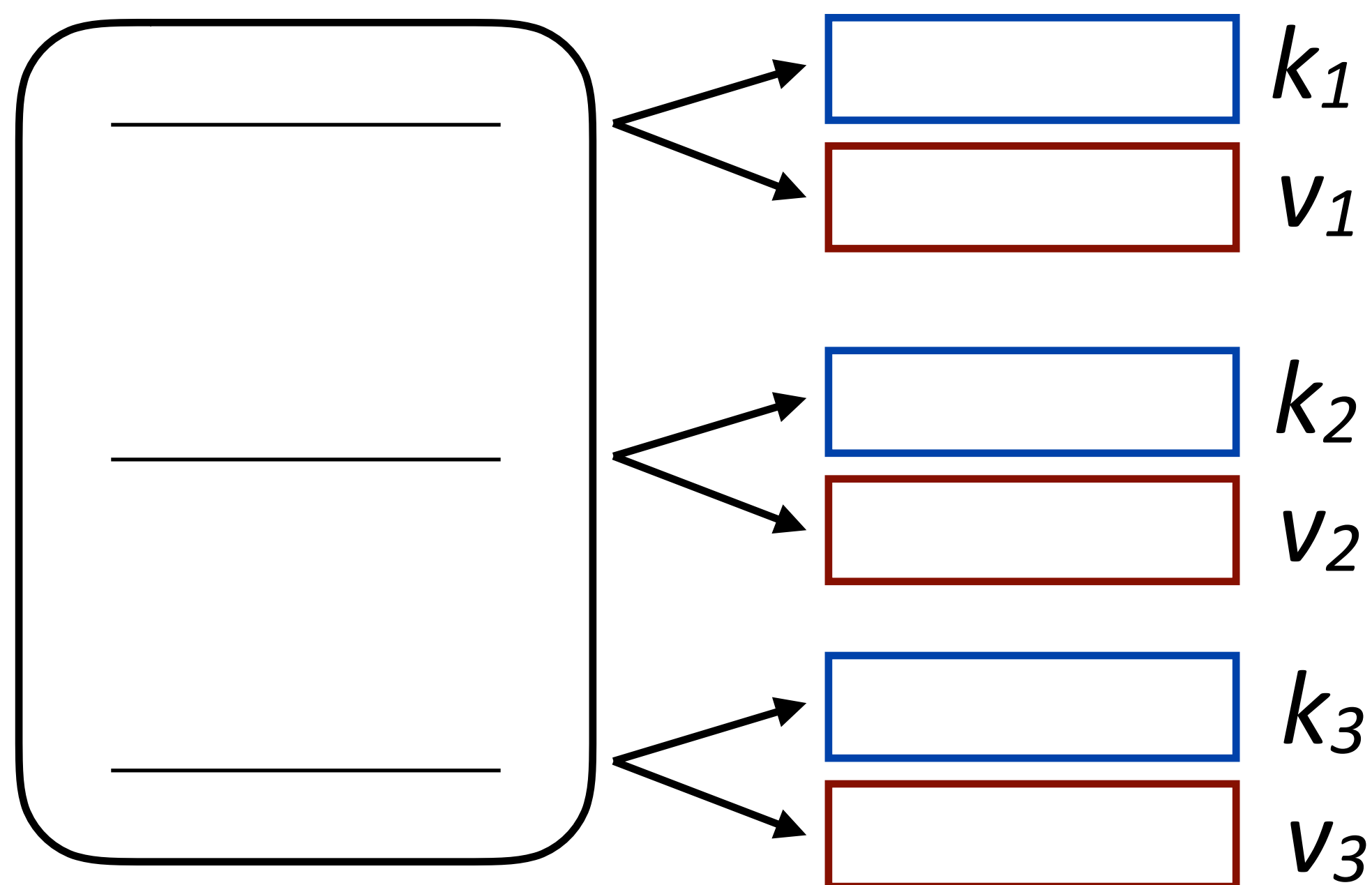
---

- ▶ Tanya Goyal: 82.46
  - ▶ Ensemble of 2-layer LSTMs with mean pooling, two-step training procedure: fine-tune vectors, then fix vectors and learn LSTM
- ▶ Su Wang: 81.89
  - ▶ 2-channel CNN with 100 feature maps, batch = 64, 25 epochs, L2 regularization and aggressive learning rate decay (0.95 per 100 epochs)
- ▶ Aditya Gupta: 81.80
  - ▶ Ensemble of (Bi?)LSTMs, hidden state dim = 200, different learning rates for different pieces of the model
- ▶ Elisa Ferracane: 81.23
  - ▶ BiLSTM with hidden state dim = 32, batch = 16, epochs = 10



# Recall: Memory Networks

- ▶ Memory networks let you reference input in an attention-like way
- ▶ Memorize input items into two vectors: a **key** and a **value**
- ▶ Keys compute attention weights given a query, weighted sum of values gives the output



$$o = \sum_i \alpha_i v_i$$

$$\alpha = \text{softmax}(e)$$

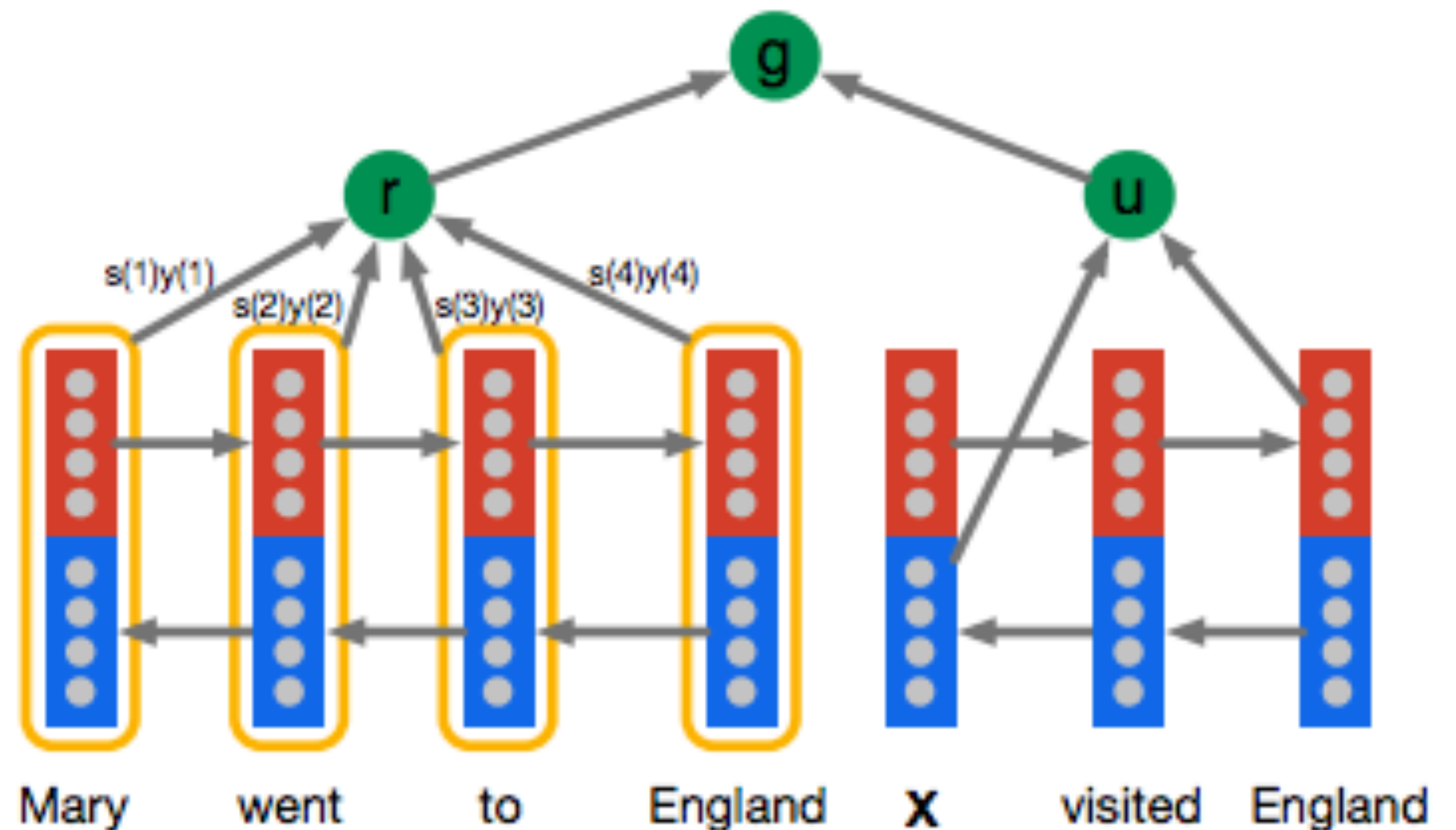
$$e_i = q \cdot k_i$$

Sukhbaatar et al. (2015)



# Recall: Attentive Reader

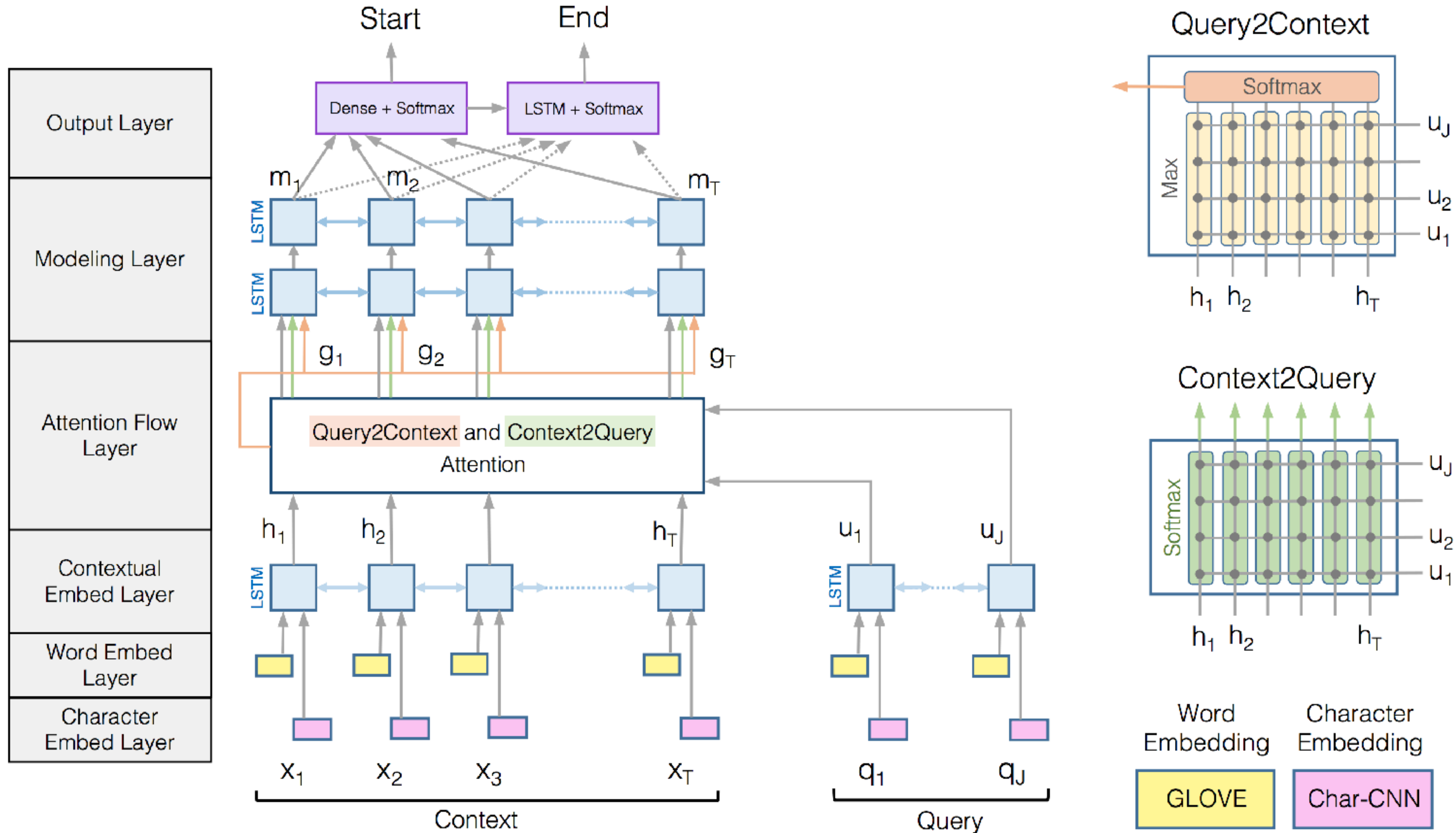
- ▶ Attentive reader: encode query, encode sentence, use attention to compute document representation, make prediction
- ▶ Uses fixed-size representations for the final prediction, multiclass classification







