# Non-autoregressive Neural Machine Translation Based on Latent-Variable Models
## 潜在変数モデルに基づいた非自己回帰型ニューラル機械翻訳

## Raphael Shu (朱 中元)
## 東京大学

# Neural Machine Translation

- NMT is conditional probabilisitic model

| 1 | ＤＥＲＳソフトウエアを用いて「ふげん発電 | | details of dose rate of `` Fugen Power Plant |
| 2 | 線量率を計算する際の状況の変化，すなわち | | the changes of conditions for computation of |
| 3 | ＤＥＲＳはこれらの変化に対応して新たな線 | | responding to these changes DERs can compute |
| 4 | このソフトウエアのＲ５バージョンの特徴， | | the characteristics of R5 version of this so: |
| 5 | 感知用と出力用の２基のコイル，増幅器，及 | | here was developed a phase shift magnetic se: |
| 6 | この回路は，入力信号位相の変化により共振 | | this is a feedback circuit shifting resonanc( |
| 7 | コイルとしては，内径６ｍｍで８００ターン | | for the coils , here were used two coils witl |
| 8 | 実験では，水道水，純水，及び磁化デバイス | | on a test , it was possible to analyze featu: |
| 9 | また，流速と流量の変化も検出できることを | | and , it was confirmed to enable also to det( |
| 10 | 無線ＩＣタグ（ＲＦＩＤ）及びセンサーネッ | | here was described a high sensitivity strain |

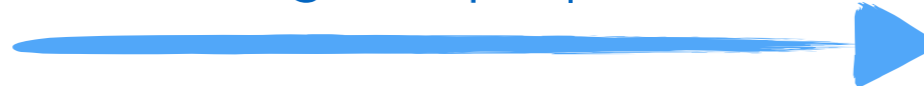Parallel corpus

training $\rightarrow$

**MT Model**
$$p(Y|X)$$

- Translation in NMT
  - Find an output to maximize the probability

そこのどんつきまでガッと行ったら右やで

$X$

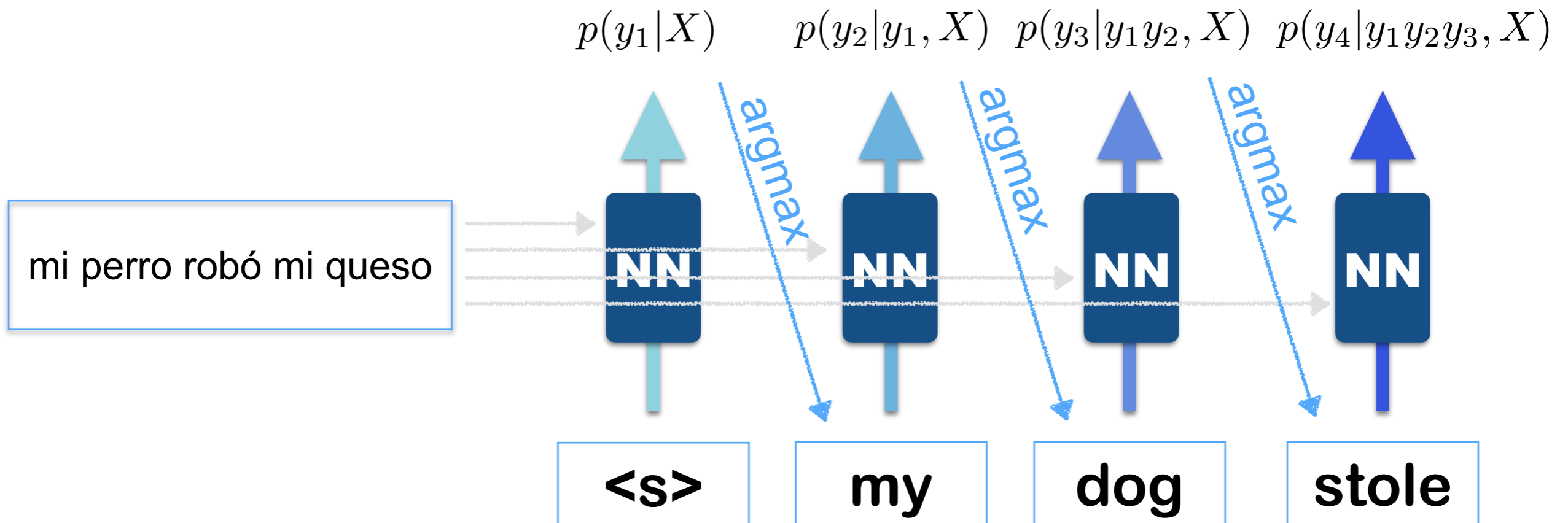what is argmax p(Y|X=そこの...) ?

