

Latent Variable Models and Iterative Refinement for **Non-Autoregressive** Neural Machine Translation (= Parallel)

Jason Lee
New York University

Based on joint work with Kyunghyun Cho, Elman Mansimov, Raphael Shu and Hideki Nakayama

Fall 2017, when I started my PhD....



So Jason, what do you want to work on?



Umm... do you have any ideas?

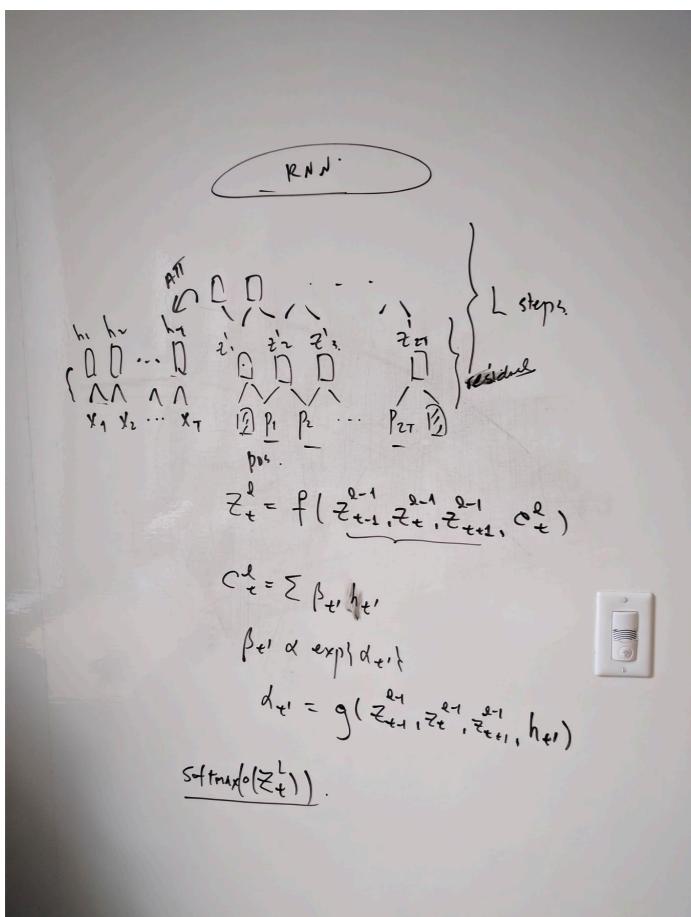
Slow decoding speed is a big problem in MT. Why don't you start looking into Neural GPU?

Published as a conference paper at ICLR 2016

NEURAL GPUs LEARN ALGORITHMS

Lukasz Kaiser & Ilya Sutskever

Google Brain {lukaszkaiser,ilyasu}@google.com



Jonas Gehring
Michael Auli
David Grangier
Denis Yarats
Yann N. Dauphin
Facebook AI Research

Convolutional Sequence to Sequence Learning

Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaiser@google.com

Model class, learning and parameterization

Inference algorithm

Latent Variable Models and Iterative Refinement for Non-Autoregressive Neural Machine Translation

Problem

Jason Lee
New York University

Machine Translation

- Conditional sequence generation:
 - Given an input sequence $(x_1, \dots, x_{t'})$, generate a corresponding output sequence (y_1, \dots, y_t) .
- A broad range of tasks can be framed similarly:
 - Dialogue Modelling
 - Image denoising, Image inpainting
 - Automatic Speech Recognition / Text-to-speech

Sequence Learning

- Most neural sequence models are trained to maximize the log-likelihood of the training data

$$\hat{\theta} = \operatorname{argmax}_{\theta} \sum_{n=1}^N \log p_{\theta}(y_{1:t}^n | x_{1:t'}^n)$$

Autoregressive Sequence Models

- Learning: factorize the joint probability of a sequence into a product of conditionals. Exact likelihood.

$$\log p(y_{1:t} | x_{1:t'}) = \sum_{i=1}^t \log p(y_i | y_{<i}, x_{1:t'})$$

$$\begin{aligned} & p(\text{I love BTS} \mid \text{저는 방탄소년단을 좋아합니다}) \\ &= p(\text{I} \mid \text{저는 방탄소년단을 좋아합니다}) \\ &\times p(\text{love} \mid \text{I, 저는 방탄소년단을 좋아합니다}) \\ &\times p(\text{BTS} \mid \text{I love, 저는 방탄소년단을 좋아합니다}) \end{aligned}$$

- Parameterization (conditional distribution):
 - RNNs (LSTMs, GRUs), ConvNet, Transformers

Autoregressive Models: Inference

$$\hat{y} = \operatorname{argmax}_{y_1, \dots, y_t} p(y_1, \dots, y_t | x_{1:t'})$$

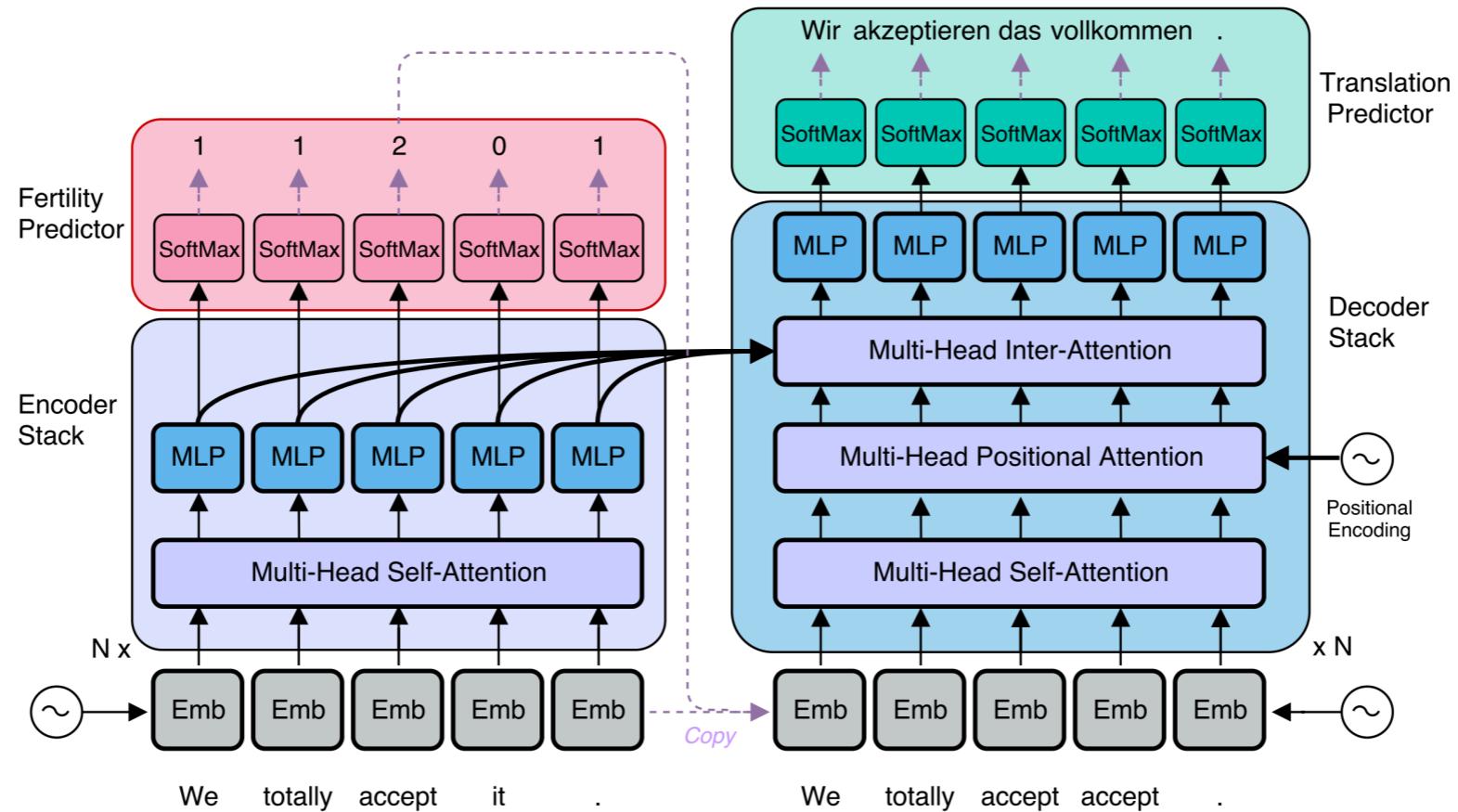
- Exact inference is intractable as the size of the search space grows exponentially w.r.t to t .
- Approximate inference:
 - Greedy search, beam search

Decoding latency grows linearly

- Decoding from an autoregressive model is slow, i.e. latency scales w.r.t the output sequence length.
- Impossible to parallelize decoding at inference time, despite modern parallel hardware
- Presents significant challenges for:
 - Document-level machine translation (~1k subword segments, Junczys-Dowmunt, 2019)
 - Text-to-speech (32,000 samples/second)

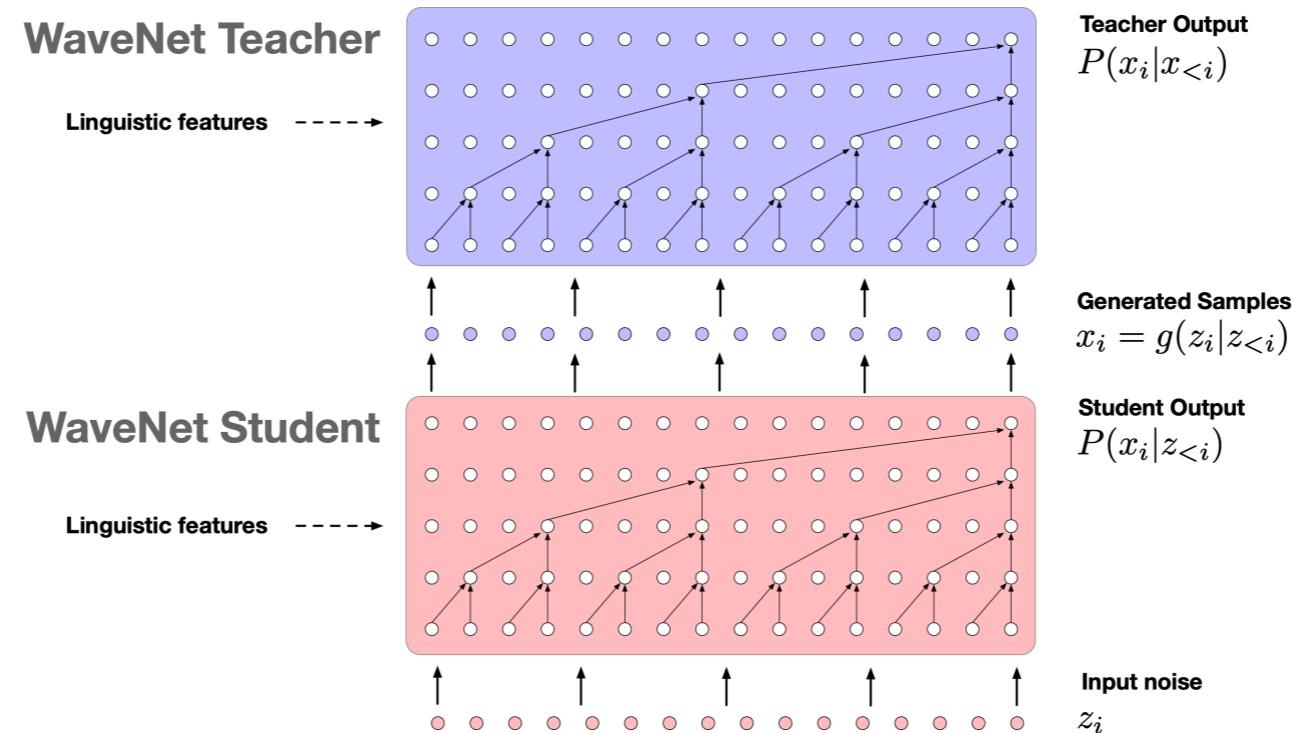
Related work on Non-Autoregressive Sequence Models