# Recall: Bidirectional Attention Flow





# This Lecture

---

- ▶ Variational autoencoders as deep generative models
  - ▶ Induce latent structure from the data by training in an unsupervised way
  - ▶ Can sample from them to produce examples
- ▶ Variational autoencoders as autoencoders

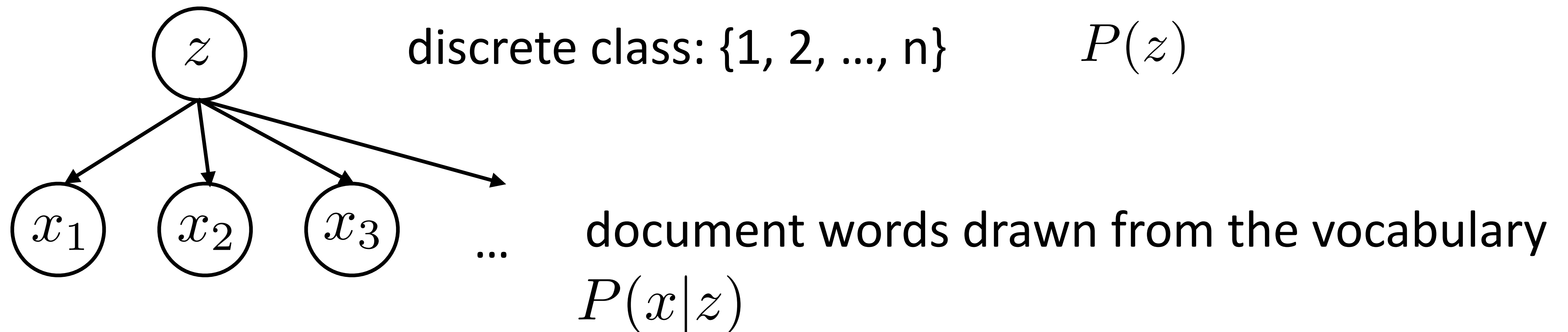
# Deep Generative Models





# Generative Models

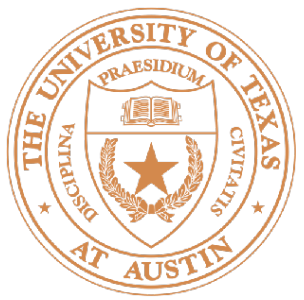
- Discrete class bag-of-words model:



Science

*too many drug trials too few patients*

$$P(z = \text{Science})P(\text{too}|z = \text{Science})P(\text{many}|z = \text{Science}) \dots$$



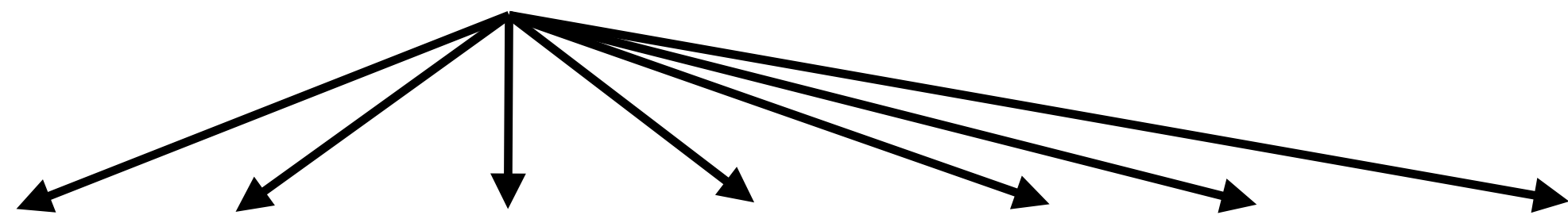
# Recall: EM for Generative Models

- ▶ Unsupervised learning: find parameters to maximize marginal likelihood

$$\log P(\mathbf{x}|\theta) = \log \sum_z P(\mathbf{x}, z|\theta)$$

- ▶ EM is a technique for doing this maximization

Science



*too many drug trials too few patients*

$$P(z = \text{Science})P(\text{too}|z = \text{Science})P(\text{many}|z = \text{Science}) \dots$$



# Recall: EM

$$\log \sum_z P(\mathbf{x}, z | \theta)$$

$$= \log \sum_z q(z) \frac{P(\mathbf{x}, z | \theta)}{q(z)}$$

► Variational approximation  $q$

$$\geq \sum_z q(z) \log \frac{P(\mathbf{x}, z | \theta)}{q(z)}$$

► Jensen's inequality (uses concavity of log)

$$= \mathbb{E}_{q(z)} \log P(\mathbf{x}, z | \theta) + \text{Entropy}[q(z)]$$

► Can optimize this lower-bound on log likelihood instead of log-likelihood



# Recall: EM

---

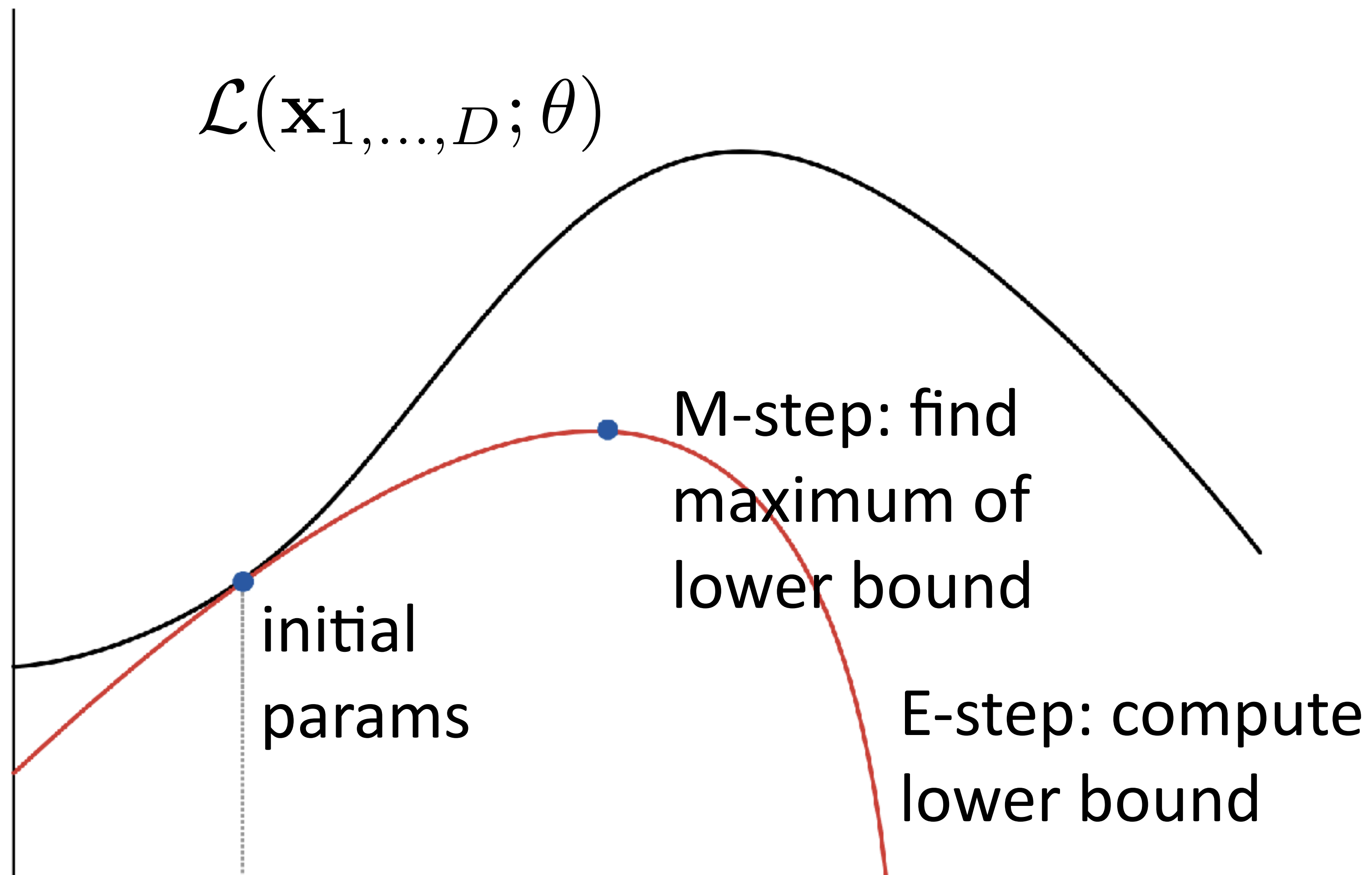
$$\log \sum_z P(\mathbf{x}, z|\theta) \geq \mathbb{E}_{q(z)} \log P(\mathbf{x}, z|\theta) + \text{Entropy}[q(z)]$$

- ▶ If  $q(z) = P(z|\mathbf{x}, \theta)$ , equality is achieved
- ▶ Expectation-maximization: alternating maximization of the lower bound over  $q$  and  $\theta$ 
  - ▶ Current timestep =  $t$ , have parameters  $\theta^{t-1}$
  - ▶ E-step: maximize w.r.t.  $q$ ; that is,  $q^t = P(z|\mathbf{x}, \theta^{t-1})$
  - ▶ M-step: maximize w.r.t.  $\theta$ ; that is,  $\theta^t = \operatorname{argmax}_{\theta} \mathbb{E}_{q^t} \log P(\mathbf{x}, z|\theta)$



# Recall: EM

$$\mathcal{L}(\mathbf{x}_1, \dots, \mathbf{x}_D) = \sum_{i=1}^D \log \sum_{\mathbf{y}} P(\mathbf{y}, \mathbf{x}_i)$$



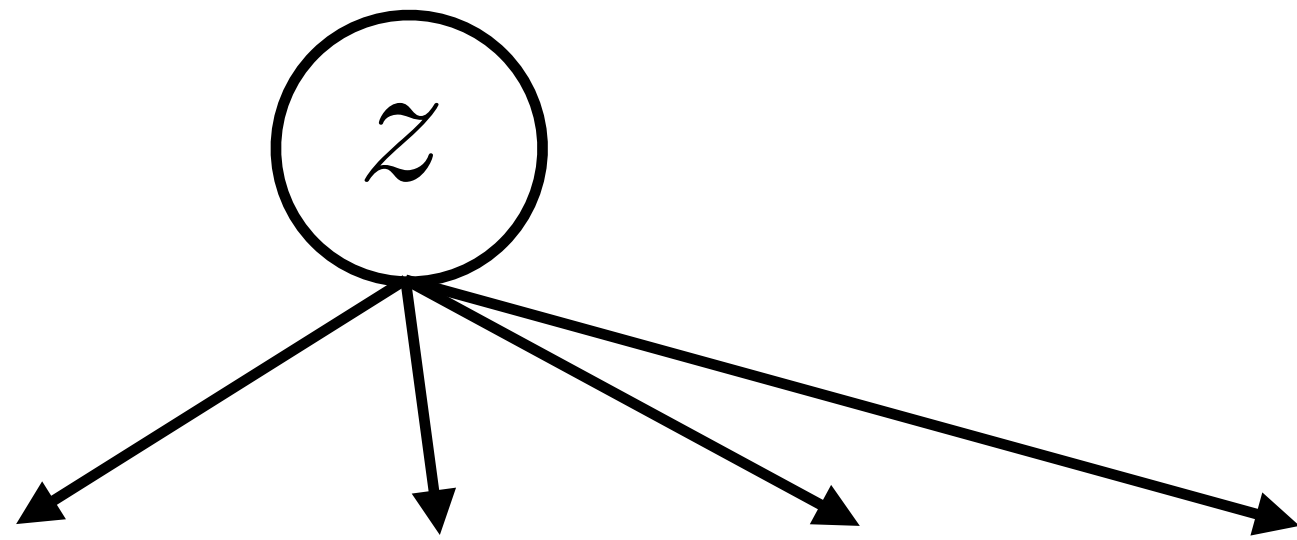


# EM for Generative Models

- ▶ What form does  $q$  take?  $q$  is a multinomial, so just a bunch of numbers

$$q(z) = P(z|\mathbf{x}, \theta) \propto P(z) \prod_i P(x_i|z)$$

- ▶ Easy to compute, easy to represent
- ▶ M-step: supervised learning problem with fractional annotation; possible because we can take the expectation:  $\theta^t = \operatorname{argmax}_{\theta} \mathbb{E}_{q^t} \log P(\mathbf{x}, z|\theta)$

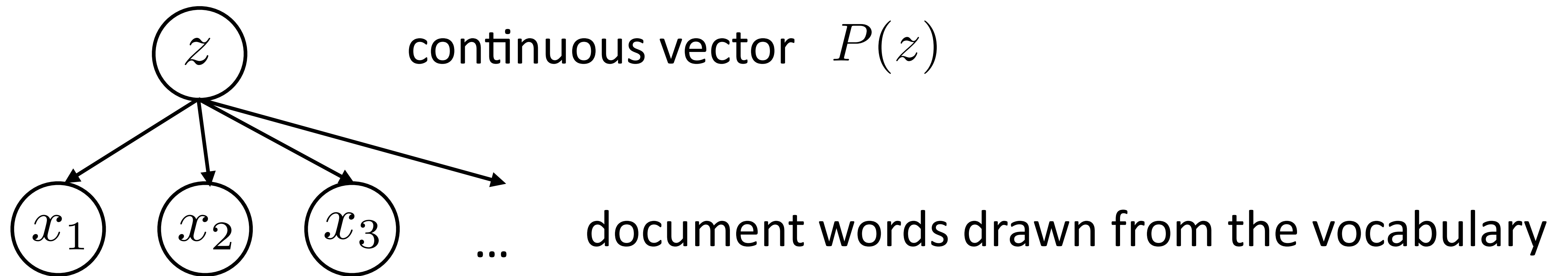


*too many drug trials too few patients*



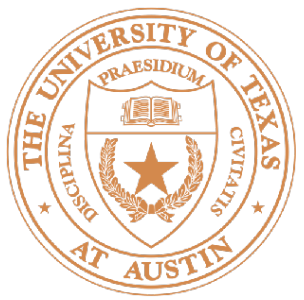


# Deep Generative Models



$$P(x|z) = \text{softmax}(\text{emb}(x)^\top z - b_x)$$

- ▶ Neural document model: probability of a word is depends on its dot product with the topic vector  $z$  (or use an even more complicated NN)
- ▶ What is  $P(z)$ ? Let's just say  $N(0, 1)$  for now...