**Go straight and turn right**

$Y$

# Autoregressive modeling

- Most of current application-level NMT models are based on autoregressive modeling
    - Breaks down p(Y|X) to word probabilities using chain rule
    - In each step, the model predicts the next word

```
p(my dog stole my cheese | mi perro robó mi queso)
      = p(my      |                    mi perro robó mi queso)
        p(dog     | my,                mi perro robó mi queso)
        p(stole   | my dog,            mi perro robó mi queso)
        p(my      | my dog stole,      mi perro robó mi queso)
        p(cheese  | my dog stole my,   mi perro robó mi queso)
```

# Obtain Translations

- Approximating the global argmax with search algorithms
    - Greedy search and beam search

$$p(y_1|X) \qquad p(y_2|y_1, X) \quad p(y_3|y_1y_2, X) \quad p(y_4|y_1y_2y_3, X)$$
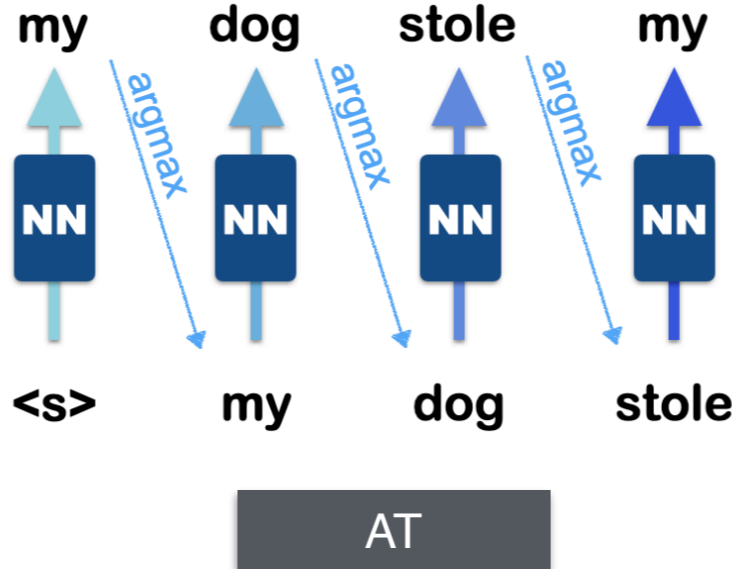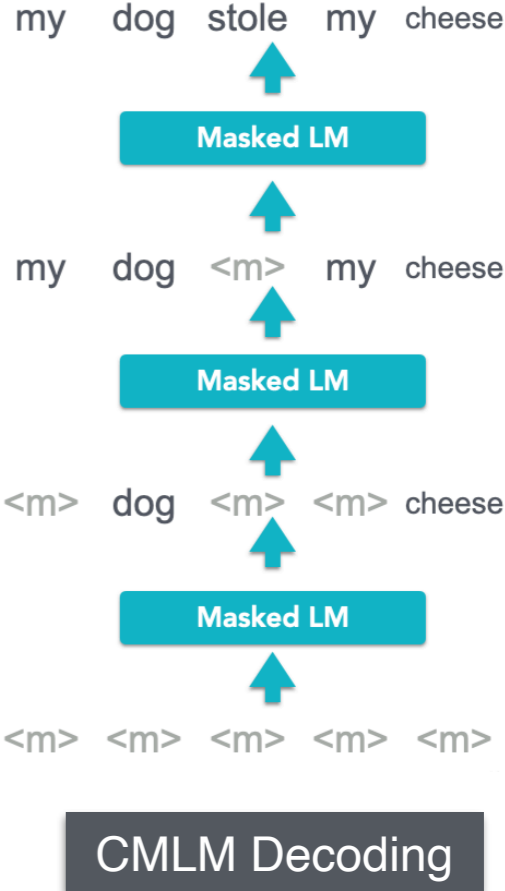


**Problems**

- <u>Low parallelizability</u>
- Worse if the model is bigger and deeper
- Require search algorithm to approximate the argmax

# Non-autoregressive Machine Translation

- Predict all output tokens in one forward pass
- Can be fully parallelized on GPU
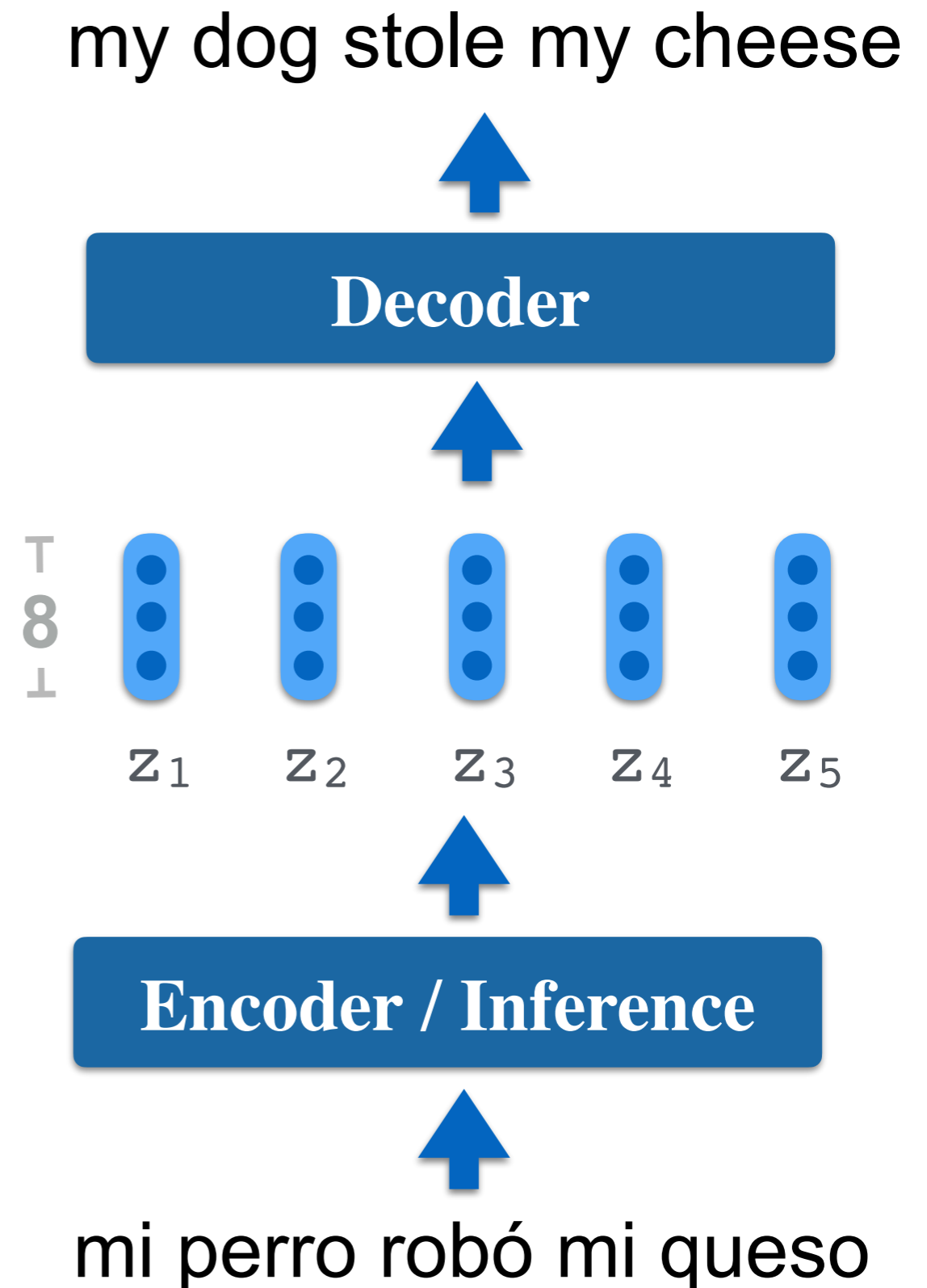