(1) Non-Autoregressive NMT



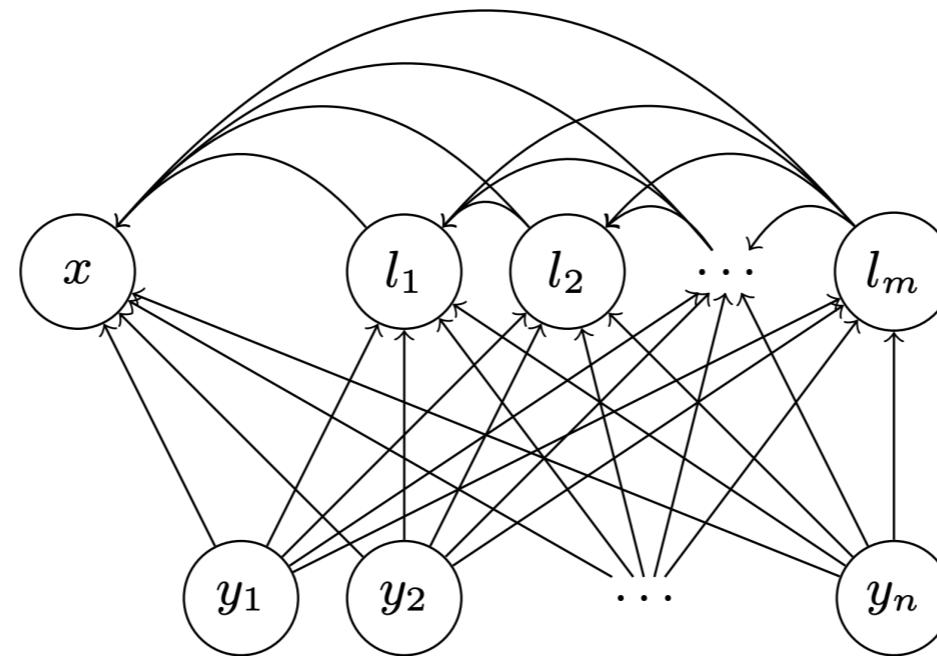
- Treat alignment (=fertility) as a latent variable
 - Supervise alignment prediction model using external model
- (1) Distillation, (2) Fertility as Latent Variable, (3) RL Finetuning + Reverse KL
- 3x speedup on WMT'14 En-De at the cost of 3-4 BLEU

(2) Parallel WaveNet for TTS



- Non-autoregressive Inverse Autoregressive Flow
 - $\log p_x(x) = \log p_z(f(x)) + \log \left| \det\left(\frac{\partial f}{\partial x}\right) \right|, \quad x_t = z_t \cdot s_\theta(z_{<t}) + \mu_\theta(z_{<t})$
 - Trained with distillation from an autoregressive WaveNet
 - Auxiliary perceptual loss terms
- 1000x speedup vs. WaveNet

(3) Semi-autoregressive NMT



- **Discrete autoencoder (non-autoregressive)** for the target sentence
 - One discrete latent variable autoencodes 8 neighbouring target tokens
- **Autoregressive prior** from the source to the discrete latent variables.
- Discretization : VQ-VAE / semantic hashing
- 4x speedup on WMT'14 En-De at the cost of 1 BLEU
 - Using distillation (Roy et al., 2018)

(4) Non-autoregressive NMT with Masked Language Models

- $P(y_{\text{mask}} | y_{\text{obs}}, x)$: predict masked tokens from observed ones
 - Training: y_{mask} selected at random
 - Decoding: mask $1 - t/T$ of target tokens with the lowest probability scores (at the t -th refinement step).
- Length prediction can be batched with l different samples
 - $T = 10, l = 2$: **27.03 BLEU** on WMT'14 En->De (1.3 speedup)
 - $T = 4, l = 2$: **25.51 BLEU** on WMT'14 En->De (3x speedup)

Our Approach

- We'd like to factorize the output distribution:

$$p(y_1, \dots, y_t | x) = \prod_{i=1}^t p(y_i | x)$$

- Main challenge: how do we capture the dependencies between the output variables y_i ?

Long range dependencies

Concrete is a versatile substance that has been used for centuries in a wide range of construction projects. Ready-mix concrete suppliers help build **home foundations, driveways, roadways, bridges, dams, buildings** and more. Self-compacting concrete improves **strength and durability**, while reducing energy use due to its ability to minimize **temperature fluctuations** in ~~buildings~~ over the course of the day.

Short dependencies

Research Questions

1. Sequence Learning

- How can we learn the distribution of output sequences (y_1, \dots, y_t) non-autoregressively?

2. Sequence Generation

- Having trained a model, how can we perform inference efficiently? (find the best output y given an input x)?

Research Questions

1. Sequence Learning

- How can we learn the distribution of output sequences (y_1, \dots, y_t) non-autoregressively?
- Let **latent variables** capture the dependencies, so the distribution of the output sequence can be factorized given the latent variables
- $\log p(y_{1:t}, z | x) = \sum_{i=1}^t \log p(y_i, z | x)$

Research Questions

2. Sequence Generation

- Having trained a model, how can we perform inference efficiently? (find the best output y given an input x)?
- $\hat{y} = \operatorname{argmax}_{y_{1:t}} \log \int_z p(y_{1:t}, z | x_{1:t'})$
- Use **iterative refinement**

Latent Variable Models and Iterative Refinement

(x, y) : (input, output) sequence

z : latent variable capturing the output sequence

$$p(y|x) = \sum_z p(y, z|x) \text{ for discrete latent variables } z$$

1. Iterative Refinement in a Discrete Space (Lee, Mansimov and Cho, 2018).

$$p(y|x) = \int p(y, z|x) dz \text{ for continuous latent variables } z$$

2. Iterative Refinement in a Hybrid Space (Shu, Lee, Nakayama and Cho, 2020).

3. Iterative Refinement in a Continuous Space (Lee, Shu and Cho 2020).

Latent Variable Models and Iterative Refinement

(x, y) : (input, output) sequence

z : latent variable capturing the output sequence

$$p(y|x) = \sum_z p(y, z|x) \text{ for discrete latent variables } z$$

1. Iterative Refinement in a Discrete Space (Lee, Mansimov and Cho, 2018).

$$p(y|x) = \int p(y, z|x) dz \text{ for continuous latent variables } z$$

2. Iterative Refinement in a Hybrid Space (Shu, Lee, Nakayama and Cho, 2020).
3. Iterative Refinement in a Continuous Space (Lee, Shu and Cho 2020).

Discrete Latent Variable Models

We introduce L discrete latent variables z_1, \dots, z_L to capture the dependencies between the output variables.

$$\begin{aligned} p(y|x) &= \sum_{z_{1:L}} p(y, z_{1:L}|x) \\ &= \sum_{z_{1:L}} \prod_{l=1}^{L+1} p(z_l | z_{<l}, x), \text{ setting } y = z_{L+1} \end{aligned}$$

Discrete Latent Variable Models

We introduce L discrete latent variables z_1, \dots, z_L to capture the dependencies between the output variables.

$$\begin{aligned} p(y|x) &= \sum_{z_{1:L}} p(y, z_{1:L}|x) \\ &= \sum_{z_{1:L}} \prod_{l=1}^{L+1} p(z_l | z_{<l}, x), \text{ setting } y = z_{L+1} \end{aligned}$$

We constrain z as discrete latent variables with the same dimensionality as the target sentence. $z \in \{0, \dots, |V_{trg}|\}^T$

Approximations to Sidestep the Marginalization

We introduce L discrete latent variables z_1, \dots, z_L to capture the dependencies between the output variables.

$$\begin{aligned} p(y | x) &= \sum_{z_{1:L}} p(y, z_{1:L} | x) \\ &= \sum_{z_{1:L}} \prod_{l=1}^{L+1} p(z_l | z_{<l}, x), \text{ setting } y = z_{L+1} \end{aligned}$$

As this marginalization is intractable, we consider two approximations:

(1) Deterministic Approximation

$$p(y|x) = \sum_{z_{1:L}} p(y, z_{1:L} | x)$$
$$= \sum_{z_{1:L}} \prod_{l=1}^{L+1} p(z_l | z_{<l}, x)$$

$$\geq \prod_{l=1}^{L+1} p(\hat{z}_l | \hat{z}_{<l}, x), \text{ where } \hat{z}_l = \operatorname{argmax}_z p(z | \hat{z}_{<l}, x)$$

(2) Stochastic Approximation

$$p(y|x) = \sum_{z_{1:L}} p(y, z_{1:L} | x)$$

$$= \sum_{z_{1:L}} \prod_{l=1}^{L+1} p(z_l | z_{<l}, x)$$

$$\approx \prod_{l=1}^{L+1} p(\hat{z}_l | \hat{z}_{<l}, x), \text{ where } \hat{z}_l \sim p(z | \hat{z}_{<l}, x)$$

This gives an unbiased estimate of the marginal.