# EM for Deep Generative Models

Expectation-maximization: alternating maximization over  $q$  and  $\theta$

E-step: maximize w.r.t.  $q$ ; that is,  ~~$q^t = P(z|\mathbf{x}, \theta^{t-1})$~~

M-step: maximize w.r.t.  $\theta$ ; that is,  $\theta^t = \operatorname{argmax}_{\theta} \mathbb{E}_{q^t} \log P(\mathbf{x}, z|\theta)$

$$P(z) = N(0, 1) \quad P(x|z) = \operatorname{softmax}(\operatorname{emb}(x)^\top z - b_x)$$

►  $P(z|\mathbf{x})$  is now a complicated distribution, can't simply use it for  $q$



# EM for Deep Generative Models

Expectation-maximization: alternating maximization over  $q$  and  $\theta$

E-step: maximize w.r.t.  $q$ ; that is,  $q^t = \operatorname{argmin}_q \operatorname{KL}(q(z) \| P(z|\mathbf{x}))$

M-step: maximize w.r.t.  $\theta$ ; that is,  $\theta^t = \operatorname{argmax}_\theta \mathbb{E}_{q^t} \log P(\mathbf{x}, z|\theta)$

$$P(z) = N(0, 1) \quad P(x|z) = \operatorname{softmax}(\operatorname{emb}(x)^\top z - b_x)$$

- ▶  $P(z|\mathbf{x})$  is now a complicated distribution, can't simply use it for  $q$
- ▶ E-step: choose a family of distributions  $q$ , find the best  $q$  in that family
$$q = N(\mu, \operatorname{diag}(\sigma^2))$$
  - ▶ Even computing the best  $\mu$  and  $\sigma$  for an example is hard!
- ▶ M-step: now we need to take an expectation over a continuous distribution



# Deep Generative Models

- ▶ EM doesn't seem to be helping...let's start over with the objective

$$\log \sum_z P(\mathbf{x}, z|\theta) = \log \sum_z q(z) \frac{P(\mathbf{x}, z|\theta)}{q(z)} \geq \sum_z q(z) \log \frac{P(\mathbf{x}, z|\theta)}{q(z)}$$

Jensen

$$\begin{aligned} &= \mathbb{E}_{q(z|\mathbf{x})} [-\log q(z|\mathbf{x}) + \log P(\mathbf{x}, z|\theta)] \\ &= \mathbb{E}_{q(z|\mathbf{x})} [\log P(\mathbf{x}|z, \theta)] - \text{KL}(q(z|\mathbf{x}) || P(z)) \end{aligned}$$

- ▶ Different arrangement of terms: KL between  $q$  and *prior* + conditional likelihood term



# Comparison of Objectives

---

- ▶ EM:  $\mathbb{E}_{q(z|\mathbf{x})} [\log P(\mathbf{x}, z|\theta)] + \text{Entropy}[q(z|\mathbf{x})] + KL(q(z|\mathbf{x})||P(z|\mathbf{x}, \theta))$

“make the data likely under  $q$ ”  
(generative)

- ▶ VAE:  $\mathbb{E}_{q(z|\mathbf{x})} [\log P(\mathbf{x}|z, \theta)] - KL(q(z|\mathbf{x})||P(z))$

“make the data likely under  $q$ ” “make  $q$  close to the prior”  
(discriminative)

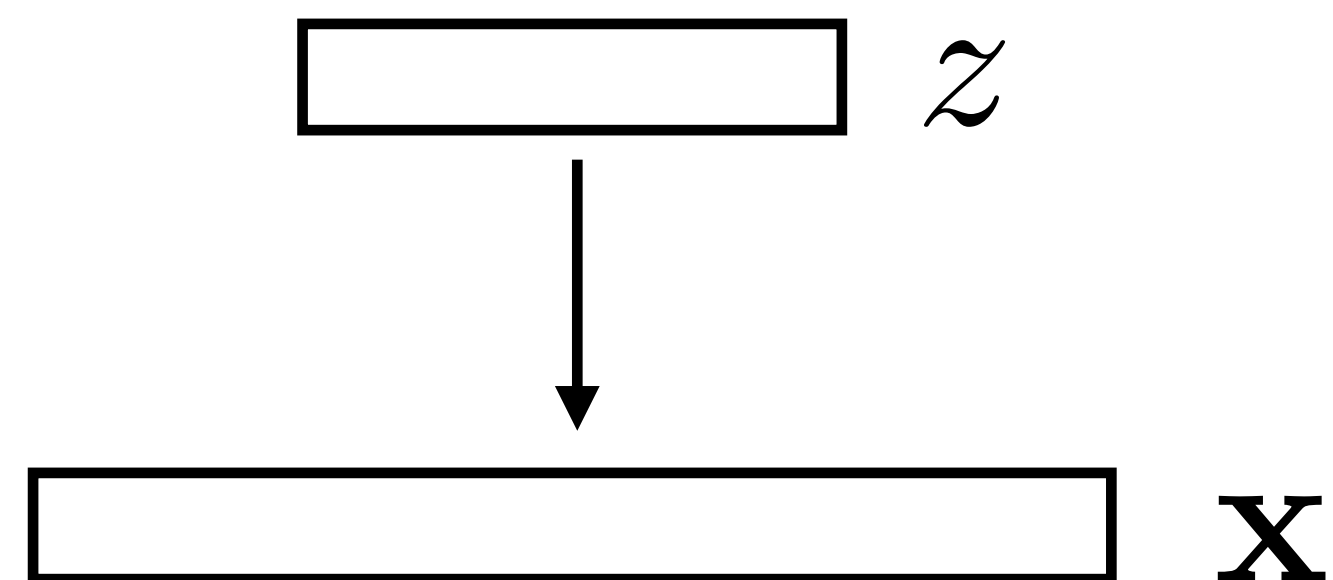
- ▶ Approximate  $q$  with a separate set of parameters, optimize  $q$  and  $\theta$  jointly with gradient descent
- ▶ Still need to reckon with that expectation over a continuous  $q(z)$ ...



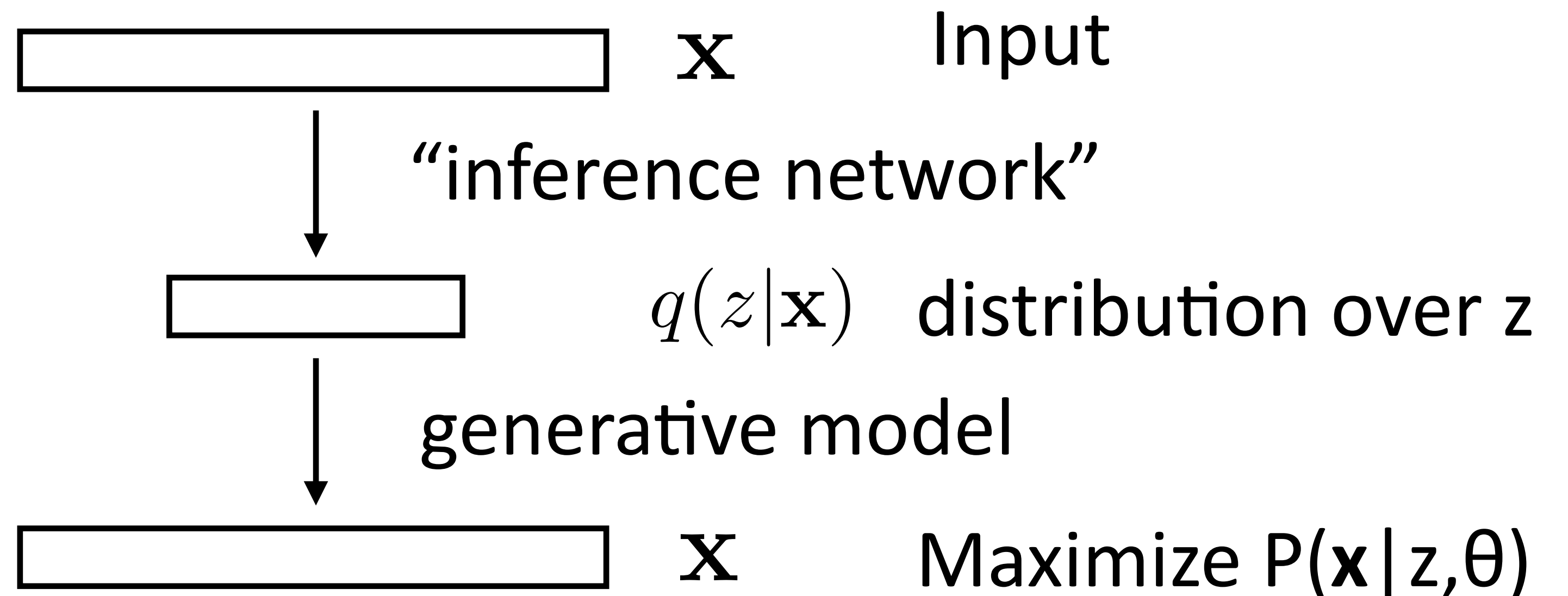
# Variational Autoencoders

$$\mathbb{E}_{q(z|\mathbf{x})} [\log P(\mathbf{x}|z, \theta)] + \text{KL}(q(z|\mathbf{x}) \| P(z))$$

Generative model (test):



Autoencoder (training):