AT



NAT

- CMLM (conditional masked language model) [Ghazvinine et al., 2018]
  - Predict the sequence and mask tokens with low confidence
  - Perform such token refinement by multiple iterations
- Drawback of token-based refinement models
  - Token prediction is time-consuming
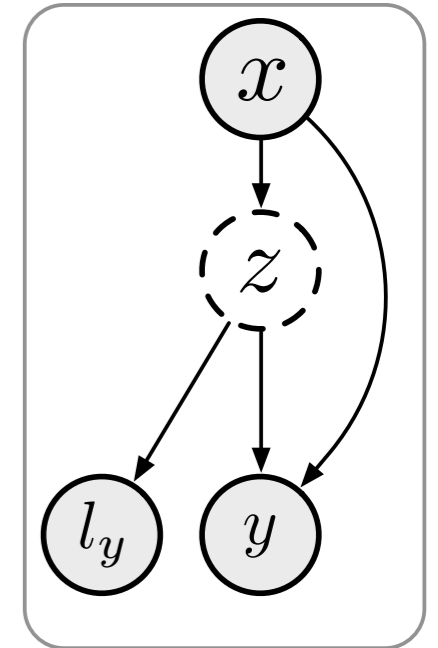


CMLM Decoding

5

# Latent-variable Non-autoregressive MT

- Core Idea:
  - Capture translation decisions with continuous latent variables
- Each source token is assigned with one latent variable
- Each latent variable is a <u>low-dimensional</u> vector
- Finding the best setting of latent variables with high-speed inference

my dog stole my cheese

**Decoder**

$z_1$  $z_2$  $z_3$  $z_4$  $z_5$

**Encoder / Inference**

mi perro robó mi queso

# Objective function



- Similar to VAE, we train our model with ELBO (evidence lower bound)