(2) Stochastic Approximation

$$p(y|x) = \sum_{z_{1:L}} p(y, z_{1:L} | x)$$

$$= \sum_{z_{1:L}} \prod_{l=1}^{L+1} p(z_l | z_{<l}, x)$$

$$\approx \prod_{l=1}^{L+1} \boxed{p(\hat{z}_l | \hat{z}_{<l}, x)}, \text{ where } \hat{z}_l \sim p(z | \hat{z}_{<l}, x)$$

Can be viewed as a denoising autoencoder.

$$\prod_{l=1}^{L+1} p(y^* | \tilde{y}_l, x)$$

Denoising Autoencoder

- $x \in \mathbb{R}^d$, L2 squared loss $l(\cdot)$, stochastic corruption function $c(\cdot)$, denoising function $d(\cdot)$

$$\mathcal{L}_{DAE} = \mathbb{E}[l(x, d(c(x)))]$$

- Alain and Bengio showed that DAE estimates the score (gradient of the log density.)
- Our corruption operations: (1) swap neighbouring tokens, (2) replace token with a random token, (3) repeat tokens

Joint Training Objective

$$J_{LVM}(\theta) = - \sum_{l=1}^{L+1} \log p_\theta(\hat{z}_l | \hat{z}_{<l}, x)$$

$$J_{DAE}(\theta) = - \sum_{l=1}^{L+1} \log p_\theta(y^* | \tilde{y}_l, x)$$

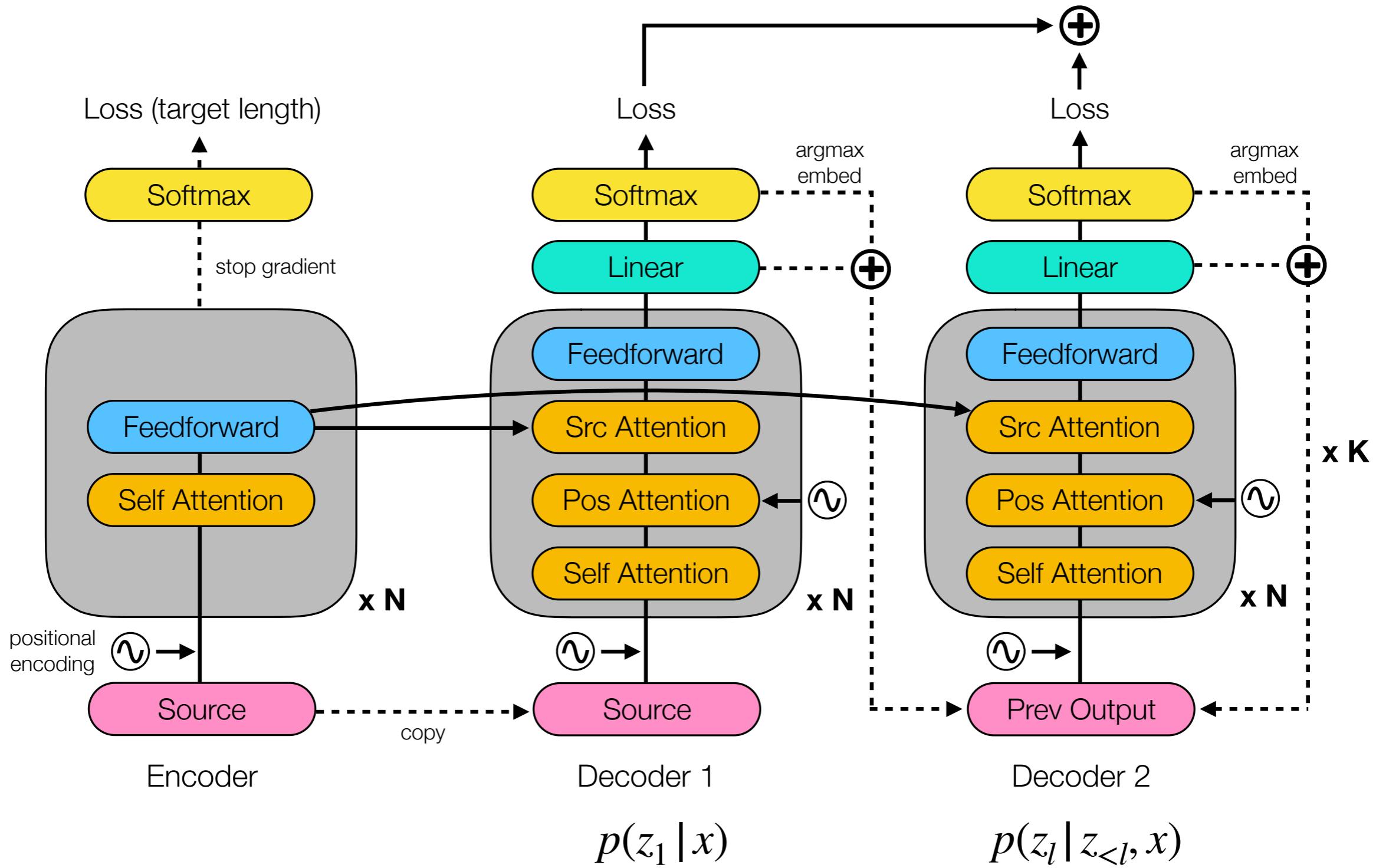
$$J(\theta) = \alpha_l \cdot J_{LVM}(\theta) + (1 - \alpha_l) \cdot J_{DAE}(\theta)$$

We train the length predictor (2-layer MLP from source hidden states) jointly with the rest of the model.

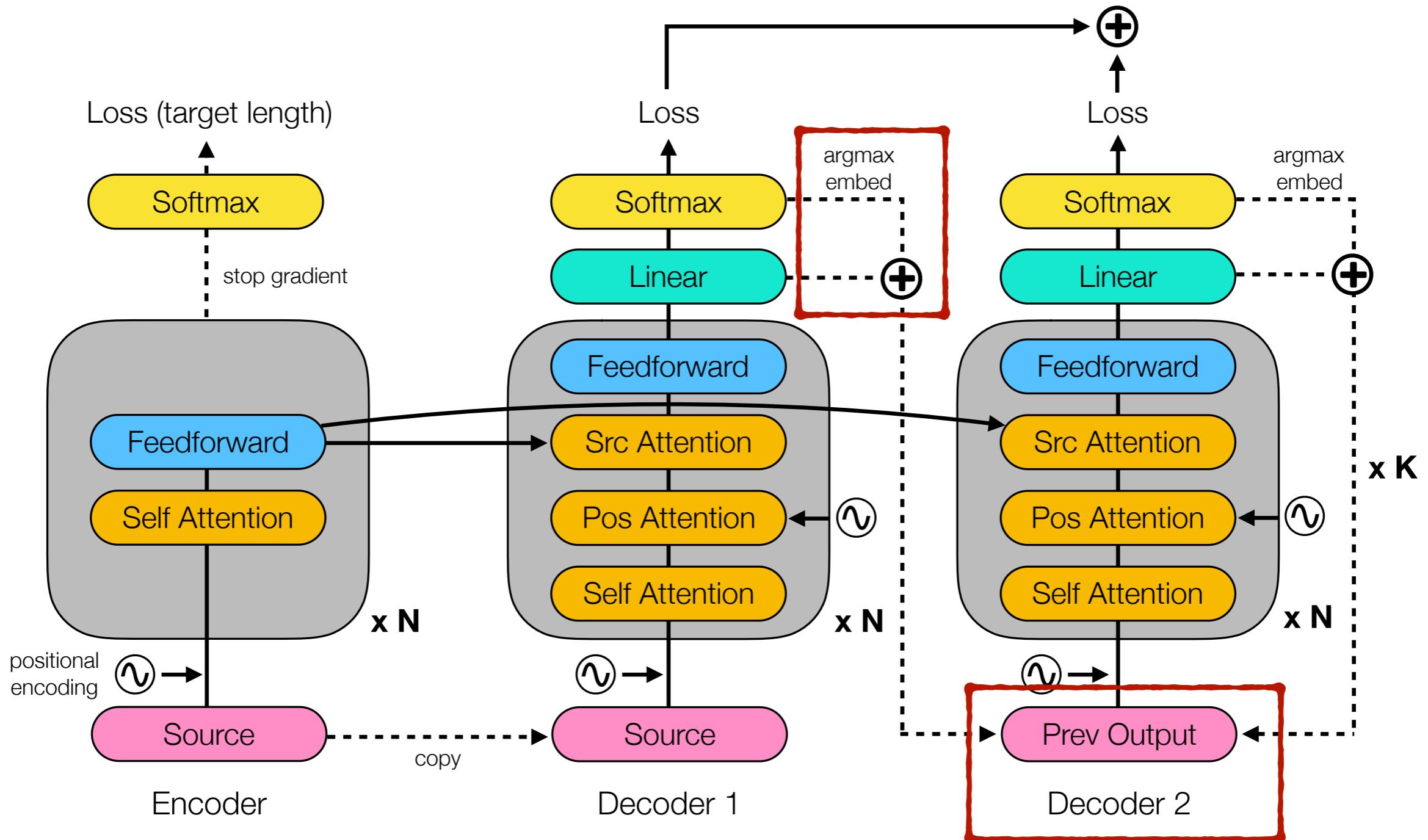
Deterministic Inference

1. $\hat{T} = \operatorname{argmax}_T \log p(T|X)$: predict the target length
2. $\hat{t}_t^0 = \operatorname{argmax}_{y_t} \log p(y_t^0|x)$: predict the first latent variables
3. $\hat{t}_t^l = \operatorname{argmax}_{y_t} \log p(y_t^l|\hat{Y}_l, x)$: refine the latent variables
4. Repeat (3) for either (1) K times, or (2) until a certain termination criterion has been met.