# Training VAEs

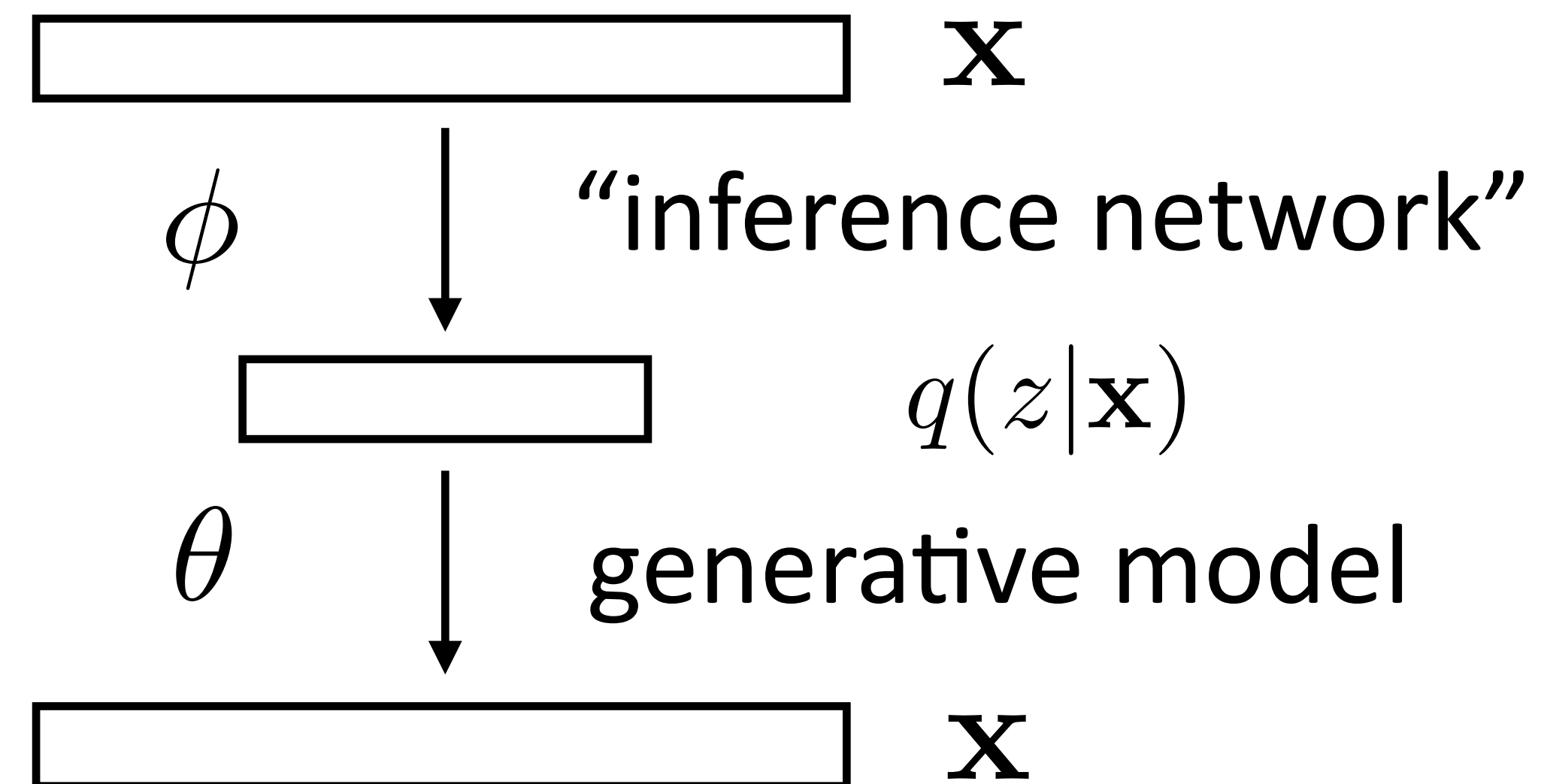
$$\mathbb{E}_{q(z|\mathbf{x})} [\log P(\mathbf{x}|z, \theta)] + \text{KL}(q(z|\mathbf{x}) \| P(z))$$

- Choose  $q$  to be Gaussian with parameters that are computed from  $\mathbf{x}$

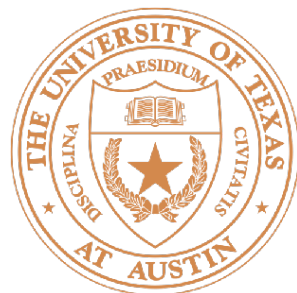
$$q = N(\mu(\mathbf{x}), \text{diag}(\sigma^2(\mathbf{x})))$$

- mu and sigma are computed from one-layer feedforward networks over  $\mathbf{x}$ , call their parameters  $\phi$

Autoencoder (training):



- How to handle the expectation? Just sample!



# Training VAEs

For each example  $\mathbf{x}$

Compute  $q$  (run forward pass to compute  $\mu$  and  $\sigma$ )

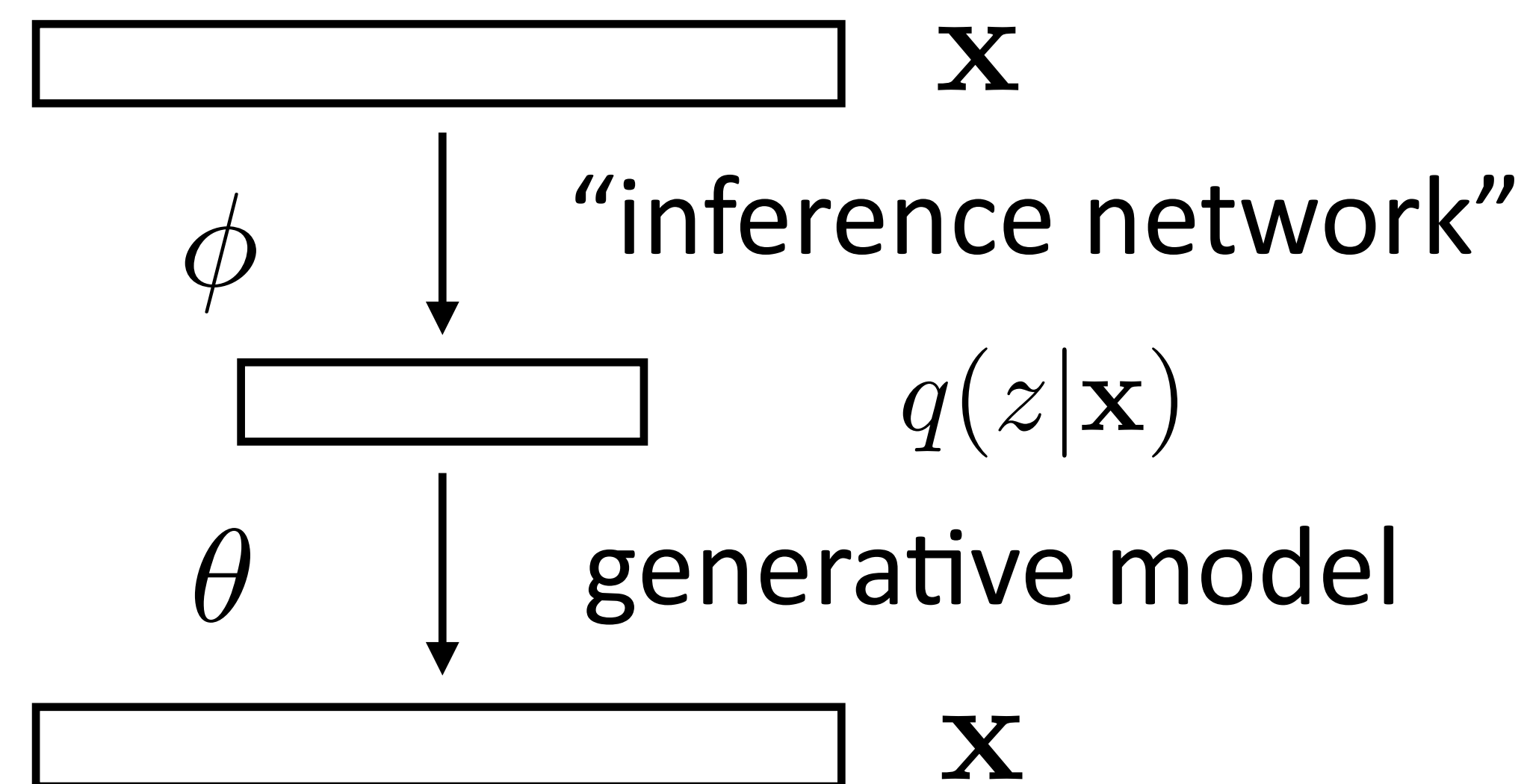
For some number of samples

Sample  $z \sim q$

Compute  $P(\mathbf{x}|z)$  and compute loss

Backpropagate to update  $\phi$ ,  $\theta$

Autoencoder (training):





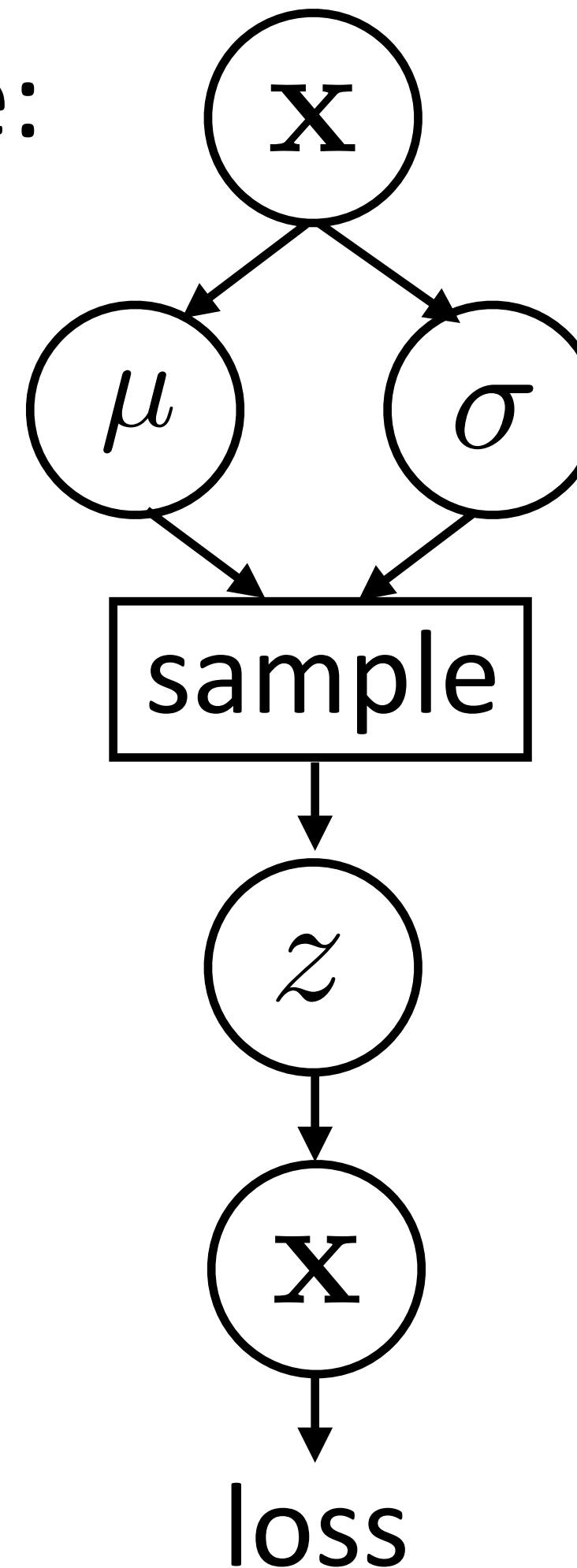
# Reparameterization Trick

- ▶ Can't backpropagate through a sampling operation

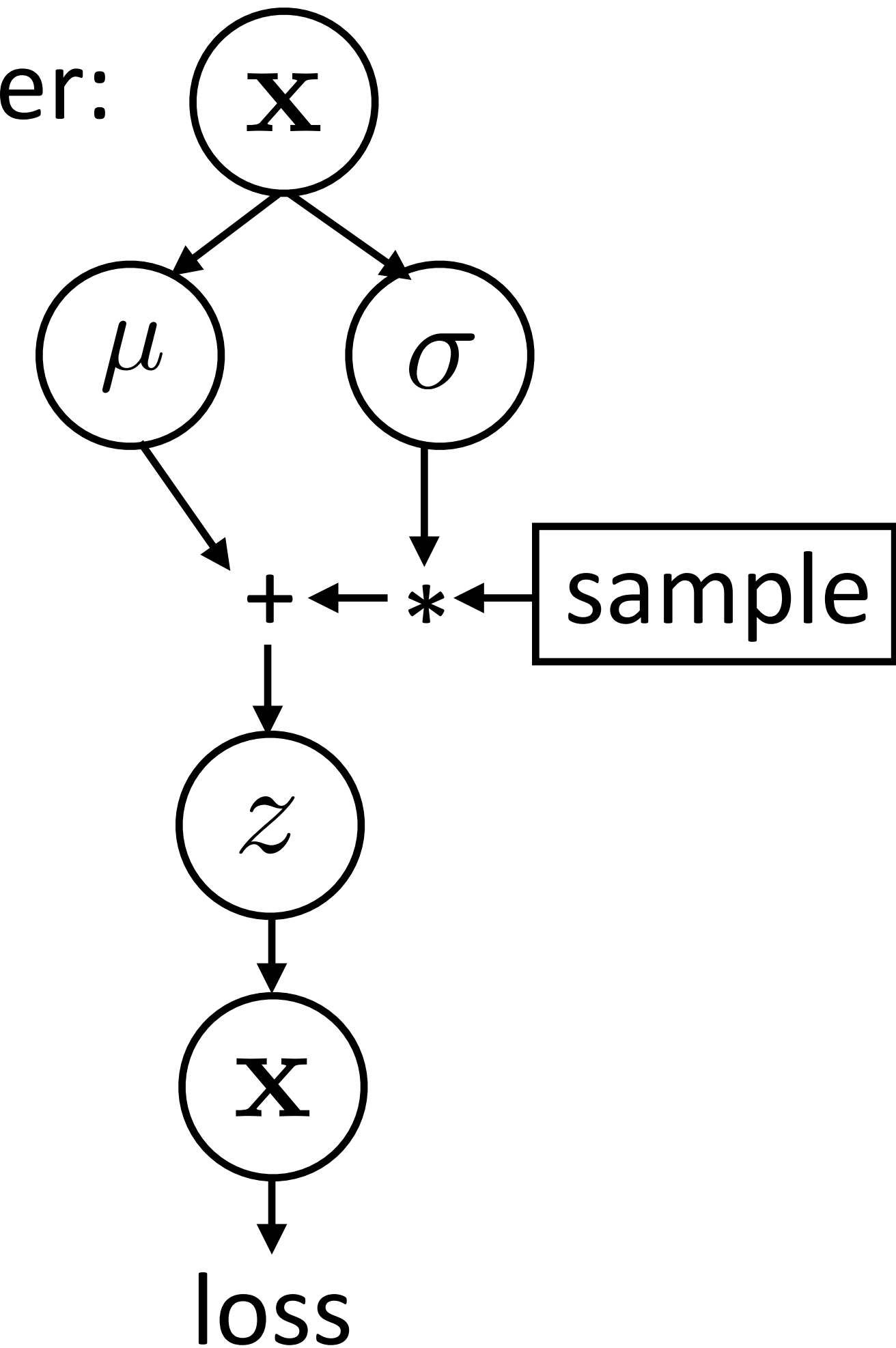
- ▶ Recall that

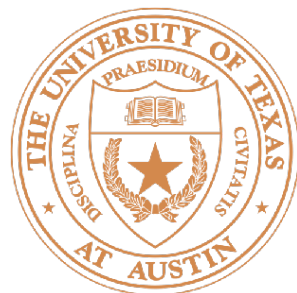
$$N(\mu, \text{diag}(\sigma^2)) = \mu + \sigma N(0, I)$$