$$\log p(Y|X) \geq \text{ELBO}(X, Y; \theta, \phi, \omega)$$

$$= \mathbb{E}_{Z \sim q_\phi} \Big[ \log p_\theta(Y|X, Z, l_Y) p_\theta(l_Y|X, Z) \Big] - \text{KL}(q_\phi(Z|X, Y) || p_\omega(Z|X))$$
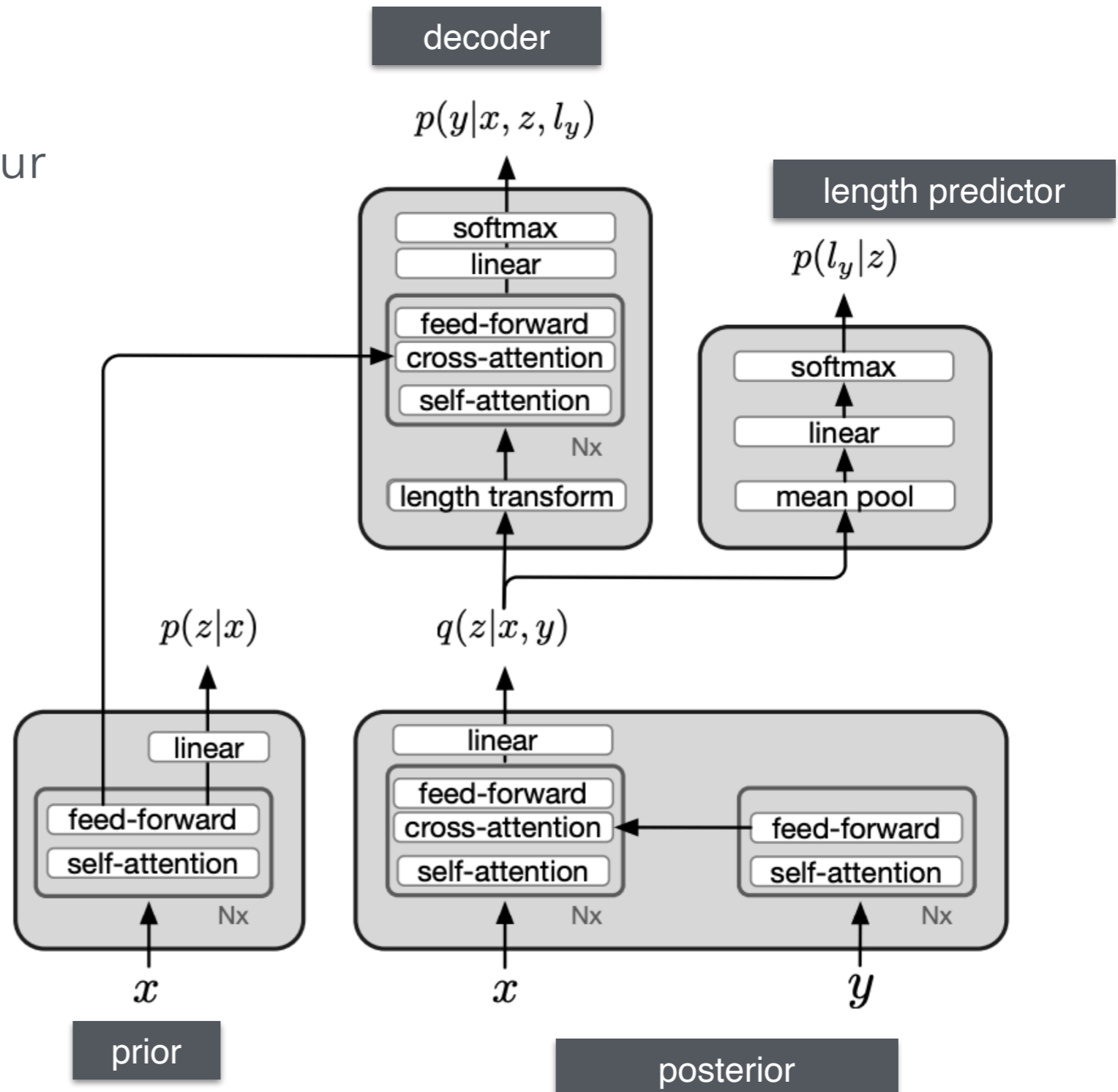
| decoder | length predictor | posterior | prior |

- Training:

$$\hat{\theta}, \hat{\phi}, \hat{\omega} = \underset{\theta, \phi, \omega}{\text{argmax}} \, \text{ELBO}(X, Y; \theta, \phi, \omega)$$

# Model architecture

- Latent NAT parameterizes four distributions
- Reuse Transformer modules
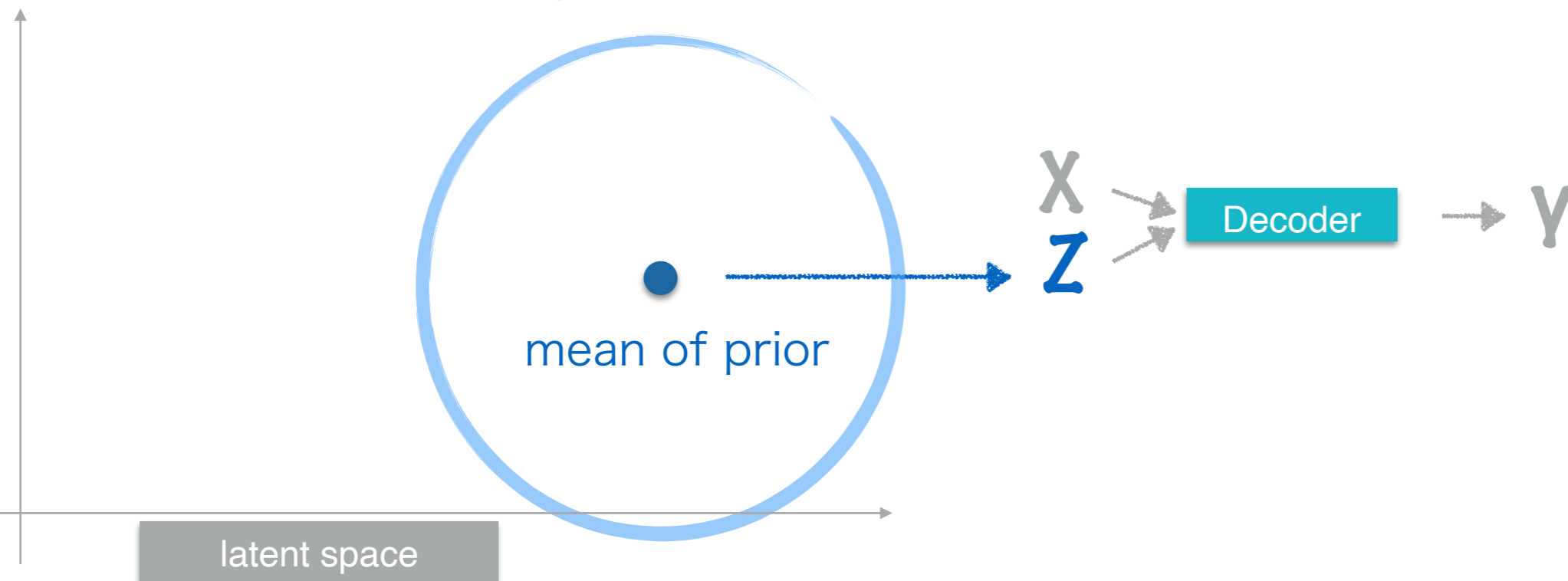- **Length Transform**: adjust |x| vectors to |y| vectors (skip the details here)