Network Architecture

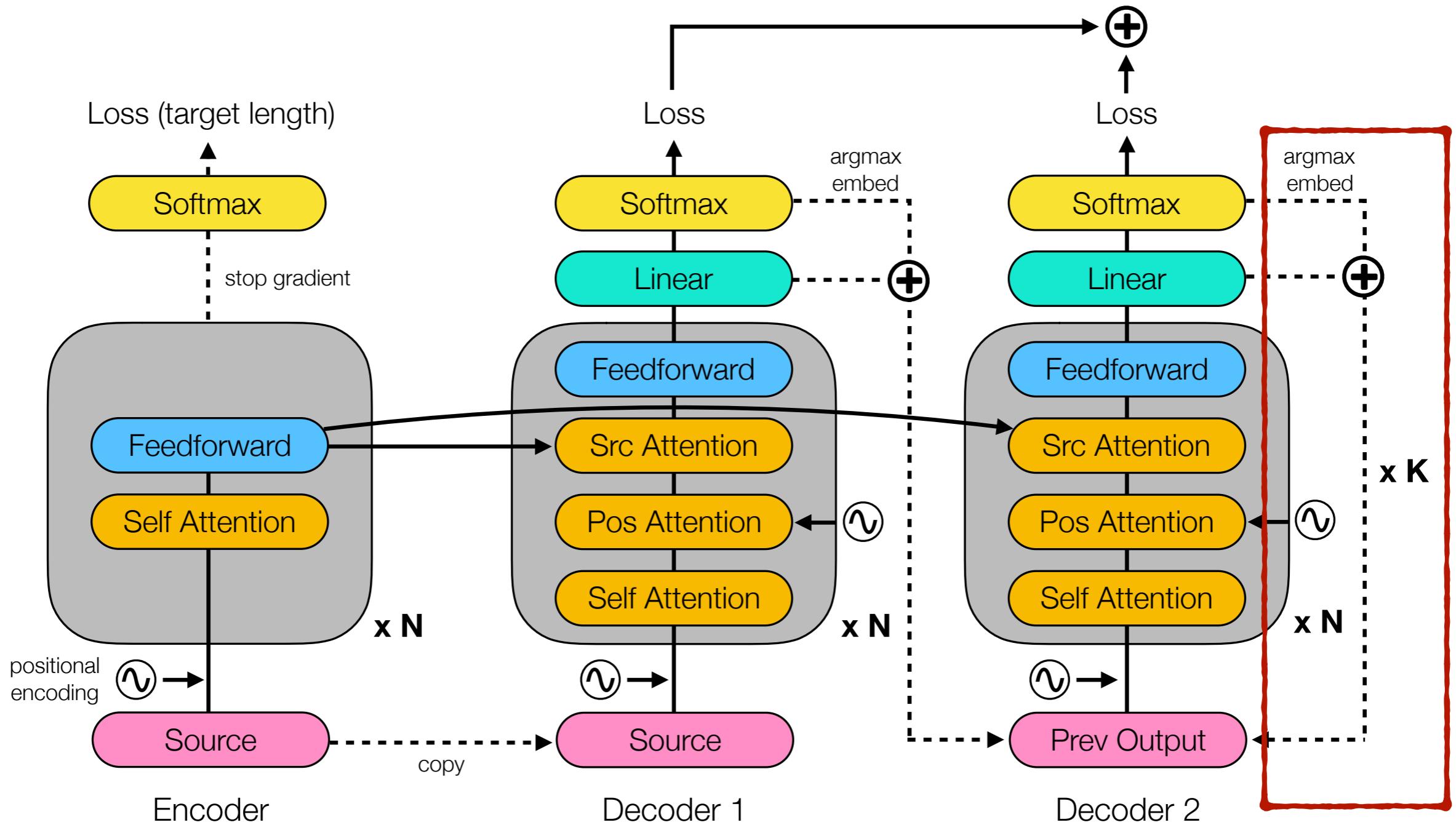


Network Architecture



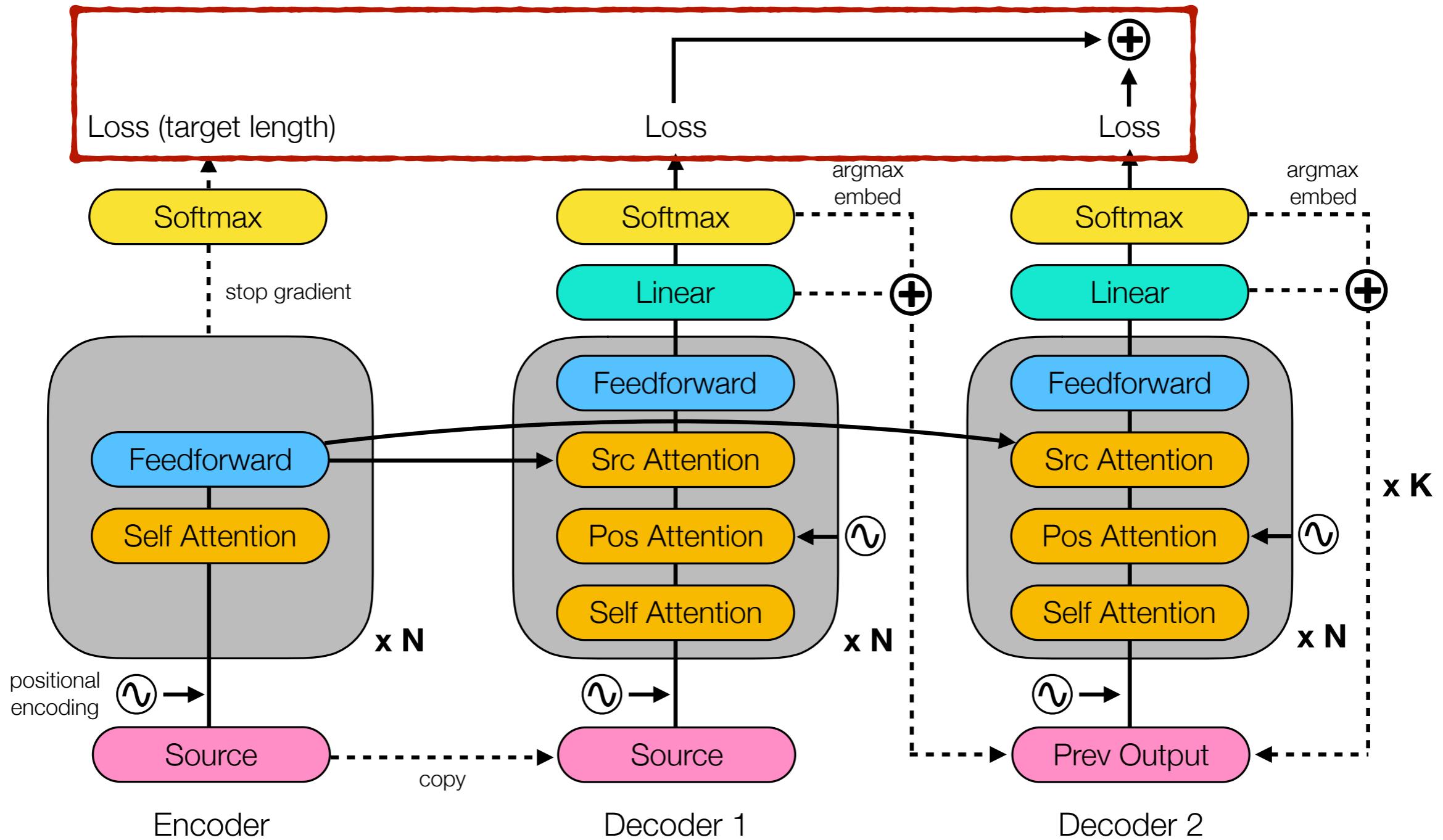
- We use embeddings + hidden activations from the previous layer as input to the next

Network Architecture



- Refine K (=4) times during training, potentially use $> K$ steps during inference.
- **Adaptive** decoding : refine until output converges

Network Architecture

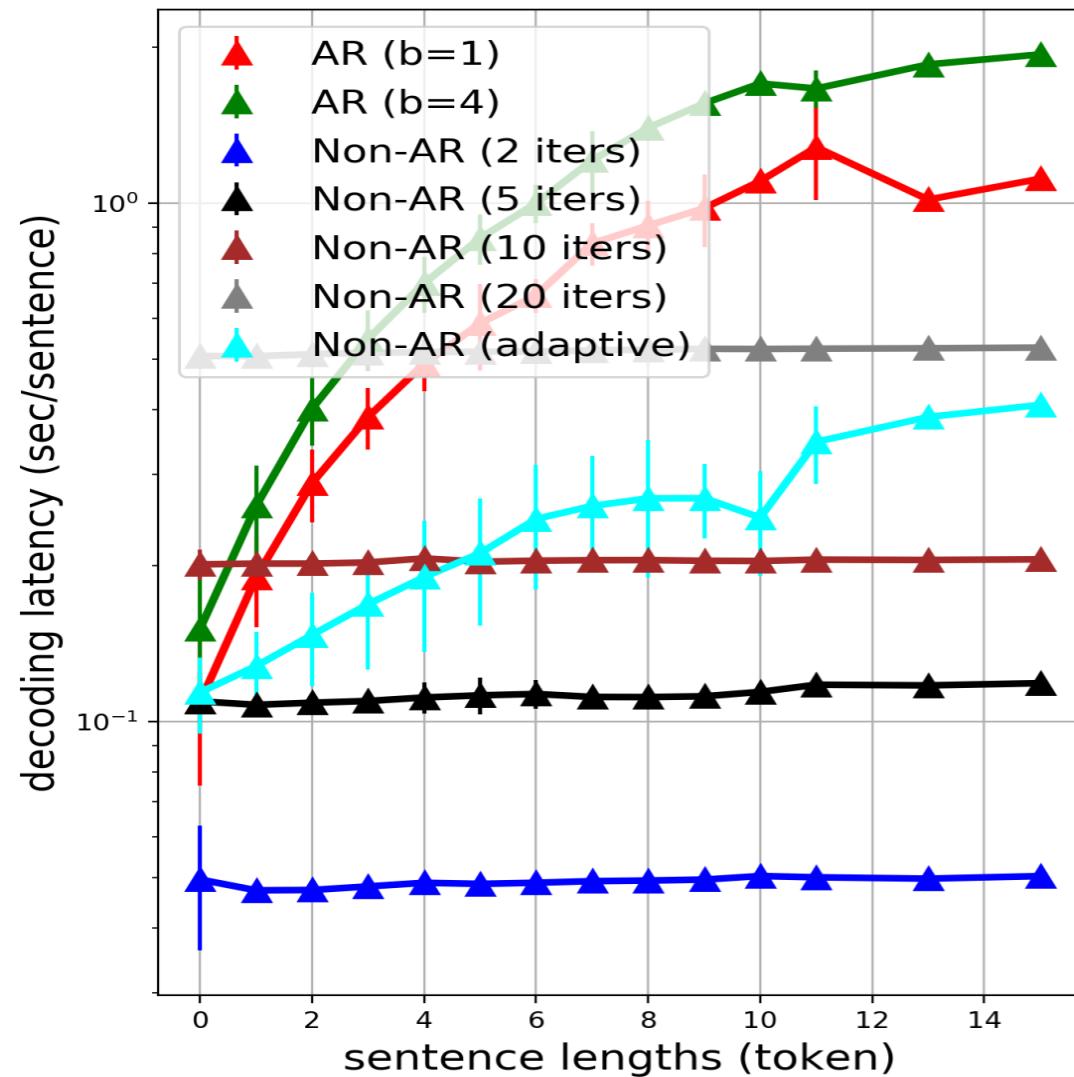


- Loss from all iterations and the length prediction model are summed

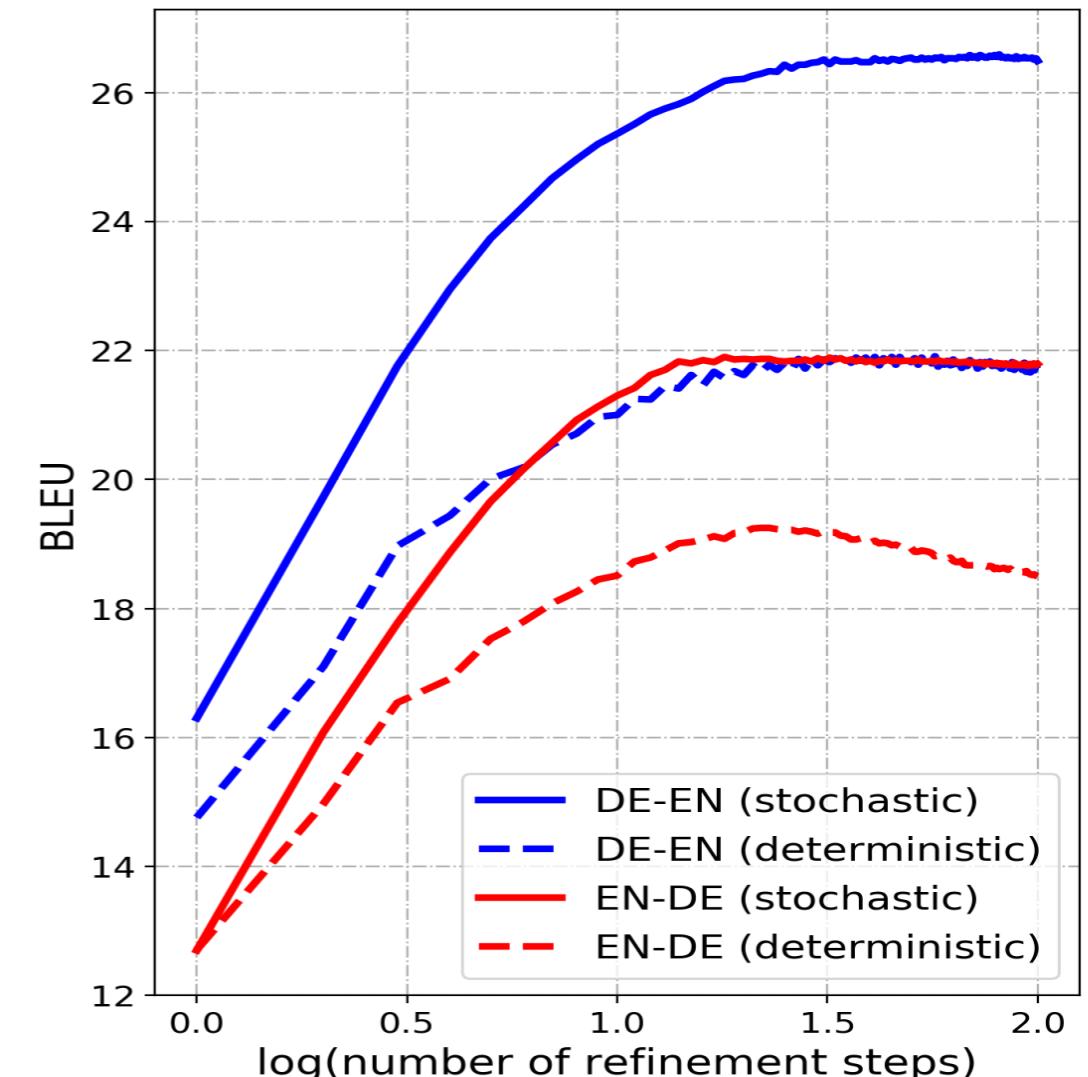
Non-autoregressive NMT on WMT'14 En-De

		BLEU(%)	Gap	Speed
(Gu et al., 2018)	Baseline	23.4		
	Gu's Model (NAT + NPD S=100)	19.1	-4.3	2.3x
(Lee et al., 2018)	Baseline	24.5		
	Lee's Model (Adaptive NAR)	21.5	-3.0	1.9x
(Kaiser et al., 2018)	Kaiser's baseline	23.5		
	Kaiser's Model (LT, Improved semhash)	19.8	-3.7	3.8x