Before:



After:





# Training VAEs

For each example  $\mathbf{x}$

Compute  $q$  (run forward pass to compute  $\mu$  and  $\sigma$ )

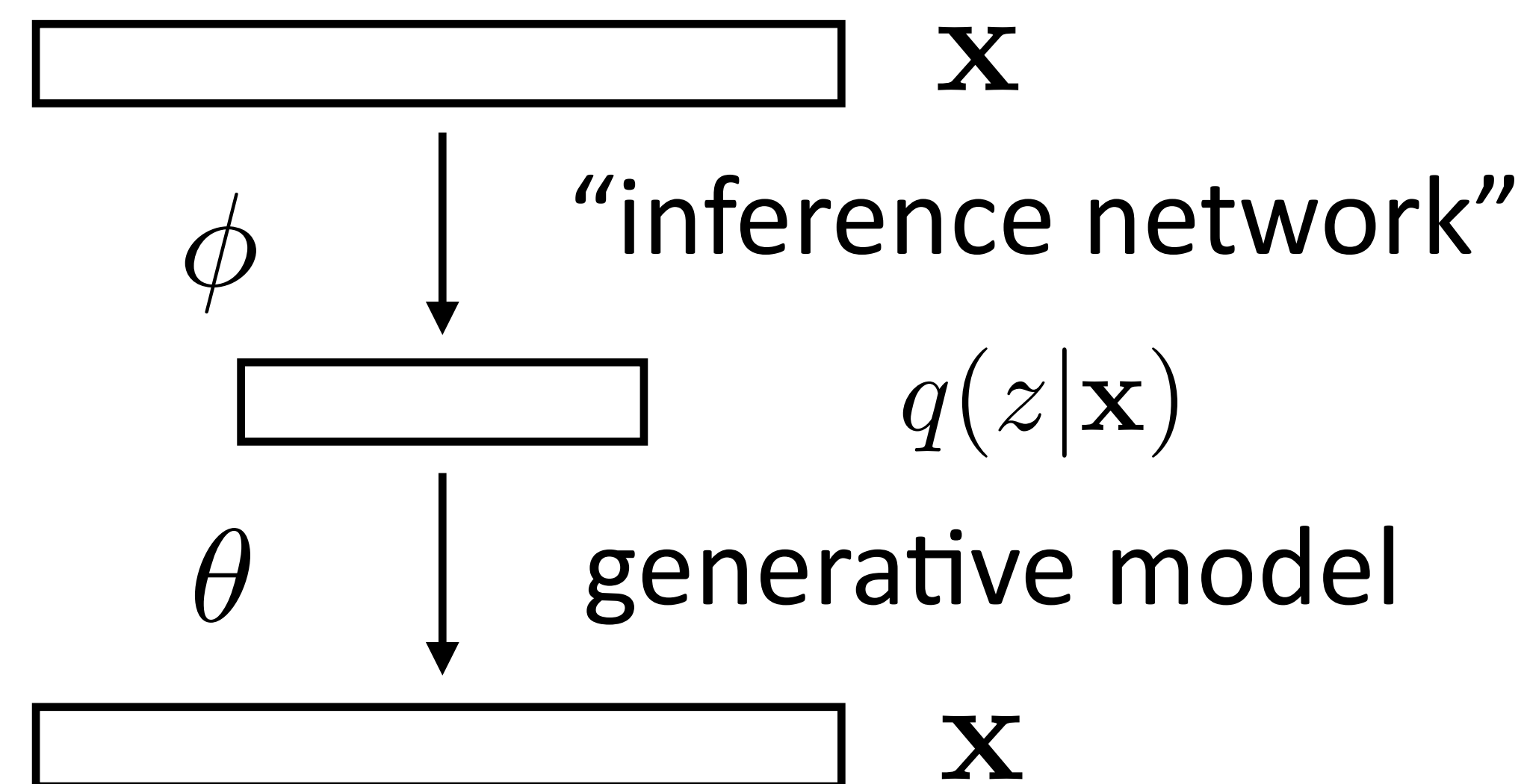
For some number of samples

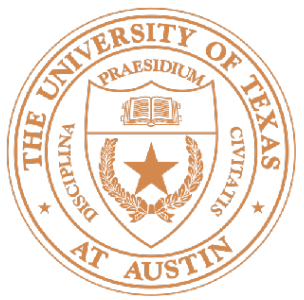
Sample  $z \sim q$

Compute  $P(\mathbf{x}|z)$  and compute loss

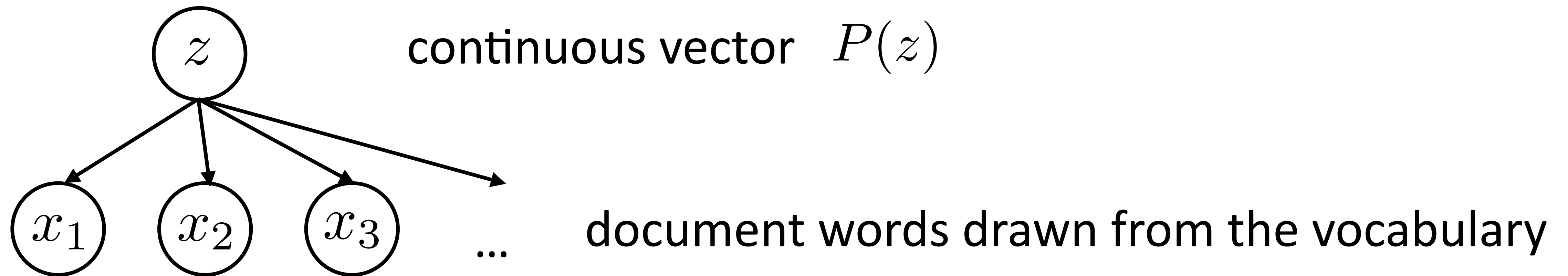
Backpropagate to update  $\phi$ ,  $\theta$

Autoencoder (training):





# VAEs as Deep Generative Models



$$P(x|z) = \text{softmax}(\text{emb}(x)^\top z - b_x)$$

- ▶ We've seen a way to train this real-valued bag-of-words model in a fully unsupervised way
- ▶ “Encoder network” looks like the E-step of EM (but has distinct parameters), backpropagate end-to-end through encoder and decoder



# Neural Variational Document Model

---

- ▶ Train this generative model on 20NewsGroups (online newsgroups) and RCV1 (newswire)
- ▶ Unsupervised learning: how to evaluate?
  - ▶ Data likelihood (perplexity)
  - ▶ See if interesting latent structure comes out





# Neural Variational Document Model

Model	Dim	20News	RCV1
LDA	50	1091	1437
LDA	200	1058	1142
RSM	50	953	988
docNADE	50	896	742
SBN	50	909	784
fDARN	50	917	724
fDARN	200	—	598
NVDM	50	<b>836</b>	563
NVDM	200	852	<b>550</b>

(a) Perplexity on test dataset.

- Randomly sample a dimension of  $z$ , see what words score highest along that axis, manually label that dimension

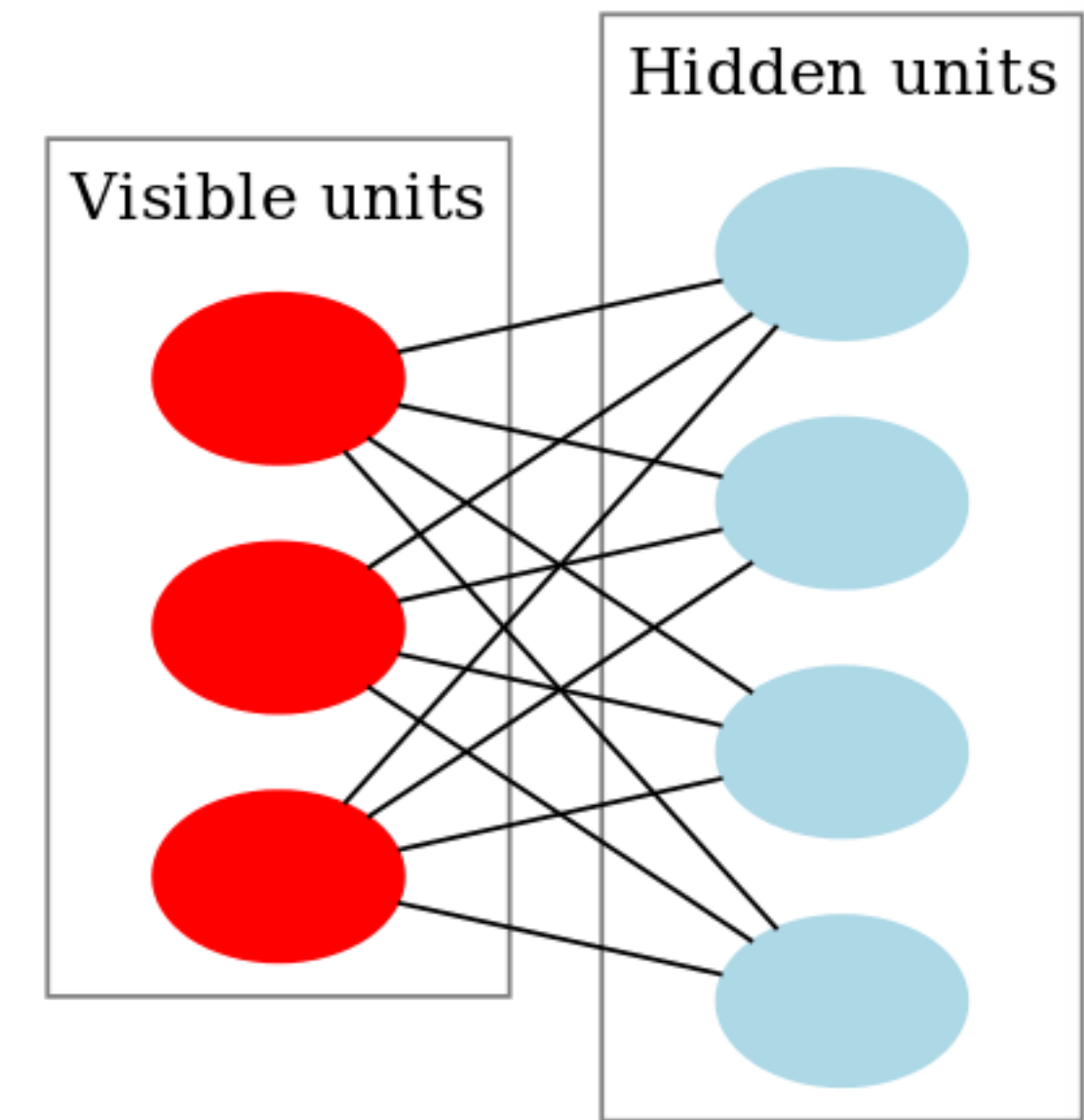
<i>Space</i>	<i>Religion</i>	<i>Encryption</i>	<i>Sport</i>	<i>Policy</i>
orbit	muslims	rsa	goals	bush
lunar	worship	cryptography	pts	resources
solar	belief	crypto	teams	charles
shuttle	genocide	keys	league	austin
moon	jews	pgp	team	bill
launch	islam	license	players	resolution
fuel	christianity	secure	nhl	mr
nasa	atheists	key	stats	misc
satellite	muslim	escrow	min	piece
japanese	religious	trust	buf	marc



# History: Restricted Boltzmann Machines

- ▶ Neural generative model with hidden (boolean) variables  $z$  and observed variables  $x$

$$P(x, z) = \frac{1}{Z} \exp(x^\top W z)$$



- ▶ Contrastive divergence:  
given  $x$ , compute  $P(z|x)$ , sample  $z$   
sample  $x' \sim P(x|z)$ , sample  $z' \sim P(z'|x)$   
update towards  $(z, x)$  away from  $(z', x')$