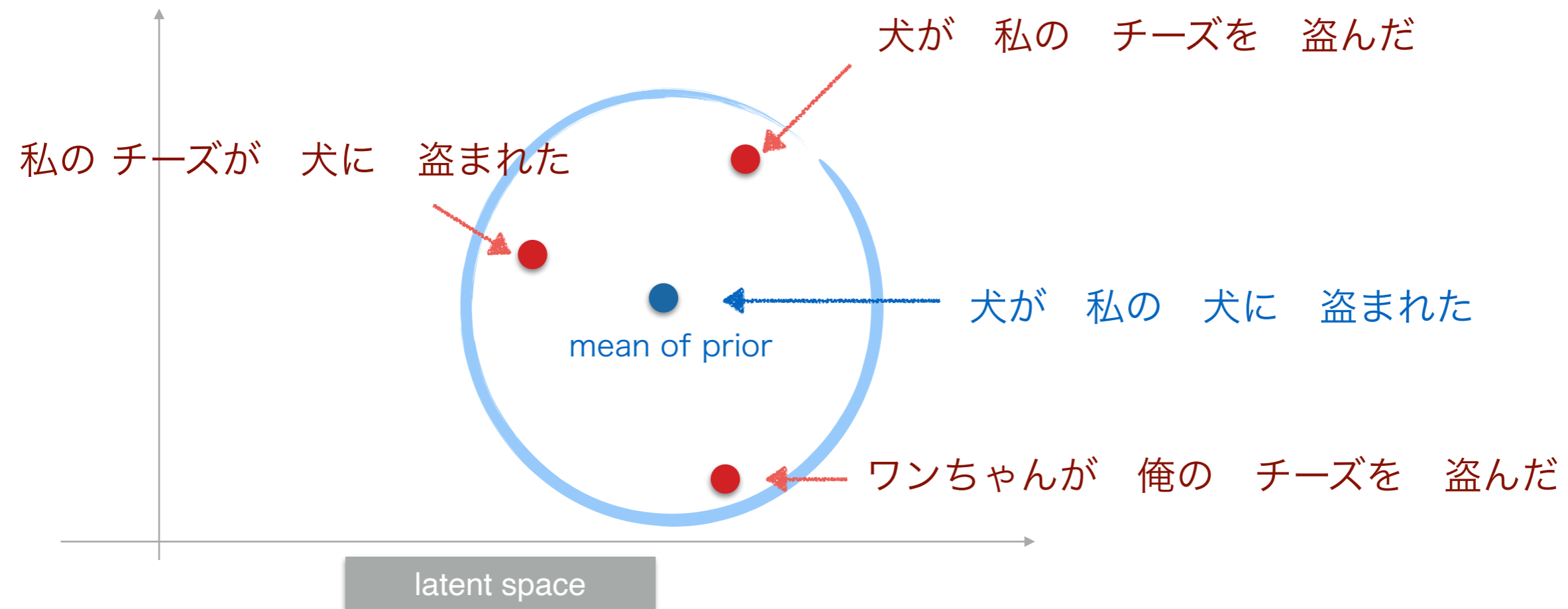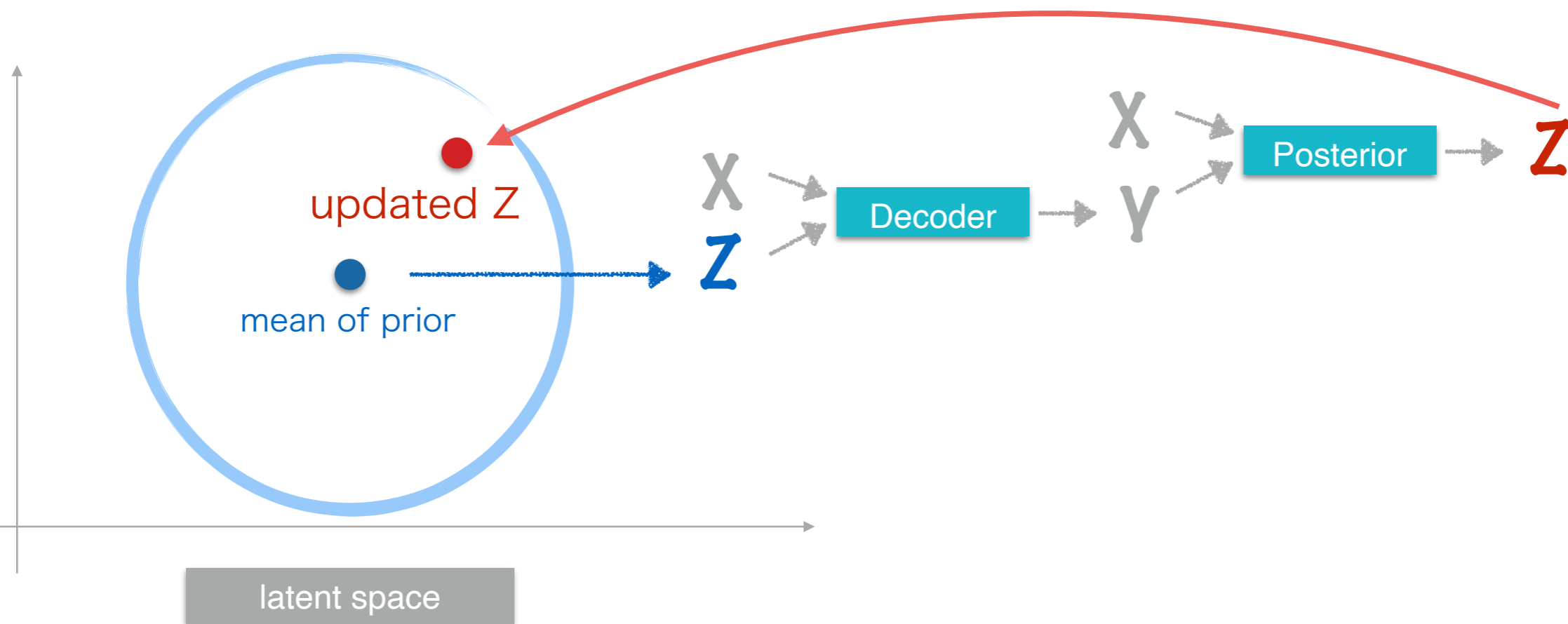# Translation with Latent NAT