Latency vs. Sequence Length



Refinement Steps vs. Performance



Qualitative Results : IWSLT'16 De -> En

Src	seitdem habe ich sieben Häuser in der Nachbarschaft mit den Lichtern versorgt und sie funktionierenen wirklich gut .
Iter 1	and I 've <u>been seven homes since</u> <u>in neighborhood</u> with the lights and they 're really functional .
Iter 2	and I 've <u>been seven homes in the</u> neighborhood with the lights , and they 're a really functional .
Iter 4	and I 've <u>been seven homes in</u> neighborhood with the lights , and they 're a really functional .
Iter 8	and I 've been providing seven homes in the neighborhood with the lights and they 're a really functional .
Iter 20	and I 've been providing seven homes in the neighborhood with the lights , and they 're a very good functional .
Ref	since now , I 've set up seven homes around my community , and they 're really working .
Src	er sah sehr glücklich aus , was damals ziemlich ungewöhnlich war , da ihn die Nachrichten meistens deprimierten .
Iter 1	he looked very happy , which was pretty <u>unusual the</u> , because <u>the news was were</u> usually depressing .
Iter 2	he looked very happy , which was pretty <u>unusual at the</u> , because <u>the news was s</u> depressing .
Iter 4	he looked very happy , which was pretty <u>unusual at the</u> , because <u>news was mostly</u> depressing .
Iter 8	he looked very happy , which was pretty <u>unusual at the time</u> because <u>the news was mostly</u> depressing .
Iter 20	he looked very happy , which was pretty <u>unusual at the time</u> , because <u>the news was mostly</u> depressing .
Ref	there was a big smile on his face which was unusual then , because the news mostly depressed him .
Src	furchtlos zu sein heißt für mich , heute ehrlich zu sein .
Iter 1	to be , <u>for me</u> , to be honest today .
Iter 2	to be <u>fearless</u> , <u>me</u> , is to be honest today .
Iter 4	to be <u>fearless for me</u> , is to be honest today .
Iter 8	to be <u>fearless for me</u> , <u>me</u> to be honest today .
Iter 20	to be <u>fearless for me</u> , is to be honest today .
Ref	so today , for me , being fearless means being honest .

Table 4: Three sample De→En translations from the non-autoregressive sequence model. Source sentences are from the dev set of IWSLT'16. The first iteration corresponds to Decoder 1, and from thereon, Decoder 2 is repeatedly applied. Sub-sequences with changes across the refinement steps are underlined.

Qualitative Results: MS COCO



Generated Caption

Iter 1 a yellow bus parked on parked in of parking road .
Iter 2 a yellow and black on parked in a parking lot .
Iter 3 a yellow and black bus parked in a parking lot .
Iter 4 a yellow and black bus parked in a parking lot .

Reference Captions

a tour bus is parked on the curb waiting
city bus parked on side of hotel in the rain .
bus parked under an awning next to brick sidewalk
a bus is parked on the curb in front of a building .
a double decked bus sits parked under an awning



Generated Caption

Iter 1 a woman standing on playing tennis on a tennis racquet .
Iter 2 a woman standing on a tennis court a tennis racquet .
Iter 3 a woman standing on a tennis court a a racquet .
Iter 4 a woman standing on a tennis court holding a racquet .

Reference Captions

a female tennis player in a black top playing tennis
a woman standing on a tennis court holding a racquet .
a female tennis player preparing to serve the ball .
a woman is holding a tennis racket on a court
a woman getting ready to reach for a tennis ball on the ground

Latent Variable Models and Iterative Refinement

(x, y) : (input, output) sequence

z : latent variable capturing the output sequence

$$p(y|x) = \sum_z p(y, z|x) \text{ for discrete latent variables } z$$

1. Iterative Refinement in a Discrete Space (Lee, Mansimov and Cho, 2018).

$$p(y|x) = \int p(y, z|x) dz \text{ for continuous latent variables } z$$

2. Iterative Refinement in a Hybrid Space (Shu, Lee, Nakayama and Cho, 2020).

3. Iterative Refinement in a Continuous Space (Lee, Shu and Cho 2020).

Continuous Latent Variable Models

$$\log p_{\theta}(y|x) = \log \int p_{\theta}(y, z|x) dz$$

Continuous Latent Variable Models

$$\log p_{\theta}(y|x) = \log \int p_{\theta}(y, z|x) dz$$

$$= \log \int \frac{p_{\theta}(y|z, x) p_{\theta}(z|x)}{q_{\phi}(z|y, x)} q_{\phi}(z|y, x) dz$$

$$= \log \mathbb{E}_{z \sim q_{\phi}} \left[\frac{p_{\theta}(y|z, x) p_{\theta}(z|x)}{q_{\phi}(z|y, x)} \right]$$

$$\geq \mathbb{E}_{z \sim q_{\phi}} [\log p_{\theta}(y|z, x) + \log p_{\theta}(z|x) - \log q_{\phi}(z|y, x)]$$

Decoder

Prior

Approximate Posterior

$$= \mathbb{E}_{z \sim q_{\phi}} [\log p_{\theta}(y|z, x)] - \text{KL}[q_{\phi}(z|y, x) || p_{\theta}(z|x)]$$

$$= \text{ELBO}(x, y, \theta, \phi)$$

Continuous Latent Variable Models

$$\log p_{\theta}(y|x) = \log \int p_{\theta}(y, z|x) dz$$

$$= \log \int \frac{p_{\theta}(y|z, x) p_{\theta}(z|x)}{q_{\phi}(z|y, x)} q_{\phi}(z|y, x) dz$$

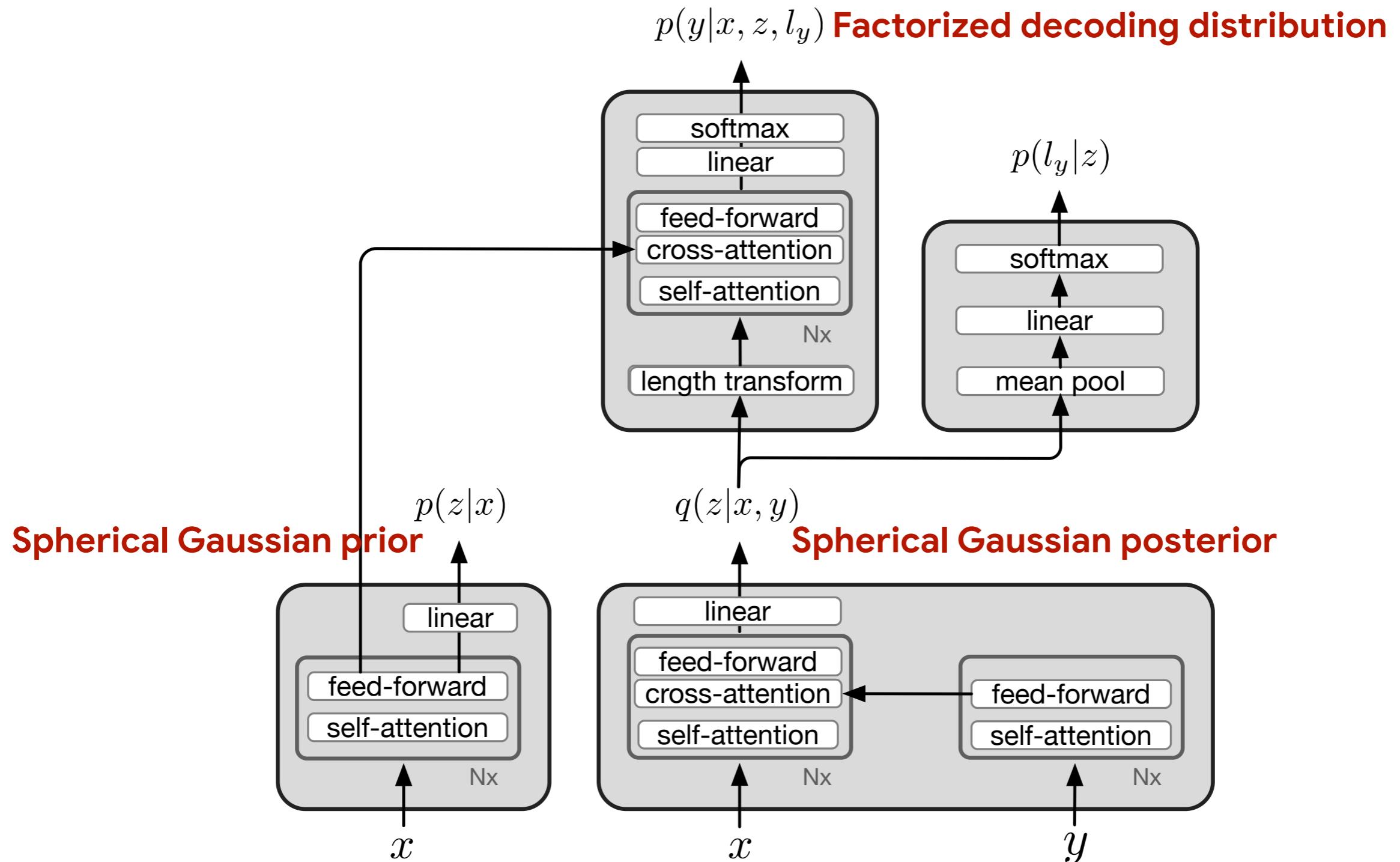
$$= \log \mathbb{E}_{z \sim q_{\phi}} \left[\frac{p_{\theta}(y|z, x) p_{\theta}(z|x)}{q_{\phi}(z|y, x)} \right]$$

$$\geq \mathbb{E}_{z \sim q_{\phi}} [\log p_{\theta}(y|z, x) + \log p_{\theta}(z|x) - \log q_{\phi}(z|y, x)]$$

$$= \mathbb{E}_{z \sim q_{\phi}} [\log p_{\theta}(y|z, x)] - \text{KL}[q_{\phi}(z|y, x) || p_{\theta}(z|x)]$$

$$= \text{ELBO}(x, y, \theta, \phi) \quad J(x, y, \theta, \phi) = \text{ELBO}(x, y, \theta, \phi) + \log p(l_y|z)$$