“inference network”

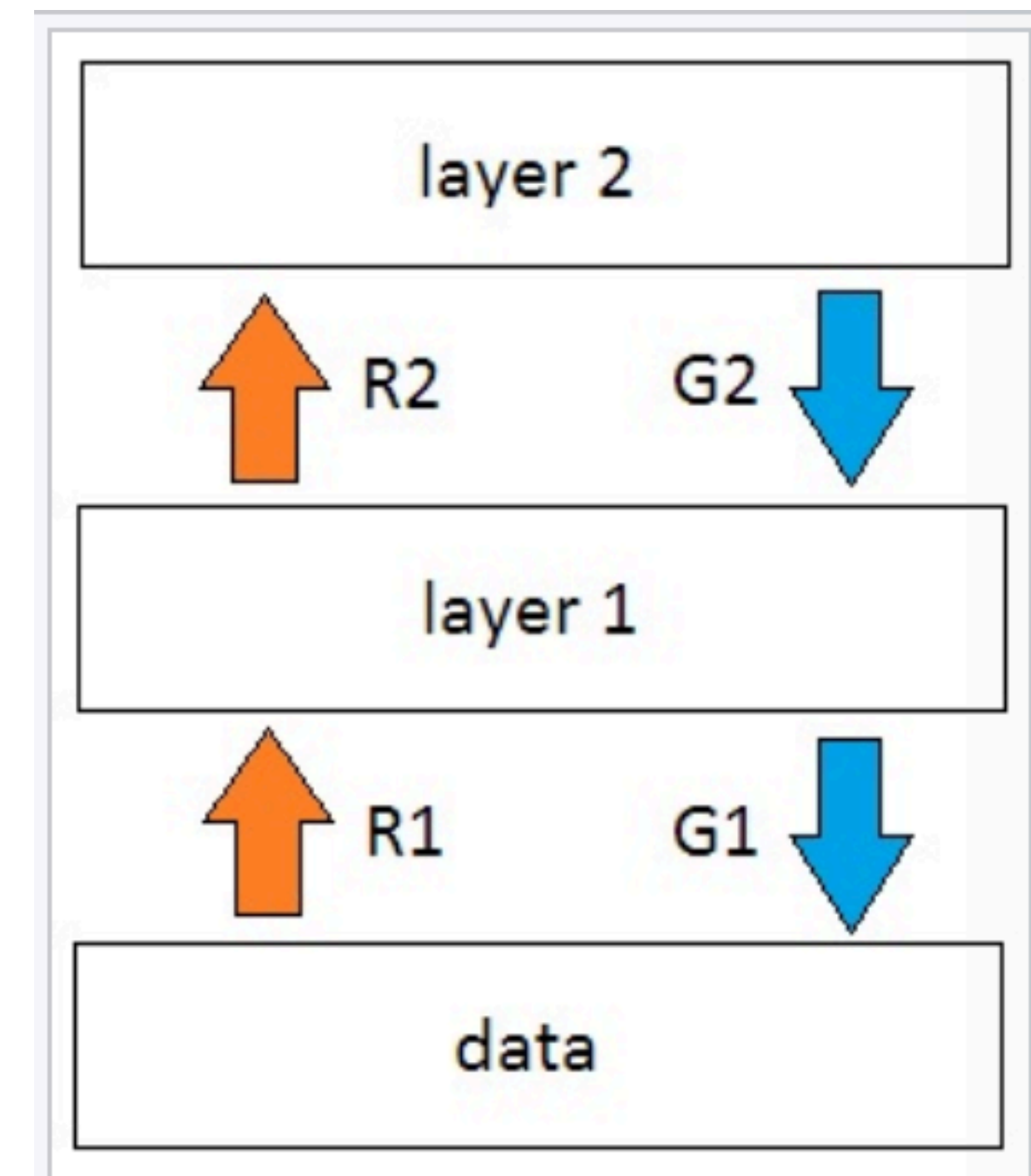
“generative network”

Smolensky (1986), Carreira-Perpiñán and Hinton (2005)



# History: Wake-Sleep Algorithm

- ▶ Deep generative model with generation parameters  $G$  and “recognition” parameters  $R$
- ▶ “Wake” phase: take data, encode it “upwards” using  $R$ , train  $G$  in a supervised way
- ▶ “Sleep” phase: generate top-down, train  $R$  in a supervised way
- ▶ One layer of this trained end-to-end looks like VAEs  
data  $\xrightarrow{R1}$  layer1  $\xrightarrow{G1}$  data



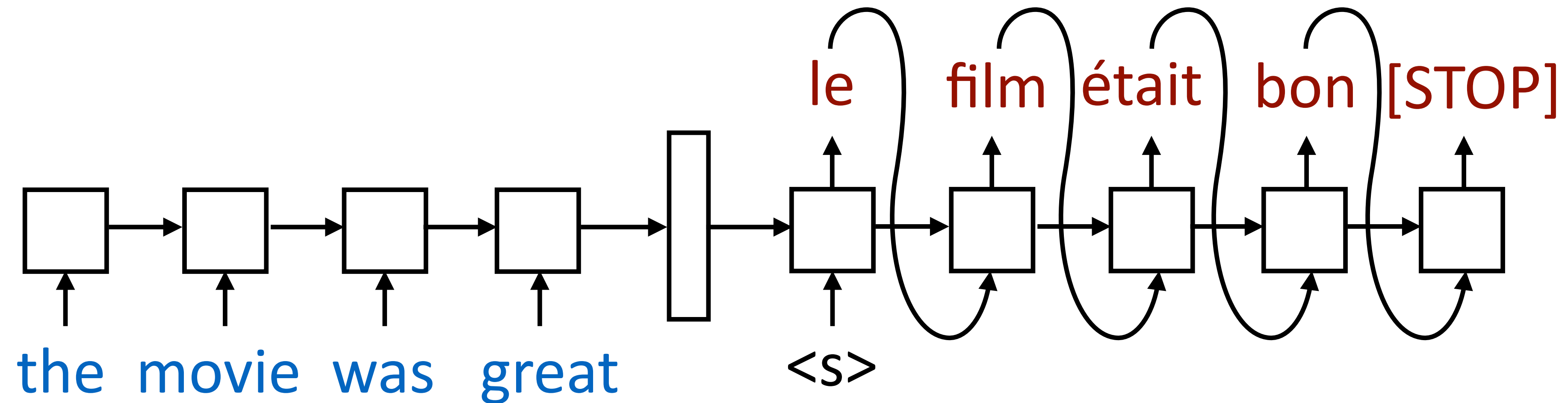
Hinton et al. (1995)

# VAEs as Autoencoders





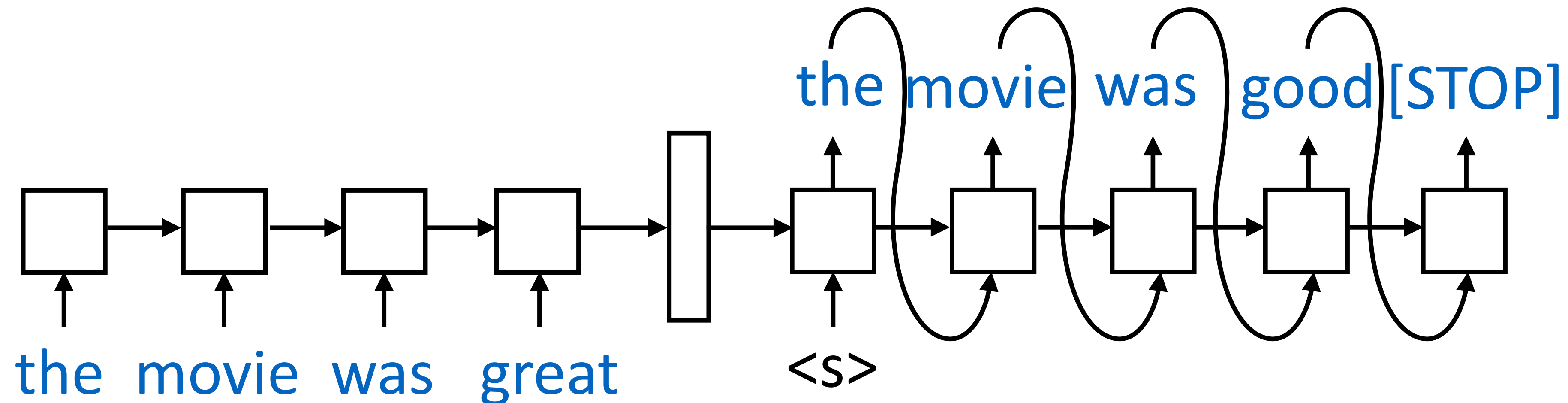
# Encoder-Decoder Models



- Encoder-decoder models without attention: compress the input into a single &\*&#! vector, unfold it to produce output



# Autoencoders



- ▶ Encoder-decoder models without attention: compress the input into a single &\*&#215;! vector, unfold it to produce output
- ▶ Autoencoder: encode input  $\mathbf{x}$  into a vector  $z$ , produce  $\mathbf{x}$  given  $z$

$$P(\mathbf{x}'|\mathbf{x}) = P(z|\mathbf{x}) \prod_i P(x'_i|z, \mathbf{x}'_{<i})$$

encoder                      decoder

- ▶ What semantics do we want the latent space to have?





- 
- A horizontal timeline with an arrow pointing to the right. Five vertical tick marks are placed along the line. Above the first tick mark is the word 'a'. Above the second tick mark is the word 'of'. Above the third tick mark is the word 'the'. Above the fourth tick mark is the word 'in'. Above the fifth tick mark are three dots '...'. The words are spaced out evenly along the timeline.

- the movie was **good** </s> </s>

$z_1$	$z_2$	$z_3$	$z_4$	$z_5$	$z_6$
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