- After model training we get:

X ⇢ [Prior] ⇢ Z     X,Z ⇢ [Decoder] ⇢ Y     X,Y ⇢ [Posterior] ⇢ Z

- Naive inference (decoding)

X ⇢ [Prior] ⇢ $p(Z|X)$

mean of prior

X, Z ⇢ [Decoder] ⇢ Y

latent space

# Problem of naive inference

- Problem of naive inference
  - the center of a Gaussian may not produce the best results

犬が　私の　チーズを　盗んだ

私の チーズが　犬に　盗まれた

mean of prior

犬が　私の　犬に　盗まれた
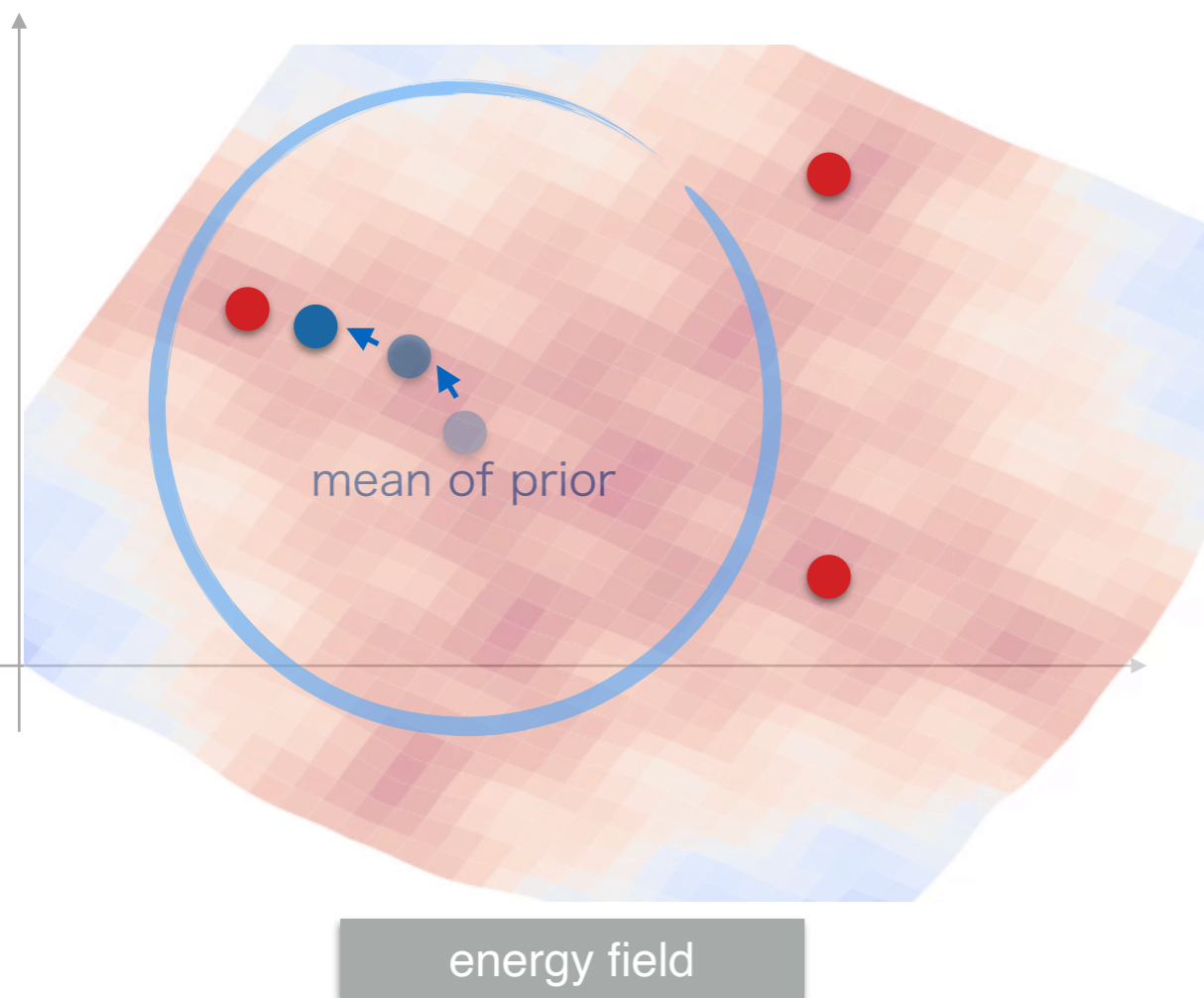
ワンちゃんが　俺の　チーズを　盗んだ

latent space

# Better inference approaches (1)

- Delta inference (Shu et al., AAAI 2020)
  - Latent-variable updating for maximizing approximated ELBO
  - ELBO is improved after iterations with rapid convergence