Network Architecture



Inference with continuous LVM

Given x , inference is the search problem:

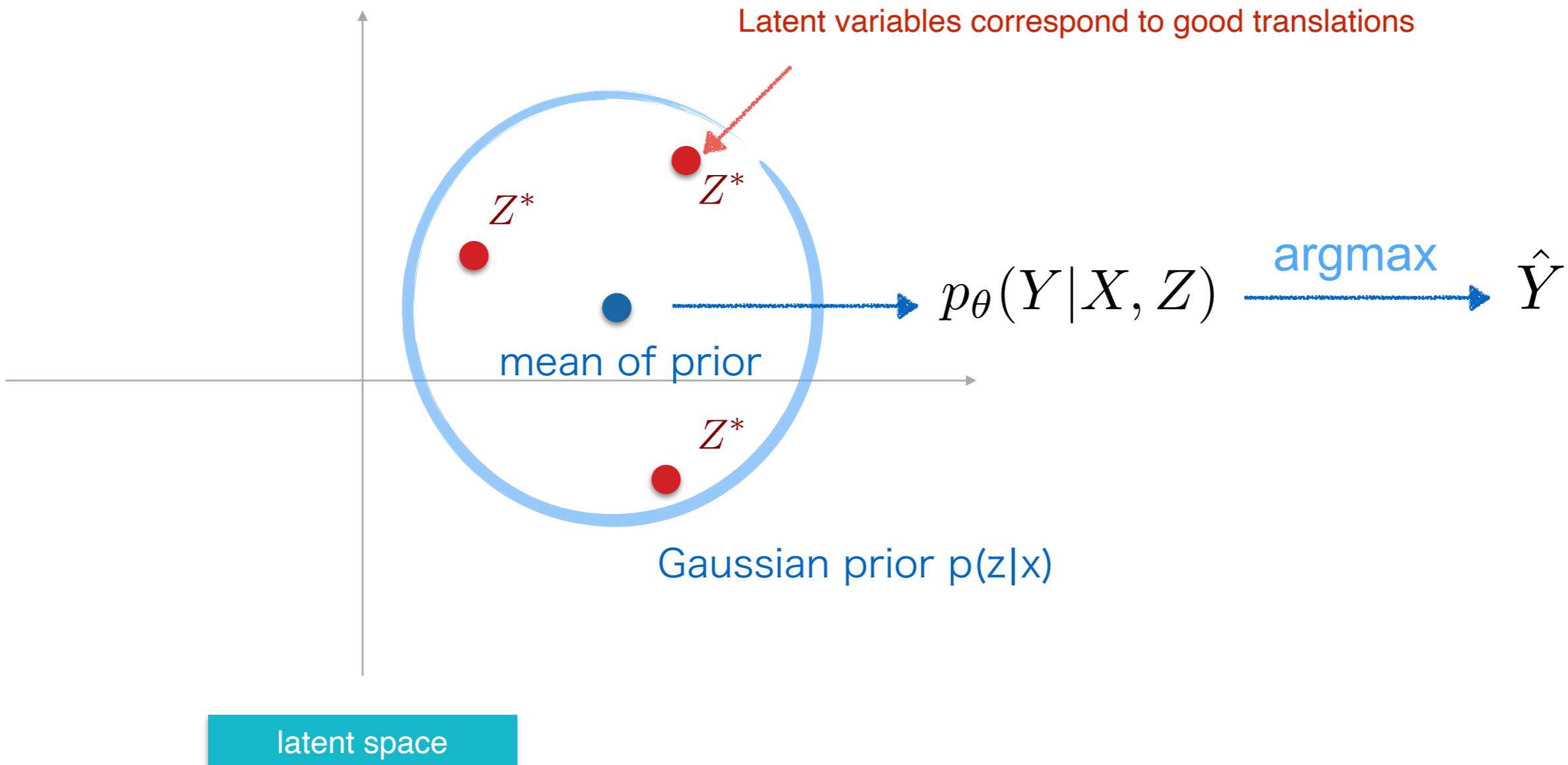
$$\begin{aligned}\hat{y} &= \operatorname{argmax}_y \int_z p(y, z | x) dz \\ &= \operatorname{argmax}_y \int_z p(y | z, x) \boxed{p(z | x)} dz\end{aligned}$$

Spherical Gaussian Prior

$\neq \operatorname{argmax}_y p(y | \mu, x)$, where μ is the mean of the prior

**Simply decoding from the mean of the Gaussian prior
may not lead to the most likely output sequence y**

Inference with continuous LVM



Iterative Inference with a Delta Posterior

Given x , inference is the search problem:

$$\hat{y} = \operatorname{argmax}_y \text{ELBO}(x, y, \theta, \phi)$$

$$= \operatorname{argmax}_y \mathbb{E}_{z \sim q_\phi} [\log p_\theta(y | z, x) + \log p_\theta(z | x) - \log q_\phi(z | y, x)]$$

The expectation over z makes solving this intractable.

Given x , inference is the search problem:

$$\hat{y} = \operatorname{argmax}_y \text{ELBO}(x, y, \theta, \phi)$$

By introducing a Dirac delta posterior:

$$\delta(z | \mu) = \begin{cases} 1 & \text{if } z = \mu \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned} \text{ELBO}(x, y, \theta, \phi) &= \mathbb{E}_{z \sim q_\phi} [\log p_\theta(y | z, x) + \log p_\theta(z | x) - \boxed{\log q_\phi(z | y, x)}] \\ &= \boxed{\log p_\theta(y | \mu, x) + \log p_\theta(\mu | x)} \end{aligned}$$

We perform approximate maximization of this **proxy lowerbound**, using a iterative procedure reminiscent to the EM algorithm.

Iterative Inference with a Delta Posterior

We can approximately maximize $\log p_\theta(y | \mu, x) + \log p_\theta(\mu | x)$ with an EM-like procedure : we call this **Delta Inference**

1. Initialization : $z_0 := \mathbb{E}_{z \sim p_\theta(z|x)}[z]$ (mean of the prior)

2. First output : $\hat{y}_0 = \operatorname{argmax}_y p(y | z_0, x)$

3. Repeat until done:

1. E-Step: match (proxy) delta posterior to the approximate posterior.

$$\hat{\mu}_i = \operatorname{argmin}_\mu \text{KL}[\delta(z | \mu) \| q_\phi(z | \hat{y}_{i-1}, x)] = \mathbb{E}_{z \sim q_\phi(z | y, x)}[z]$$

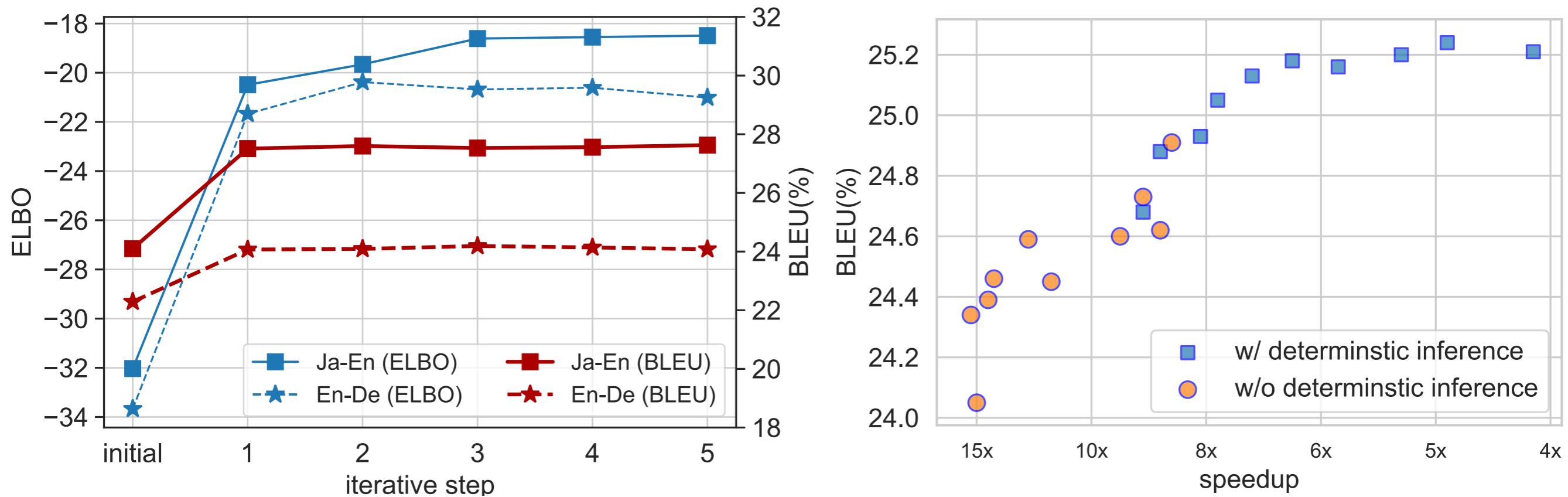
Continuous optimization

2. M-Step: maximize the lowerbound w.r.t y :

$$\hat{y}_i = \operatorname{argmax}_y p(y | \hat{\mu}_i, x)$$

Discrete search

Delta inference gives higher BLEU and ELBO



- Both ELBO and BLEU improve with more refinement steps
- Quality vs. Speed tradeoff.
- Iterative inference improves translation quality in all settings.

Non-autoregressive NMT on WMT'14 En-De

		BLEU(%)	Gap	Speed
(Gu et al., 2018)	Baseline	23.4		
	Gu's Model (NAT + NPD S=100)	19.1	-4.3	2.3x
(Lee et al., 2018)	Baseline	24.5		
	Lee's Model (Adaptive NAR)	21.5	-3.0	1.9x
(Kaiser et al., 2018) (Roy et al., 2018)	Kaiser's baseline	23.5		
	Kaiser's Model (LT, Improved semhash)	19.8	-3.7	3.8x
	Roy's baseline	28.1		
(Ghaz et al., 2019)	Roy's Model (VQVAE + EM + distillation)	26.7	-1.4	4.1x
	Baseline	27.8		
	Ghaz's Model (CMLM with 4 iterations)	26.0	-1.8	
(Shu et al., 2020)	Ghaz's Model (CMLM with 10 iterations)	26.9	-0.9	2~3x
	Baseline	28.3		
	LANMT + Delta Inference	26.1	-2.2	6.3x

Latent Variable Models and Iterative Refinement

(x, y) : (input, output) sequence

z : latent variable capturing the output sequence

$$p(y|x) = \sum_z p(y, z|x) \text{ for discrete latent variables } z$$

1. Iterative Refinement in a Discrete Space (Lee, Mansimov and Cho, 2018).

$$p(y|x) = \int p(y, z|x) dz \text{ for continuous latent variables } z$$

2. Iterative Refinement in a Hybrid Space (Shu, Lee, Nakayama and Cho, 2020).

3. Iterative Refinement in a Continuous Space (Lee, Shu and Cho 2020).

Latent Variable Models and Iterative Refinement

(x, y) : (input, output) sequence

z : latent variable capturing the output sequence

$$p(y|x) = \sum_z p(y, z|x) \text{ for discrete latent variables } z$$

1. Iterative Refinement in a Discrete Space (Lee, Mansimov and Cho, 2018).

$$p(y|x) = \int p(y, z|x) dz \text{ for continuous latent variables } z$$

2. Iterative Refinement in a Hybrid Space (Shu, Lee, Nakayama and Cho, 2020).

3. Iterative Refinement in a Continuous Space (Lee, Shu and Cho 2020).

Motivation

- Delta inference is **unsatisfactory** as the M-step requires computing the softmax (slow operation).