the movie was great </s> </s> ☐ ☐ ☐ ☒ ☐ ☐

I thought the film was **good** ☐ ☐ ☐ ☐ ☐ ☒

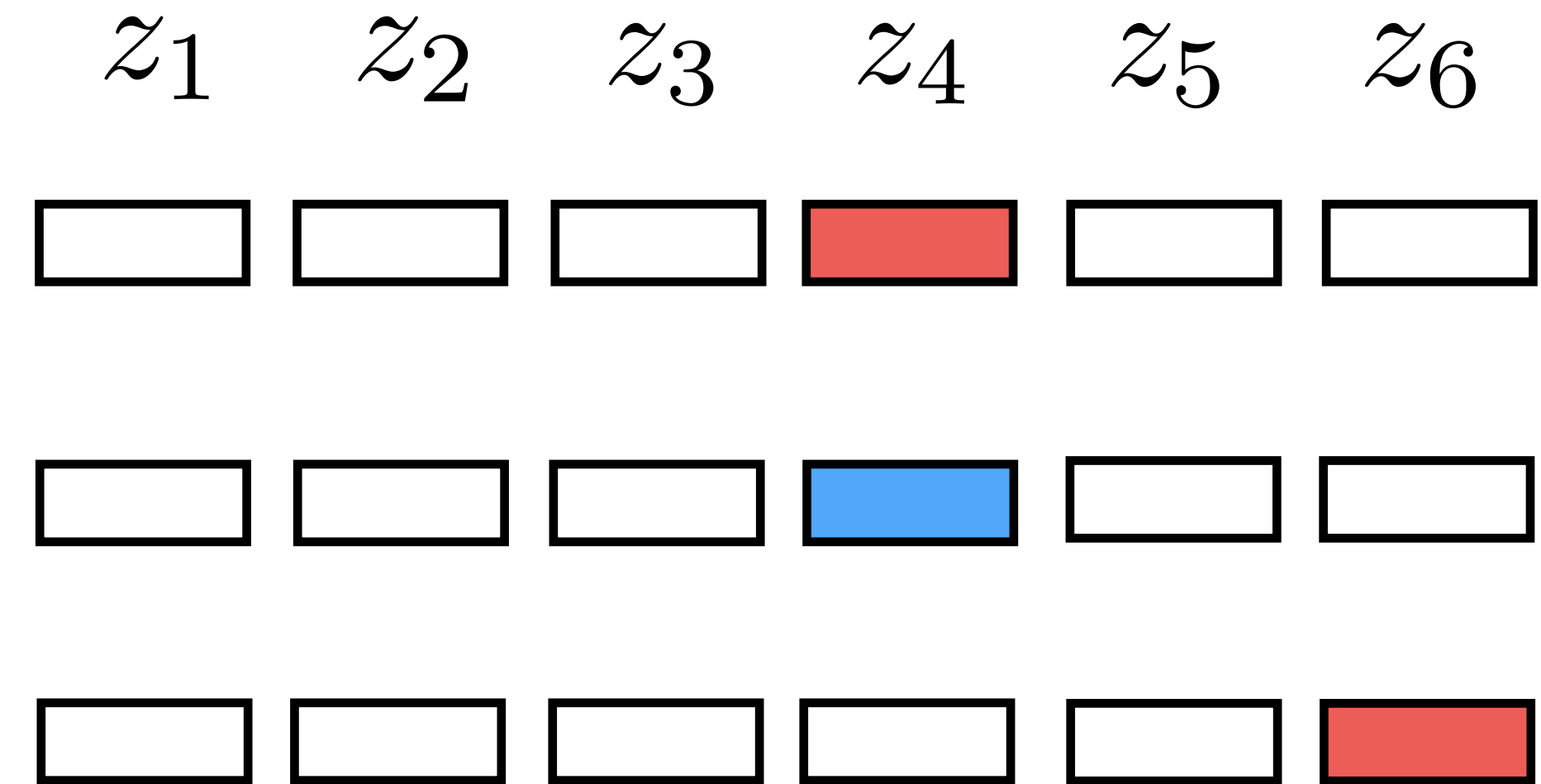


# Autoencoders

the movie was **good** </s> </s>

the movie was **great** </s> </s>

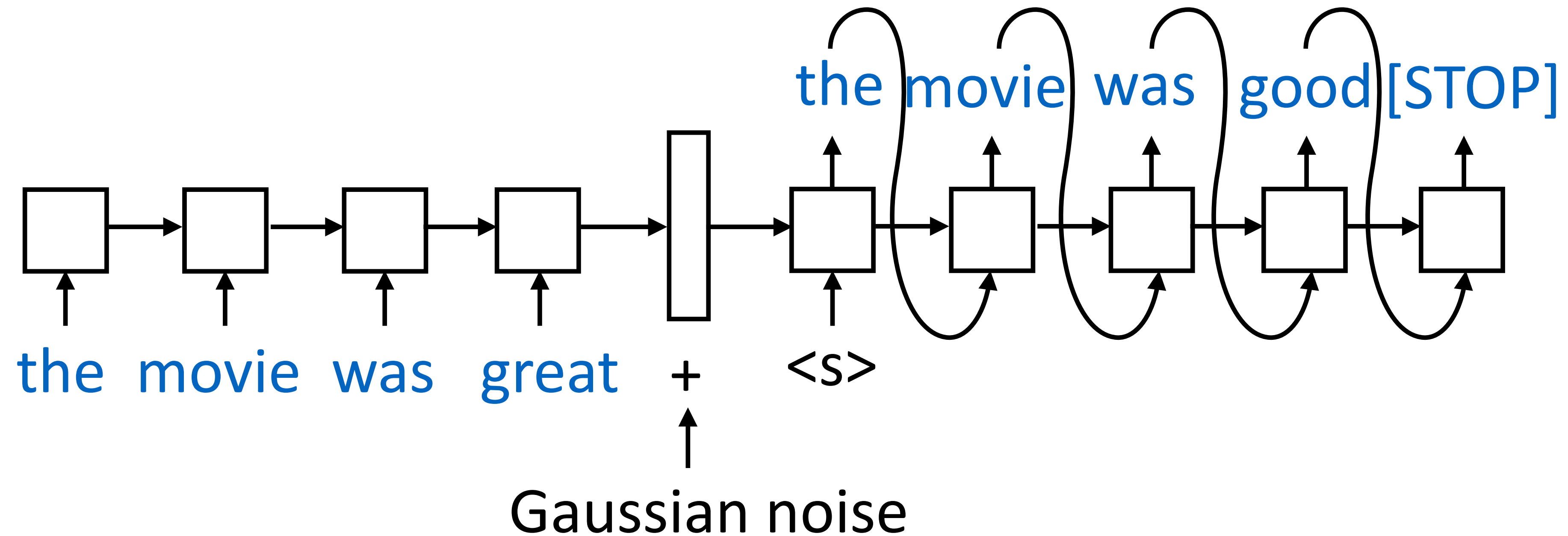
I thought the film was **good**



- ▶ Can an LSTM learn to do this?
- ▶ Yes!
- ▶ Want continuous semantic structure in the latent space: nearby points should have similar meaning



# Autoencoders



- ▶ During training, add Gaussian noise and force the network to predict
- ▶ Same computation graph as VAE with reparameterization, add KL term to make the objective the same
- ▶ Inference network ( $q$ ) is the encoder and generator is the decoder



# LSTM VAEs

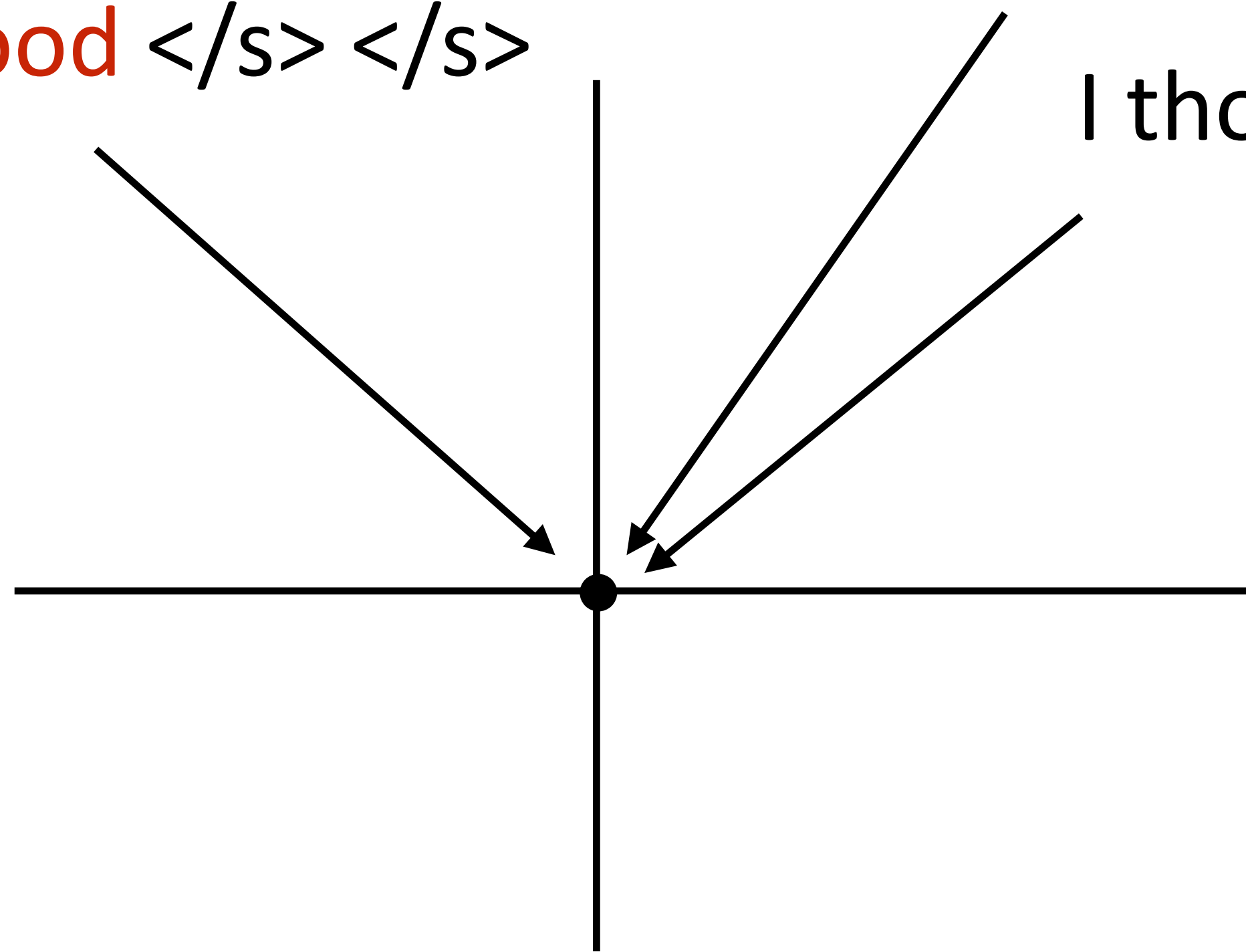
$$\mathbb{E}_{q(z|\mathbf{x})} [\log P(\mathbf{x}|z, \theta)] + \text{KL}(q(z|\mathbf{x}) || P(z))$$

- ▶ Train this up; what happens?

the movie was **great** </s> </s>

the movie was **good** </s> </s>

I thought the film was **good**

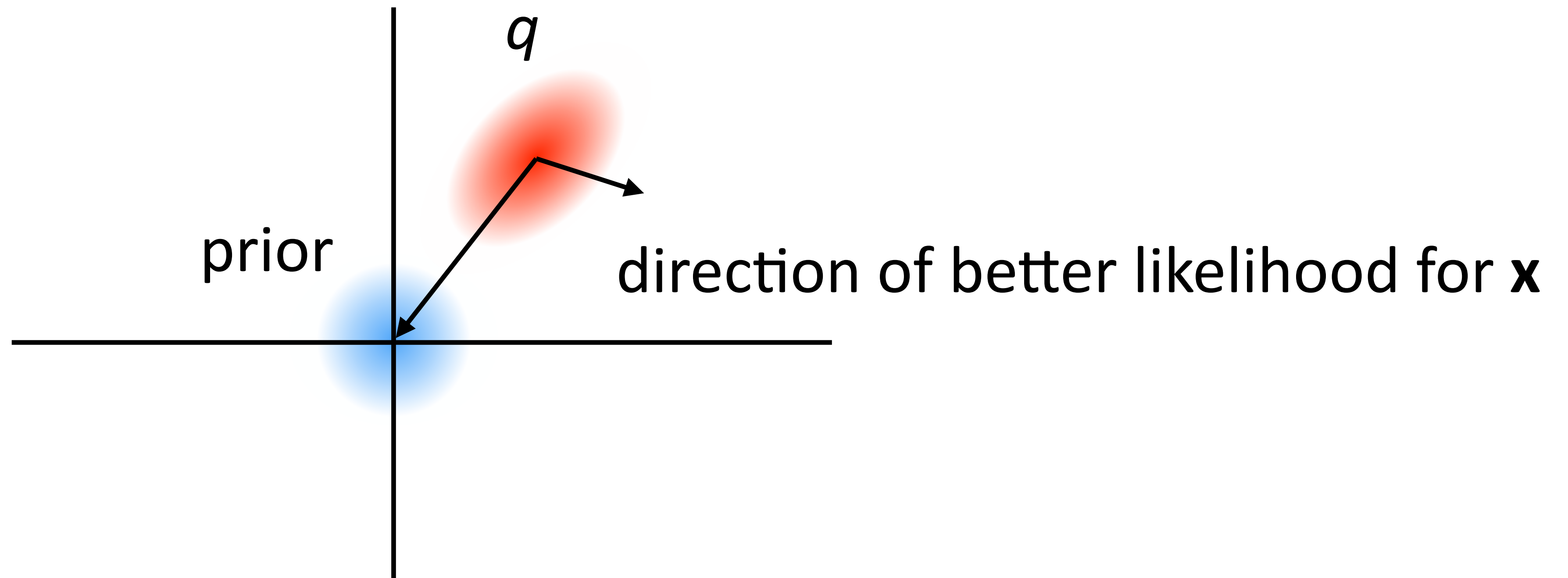




# KL Collapse

$$\mathbb{E}_{q(z|\mathbf{x})} [\log P(\mathbf{x}|z, \theta)] + \text{KL}(q(z|\mathbf{x}) || P(z))$$

- What does gradient encourage latent space to do?