# Better inference approaches (2)

- Energy-based inference (Jason et al., EMNLP 2020)
  - Build an energy model $E(Z)$
  - High-quality latent vectors get low energy

- Update latent variables with the energy gradient

$$Z_{t+1} = Z_t - \alpha \nabla_{Z_t} E(Z_t)$$

**Energy-driven gradient descent**

mean of prior

energy field

# Better inference approaches (2)

- Energy-based inference (Jason et al., EMNLP 2020)
  - Build an energy model $E(Z)$
  - High-quality latent vectors get low energy

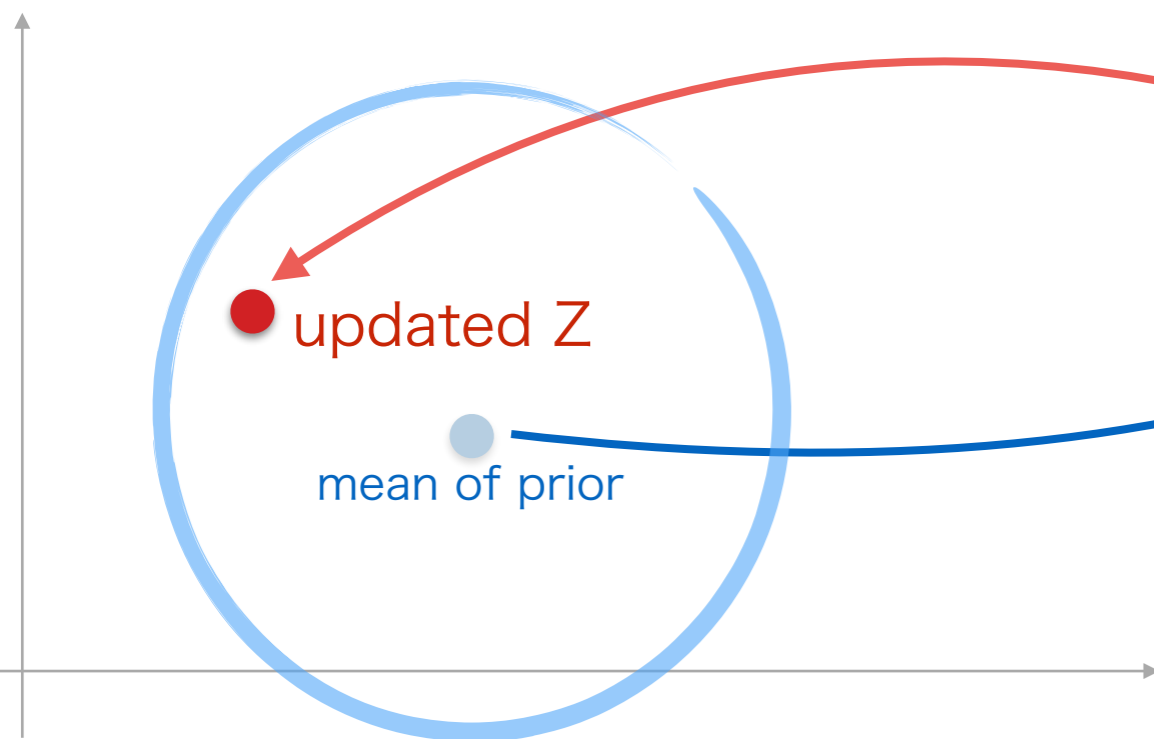- Update latent variables with the energy gradient