2. M-Step: update our parameters to maximize the lowerbound:

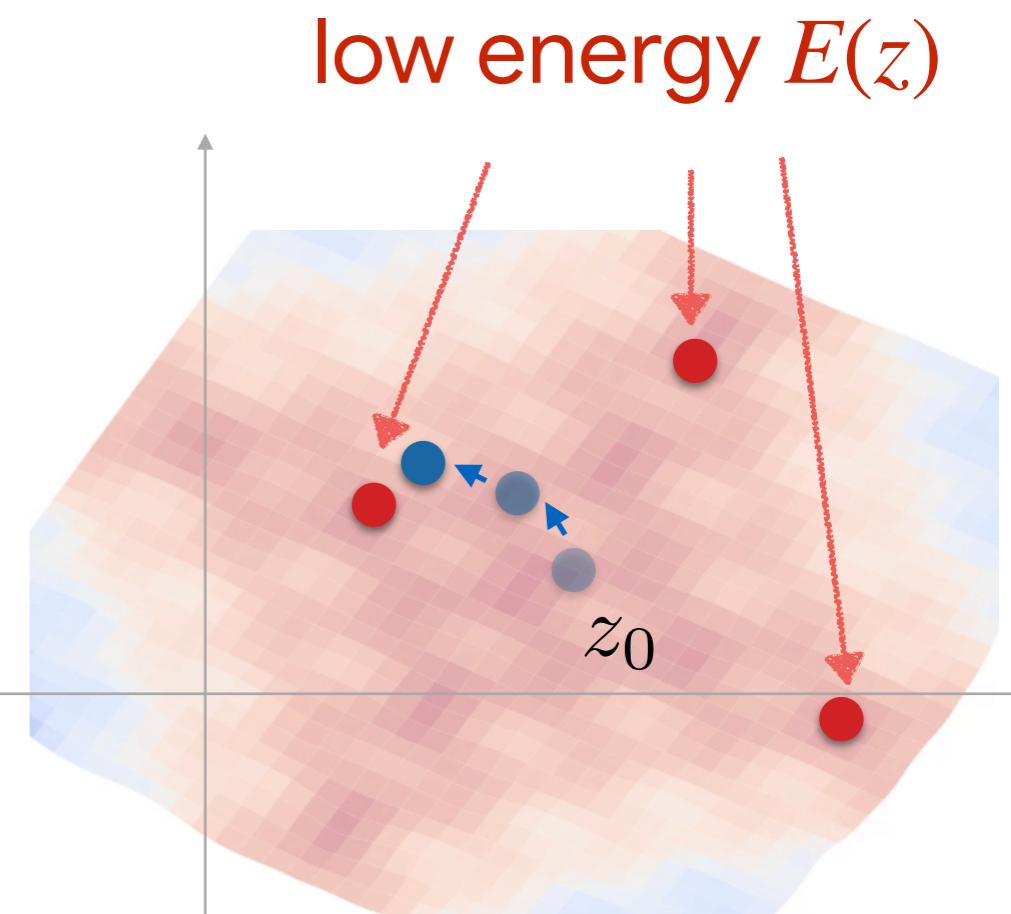
$$\hat{y}_i = \operatorname{argmax}_y p(y | \hat{\mu}_i, x)$$

Discrete search

- If iterative refinement was performed in a continuous space only, we'd see big gains in speed.

Faster Inference with Energy/Score Model

- We can imagine an energy-based model that assigns low energies to local optimas of log-likelihood
- If we can train such an energy-based model $E(z)$, then we can cast a discrete combinatorial search problem into a **continuous optimization** problem.



$$z_{t+1} = z_t - \alpha \cdot \nabla_z E(z)$$

Objective for Inference Network

- Given a latent variable model $p_\theta(y, z | x)$ (trained, fixed):
- Let $\tau_\theta(z; x)$ be a function that approximates the quality of any latent variable z w.r.t the source sentence x .
 - $\tau_\theta(z; x) := \log p_\theta(\hat{y} | x), \quad \hat{y} = \operatorname{argmax}_y p_\theta(y | z, x)$
- Our goal : Find a function $E_\psi(z; x)$ where

$$\operatorname{argmin}_z E_\psi(z; x) \approx \operatorname{argmax}_z (\tau_\theta(z; x))$$

Learning

$$\tau_\theta(z; x) := \log p_\theta(\hat{y} | x), \quad \hat{y} = \operatorname{argmax}_y p_\theta(y | z, x)$$

$$\operatorname{argmin}_z E_\psi(z; x) \approx \operatorname{argmax}_z (\tau_\theta(z; x))$$

Instead of directly approximating τ_θ , we instead train $-E_\psi$ to learn the difference of τ_θ between a pair of latent variable configurations (easier task, this is all we need for inference).

We solve the following for a pair of latent variables $z \neq \bar{z}$:

$$\begin{aligned} & \min_\psi \left\| (-E_\psi(\bar{z}) + E_\psi(z)) - (\tau(\bar{z}) - \tau(z)) \right\|^2 \\ & \approx \min_\psi \left\| \nabla_z E_\psi(z) \right\|^2 + 2 \left((\nabla_z E_\psi(z))^\top \cdot \nabla_z \tau(z) \right) \end{aligned}$$

1st order
Taylor expansion

Learning

$$\tau_\theta(z; x) := \log p_\theta(\hat{y} | x), \quad \hat{y} = \operatorname{argmax}_y p_\theta(y | z, x)$$

$$\operatorname{argmin}_z E_\psi(z; x) \approx \operatorname{argmax}_z (\tau_\theta(z; x))$$

Instead of directly approximating τ_θ , we instead train $-E_\psi$ to learn the difference of τ_θ between a pair of latent variable configurations (easier task, this is all we need for inference).

We solve the following for a pair of latent variables $z \neq \bar{z}$:

$$\begin{aligned} & \min_\psi \left\| (-E_\psi(\bar{z}) + E_\psi(z)) - (\tau(\bar{z}) - \tau(z)) \right\|^2 \\ & \approx \min_\psi \left\| \nabla_z E_\psi(z) \right\|^2 + 2 \left((\nabla_z E_\psi(z))^\top \cdot \boxed{\nabla_z \tau(z)} \right) \end{aligned}$$

1st order
Taylor expansion

Learning

$$\tau_\theta(z; x) := \log p_\theta(\hat{y} | x), \quad \hat{y} = \operatorname{argmax}_y p_\theta(y | z, x)$$

$$\operatorname{argmin}_z E_\psi(z; x) \approx \operatorname{argmax}_z (\tau_\theta(z; x))$$

As $\tau_\theta(z; x)$ is not differentiable with respect to z due to the **argmax** operation in \hat{y} , $\nabla_z \tau_\theta(z; x)$ is not defined.

We therefore use proxy gradient from delta inference, and our final training objective is the following:

$$\mathbb{E}_{z \sim p_\theta(z|x)} \left[\left\| \nabla_z E_\psi(z; x) \right\|^2 + 2 \left((\nabla_z E_\psi(z; x))^\top \cdot (\tilde{z} - z) \right) \right]$$

where \tilde{z} is the output of applying k steps of delta inference on z

Parameterization

$$\mathbb{E}_{z \sim p_\theta(z|x)} \left[\left\| \nabla_z E_\psi(z; x) \right\|^2 + 2 \left((\nabla_z E_\psi(z; x))^\top \cdot (\tilde{z} - z) \right) \right]$$

We either parameterize $\nabla_z E_\psi(z; x)$ as:

1. Gradient of a scalar-valued function E (**Energy** function; LeCun et al., 2006).
2. Train a model $S_\psi(z; x)$ to directly approximate the gradients (**Score** function; Hyvärinen, 2005), bypassing energy estimation.

Gradient-based Inference

Energy-based Inference

1. Initialize: $z = \mathbb{E}_{p_\theta(z|x)}[z]$

2. Repeat:

Backprop

$$z = z - \alpha \cdot \boxed{\nabla_z} E_\psi(z; x)$$

3. Output:

$$\hat{y} = \operatorname{argmax}_y p(y | z, x)$$

Softmax only computed once at the end.
Not on every iteration

Score-based Inference

1. Initialize: $z = \mathbb{E}_{p_\theta(z|x)}[z]$

2. Repeat:

$$z = z - \alpha \cdot S_\psi(z; x)$$

3. Output:

$$\hat{y} = \operatorname{argmax}_y p(y | z, x)$$

Results

	WMT'14 En-De		WMT'16 Ro → En		IWSLT'16 De → En	
	BLEU	Speed	BLEU	Speed	BLEU	Speed
Transformer baseline, beam = 3	28.3	1x	31.5	1x	31.5	1x
Transformer baseline, beam = 1	27.5	1.1x	30.9	1.1x	31.1	1.1x
Latent-variable non-autoregressive	25.7	15x	28.4	34x	27.0	19x
+ Delta Inference	26.1	6.3x	29.0	19x	28.3	11x
+ Energy Inference	26.1	5.8x	28.8	17x	28.6	9.5x
+ Score Inference	26.3	10x	29.1	24x	28.8	13x
+ Score Inference + Latent Search	27.4	6.2x	30.4	15x	30.2	6.3x

- Score inference gives slightly better results than energy inference
 - While being much faster, as it avoids backprop.