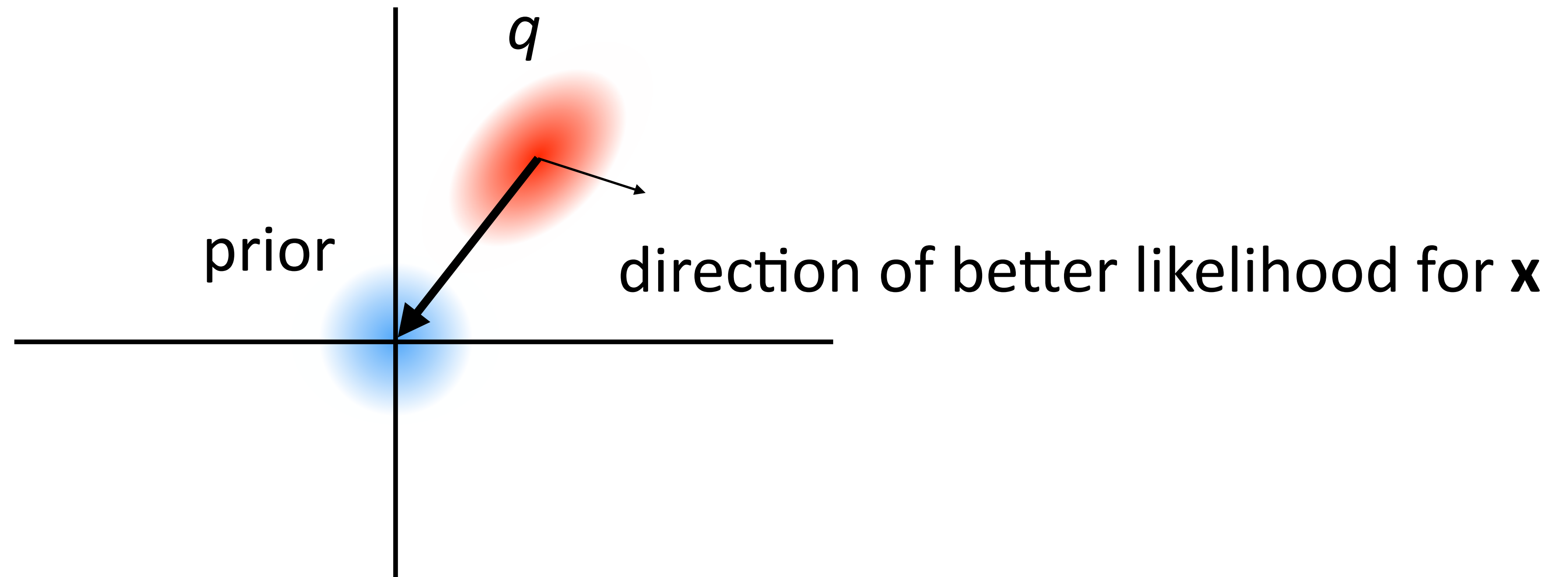




# KL Collapse

$$\mathbb{E}_{q(z|\mathbf{x})} [\log P(\mathbf{x}|z, \theta)] + \text{KL}(q(z|\mathbf{x}) || P(z))$$

- What does gradient encourage latent space to do?



- In reality, the likelihood signal is very weak,  $z$  is set to 0

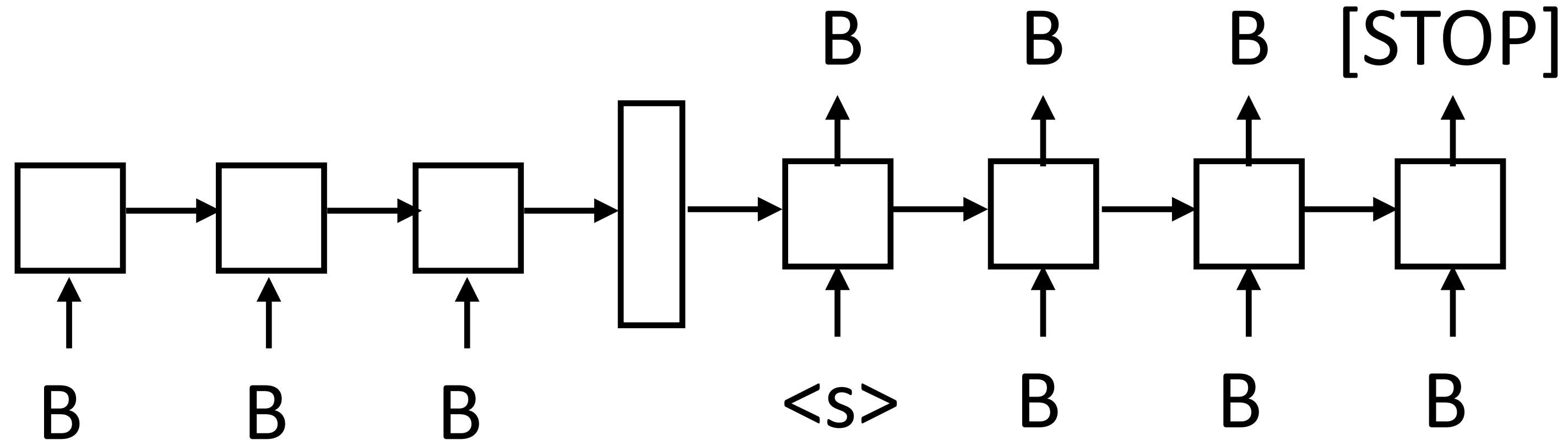
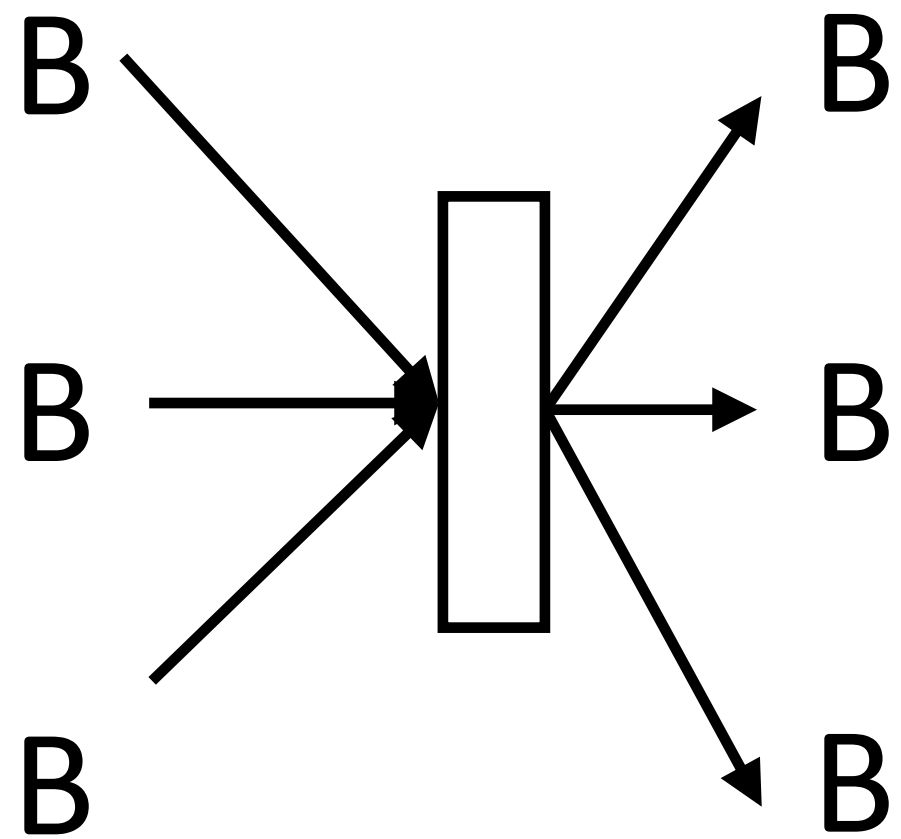




# A Tale of Two Decoders

$$\mathbb{E}_{q(z|\mathbf{x})} [\log P(\mathbf{x}|z, \theta)] + \text{KL}(q(z|\mathbf{x}) || P(z))$$

- Suppose vocab is {A, B, C}. Sentences are either AAA, BBB, or CCC



LL =  $(1/3)^3$  if input is ignored

LL = 1/3 if input is ignored

- NVDM: using  $z$  can help a lot

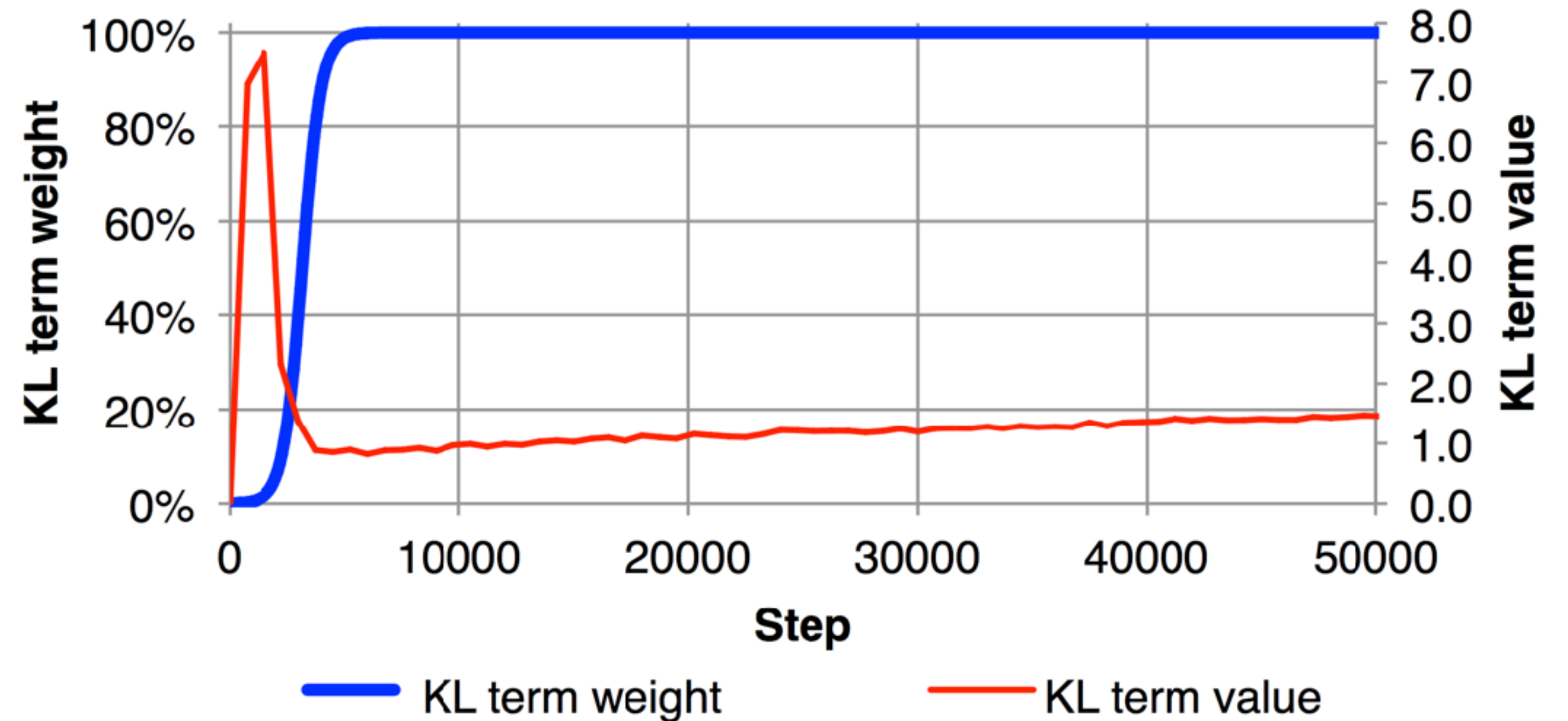
- LSTM: can get decent likelihood ignoring  $z$  entirely



# A Tale of Two Decoders

$$\mathbb{E}_{q(z|\mathbf{x})} [\log P(\mathbf{x}|z, \theta)] + \text{KL}(q(z|\mathbf{x}) || P(z))$$

- Solution: anneal KL term during learning



- Model initially uses  $z$  a lot ( $q$  gets far from the prior), then as KL term is turned up the prior balances it more

Bowman et al. (2016)



# Results

- ▶ Train autoencoder on the Penn Treebank (pretty small corpus for language modeling purposes)

Model	Standard			
	Train NLL	Train PPL	Test NLL	Test PPL
<b>RNNLM</b>	100 —	95	<b>100</b> —	<b>116</b>
<b>VAE</b>	98 (2)	100	101 (2)	119

- ▶ Doesn't really improve perplexities over RNNLM: confirms that RNN is pretty good at modeling the space





# Results

---

INPUT	<b>we looked out at the setting sun .</b>	<b>i went to the kitchen .</b>	<b>how are you doing ?</b>
MEAN	<i>they were laughing at the same time .</i>	<i>i went to the kitchen .</i>	<i>what are you doing ?</i>
SAMP. 1	<i>ill see you in the early morning .</i>	<i>i went to my apartment .</i>	<i>“ are you sure ?</i>
SAMP. 2	<i>i looked up at the blue sky .</i>	<i>i looked around the room .</i>	<i>what are you doing ?</i>
SAMP. 3	<i>it was down on the dance floor .</i>	<i>i turned back to the table .</i>	<i>what are you doing ?</i>

---

- Encode sentence, sample from  $q$ , generate from those samples

---

**“ 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 .**

---

**he was silent for a long moment .**

*he was silent for a moment .*

*it was quiet for a moment .*

*it was dark and cold .*

*there was a pause .*

**it was my turn .**

---

- Encode two samples, generate from points interpolated between the two samples



# Takeaways

---

- ▶ VAE is a framework for training deep generative models
- ▶ VAE can be seen as either a principled variational method motivated by a lower bound or simply an ad-hoc trick to make latent spaces more continuous
- ▶ Some tricks to get these models to train well
- ▶ Generative objective ensures that the latent space  $z$  has interesting and coherent semantics; lets us sample from  $z$  and generate instances from the data manifold