$$Z_{t+1} = Z_t - \alpha \nabla_{Z_t} E(Z_t)$$

**Energy-driven gradient descent**

updated Z

mean of prior

energy field

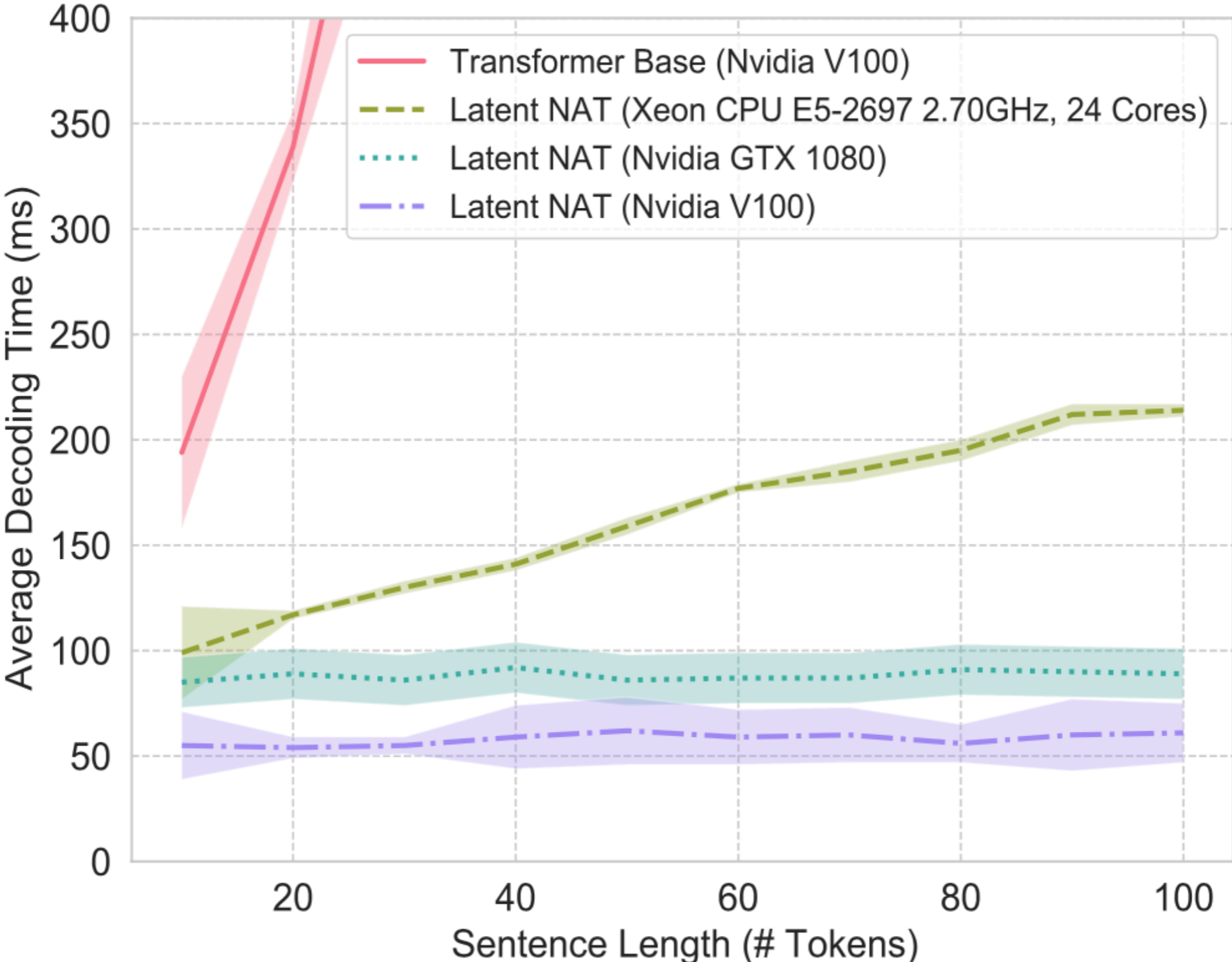# Experiments

# Experiment settings

- Dataset:
  - WMT'14 English -> German
  - IWSLT'16 Romanian -> English
  - IWSLT'16 German -> English
- Evaluation
  - Translation quality: BLEU
  - Translation speed: averaged decoding time for one sentence
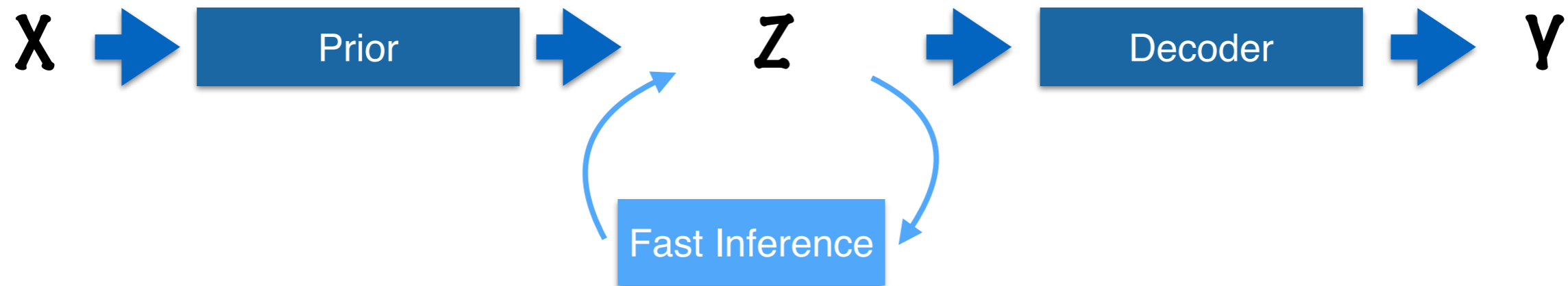
# Experiment results on machine translation

| | 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 NAT (Naive Inference)** | 25.7 | 15x | 28.4 | 34x | 27.0 | 19x |
| **+ Delta Inference** | 26.1 | 6.3x | 29.0 | 19x | 28.3 | 11x |
| **+ Energy Inference** (w/ approximation) | 26.3 | 10x | 29.1 | 24x | 28.8 | 13x |
| **+ Score Inference + Latent Search** | 27.4 | 6.2x | 30.4 | 15x | 30.2 | 6.3x |

- Latent Search: parallel decoding by sampling multiple latent variable

# Translation speed related to computational capacity

# Conclusion

X ➡ [ Prior ] ➡ Z ➡ [ Decoder ] ➡ Y

[ Fast Inference ]

- We show a novel sequence generation framework
  - sequence prediction problem is solved by latent-variable inference in the continuous space
- Continuous setting and low dimensionality enable us to updating efficiently
- Fit for on-device computing

# Thanks