Non-autoregressive NMT on WMT'14 En-De

		BLEU(%)	Gap	Speed
(Gu et al., 2018)	Baseline	23.4		
	Gu's Model (NAT + NPD S=100)	19.1	-4.3	2.3x
(Lee et al., 2018)	Baseline	24.5		
	Lee's Model (Adaptive NAR)	21.5	-3.0	1.9x
(Kaiser et al., 2018) (Roy et al., 2018)	Kaiser's baseline	23.5		
	Kaiser's Model (LT, Improved semhash)	19.8	-3.7	3.8x
	Roy's baseline	28.1		
	Roy's Model (VQVAE + EM + distillation)	26.7	-1.4	4.1x
(Ghaz et al., 2019)	Baseline	27.8		
	Ghaz's Model (CMLM with 4 iterations)	26.0	-1.8	
	Ghaz's Model (CMLM with 10 iterations)	26.9	-0.9	2~3x
(Shu et al., 2020) (Lee et al., 2020)	Baseline	28.3		
	LANMT + Delta Inference	26.1	-2.2	6.3x
	LANMT + Score Inference + Latent Search	27.4	-0.9	6.2x

Gradient and latent space visualization

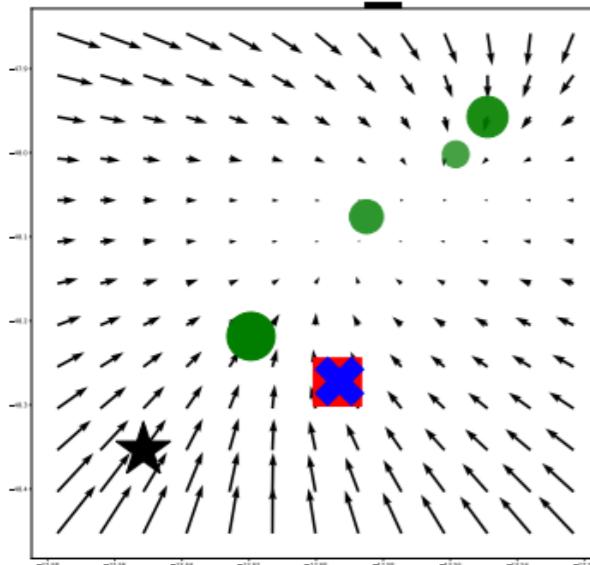
Token 1

Post: _So

Prior: _So

Delta: _So

Score: _So



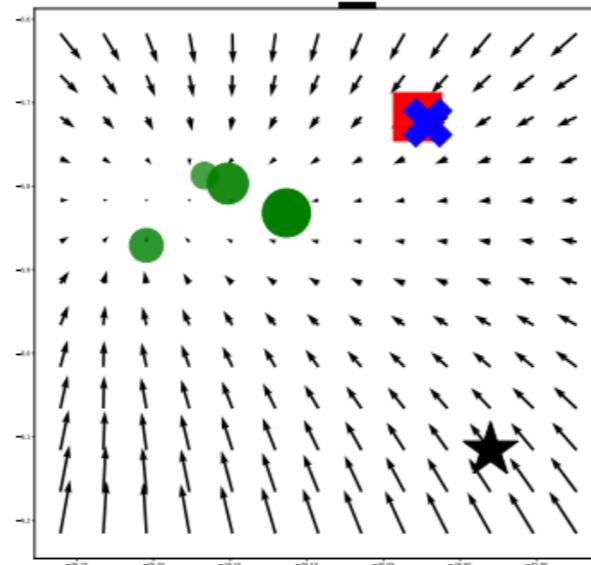
Token 2

Post: _what

Prior: _what

Delta: _what

Score: _what



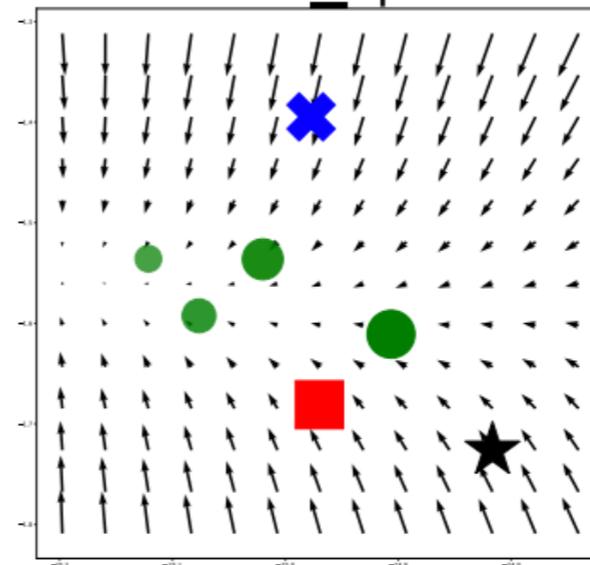
Token 3

Post: _opened

Prior: _opened

Delta: _opened

Score: _opened



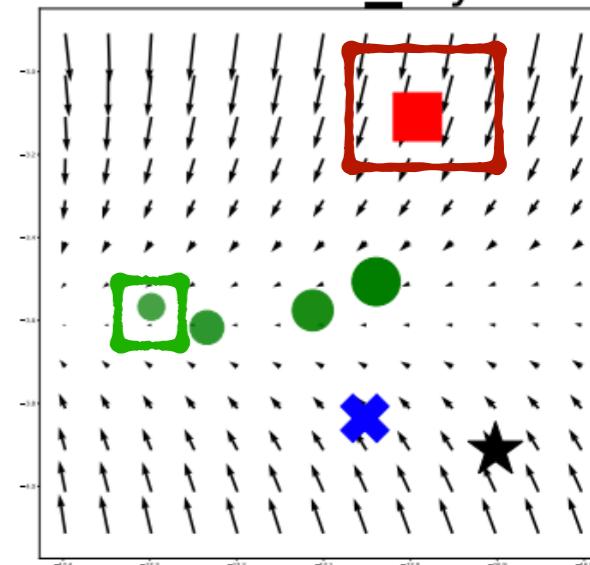
Token 4

Post: _my

Prior: _me

Delta: _me

Score: _my



★ : prior mean

● : score inference

✖ : delta inference

■ : posterior mean

Source Reference	Was öffnete mir also die Augen? So what opened my eyes ?
Posterior	So what opened my eyes ?
Prior	So what opened me eyes ?
Delta	So what opened me eyes ?
Score	So what opened my eyes ?

Non-trivial, non-local changes

Example 1

Source	There aren 't many doctors in the west African country ; just one for every 5,000 people
Reference	In dem westafrikanischen Land gibt es nicht viele rzte, nur einen fr 5.000 Menschen
Original	Es gibt nicht viele rzte im westafrikanischen Land, nur eine fr 5.000 Menschen.
Refined	Im westafrikanischen Land gibt es nicht viele rzte, nur eine fr 5.000 Menschen.

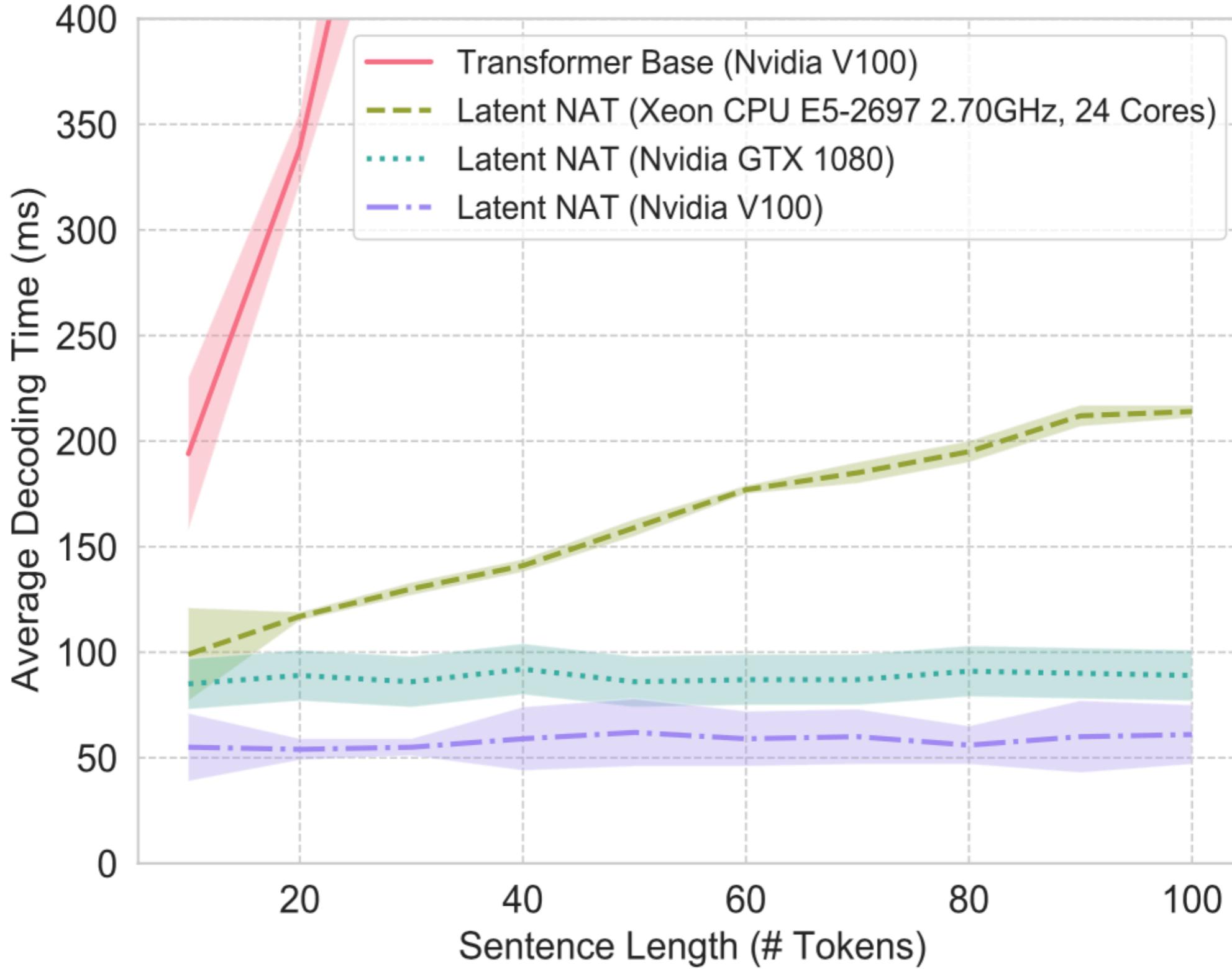
Example 2

Source	Costumes are expected to account for \$ 1.2 billion dollars out of the \$ 6.9 billion spent , according to the NRF .
Reference	Die Kostme werden etwa 1,2 Milliarden der 6,9 Milliarden ausgegebenen US-Dollar ausmachen, so der NRF.
Original	Es wird von, Kostme, dass sie die dem NRF ausgegebenen 6,9 Milliarden Dollar 1,2 Milliarden Dollar ausmachen.
Refined	Es wird erwartet, dass die Kostme nach Angaben des NRF 1,2 Milliarden Dollar aus den 6,9 Milliarden Dollar ausmachen.

Example 3

Source	It was with this piece of Bedouin wisdom that the first ever chairman Wolfgang Henne described the history and fascination behind the “Helping Hands” society .
Reference	Mit dieser Beduinenweisheit beschrieb der erste Vorsitzende Wolfgang Henne die Geschichte und Faszination des Vereins “Helfende Hnde”.
Original	Der erste Vorsitzende Wolfgang Henne beschrieb mit dieser erste Weisheit in Bedouin” die Geschichte und Faszination hinter der “Helenden Hands” Gesellschaft
Refined	Mit diesem Stck Bedouin-Weisheit beschrieb der erste Vorsitzende Wolfgang Henne jemals die Geschichte und Faszination hinter der “Heling Hands” Gesellschaft

Inference Latency (lower = better) w.r.t Length and Computational Capacity



Is speed the only reason for going non-auto?

- Autoregressive models suffer from:
 - Exposure bias (due to teacher forcing)
 - Pathological degenerate generations (infinite repetitions)

```
: response = openai.Completion.create(engine="davinci", prompt="roses are red, elephants are gray,  
print(response["choices"][0]["text"])
```

gray, and so on.

The color of an object is determined by the wavelength of light that it reflects. The color of a
rbs.

The color of an object is not determined by the wavelength of light that it emits.

The color of an object is not determined by the wavelength of light that it transmits.

The color of an object is determined by the wavelength of light that it reflects.

Collaborators



Kyunghyun



Raphael



Elman



Hideki

Thank you. Questions?

jasonleeinf.github.io