

ENGINE: Energy-Based Inference Networks for Non-Autoregressive Neural Machine Translation

Lifu Tu¹ Richard Yuanzhe Pang² Sam Wiseman¹ Kevin Gimpel¹

¹Toyota Technological Institute at Chicago

²New York University

ACL 2020

Motivation

Autoregressive Neural Machine Translation

- state-of-art models for neural machine translation
- left-to-right decoding

Non-Autoregressive Neural Machine Translation [Gu et al., 2018]

- parallel decoding, much faster
- large performance gap with autoregressive models

Motivation

Autoregressive Neural Machine Translation

- state-of-art models for neural machine translation
- left-to-right decoding

Non-Autoregressive Neural Machine Translation [Gu et al., 2018]

- parallel decoding, much faster
- large performance gap with autoregressive models

In this work, train a non-autoregressive model to minimize autoregressive energy!

Inference for Structured Models

$$\mathbf{y}^* = \arg \min_{\mathbf{y}} E_{\Theta}(\mathbf{x}, \mathbf{y})$$

where energy function $E(\mathbf{x}, \mathbf{y})$ is a scalar that measures the compatibility of each configuration \mathbf{x} and \mathbf{y} [LeCun et al., 2006; Belanger and McCallum, 2016]

Inference for Structured Models

$$\mathbf{y}^* = \arg \min_{\mathbf{y}} E_{\Theta}(\mathbf{x}, \mathbf{y})$$

where energy function $E(\mathbf{x}, \mathbf{y})$ is a scalar that measures the compatibility of each configuration \mathbf{x} and \mathbf{y} [LeCun et al., 2006; Belanger and McCallum, 2016]

Inference Networks [Tu and Gimpel (2018,2019)]

An inference network $\mathbf{A}_{\Psi} : \mathcal{X} \rightarrow \mathcal{Y}_R$ is parameterized by Ψ and trained with the goal that

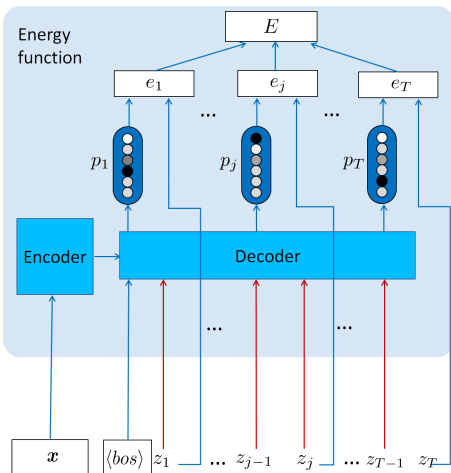
$$\mathbf{A}_{\Psi}(\mathbf{x}) \approx \arg \min_{\mathbf{y} \in \mathcal{Y}_R(\mathbf{x})} E_{\Theta}(\mathbf{x}, \mathbf{y})$$

$\mathcal{Y}_R(\mathbf{x})$ is relaxed **continuous** output space.

The objective for inference networks:

$$\hat{\Psi} = \arg \min_{\Psi} \sum_{\langle \mathbf{x}, \mathbf{y} \rangle \in \mathcal{D}} E_{\Theta}(\mathbf{x}, \mathbf{A}_{\Psi}(\mathbf{x}))$$

ENGINE

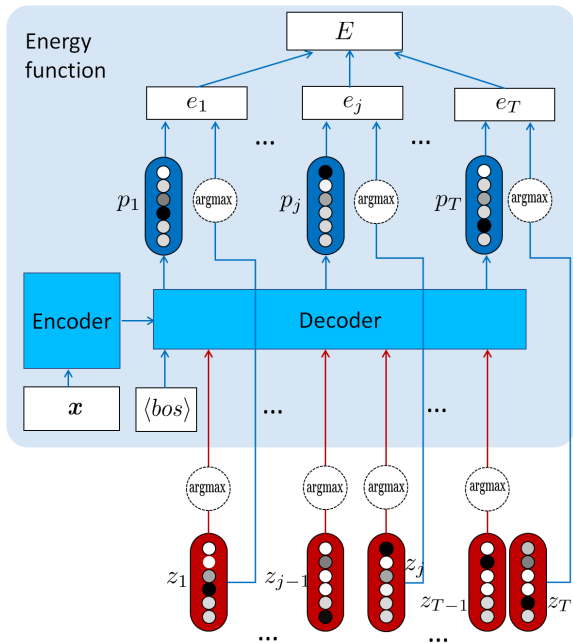


score a sequence of words \mathbf{z}

$$e_j = -\log p_{\Theta}(z_t | \mathbf{z}_{0:j-1}, \mathbf{x})$$

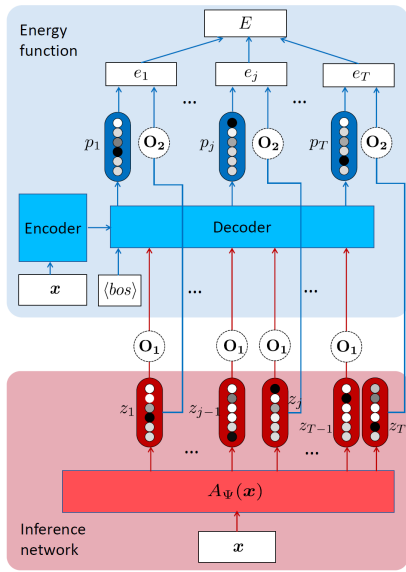
$$E = \sum_{j=1}^T e_j$$

The beam search algorithm is to minimize the above energy.



we introduce the generalized energy E

Here z is a sequence of probability vectors!

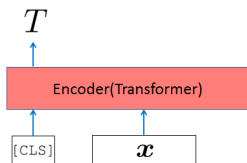


How to allow gradient-based optimization?

argmax operation is replaced by O_1 or O_2

Inference Networks Architecture

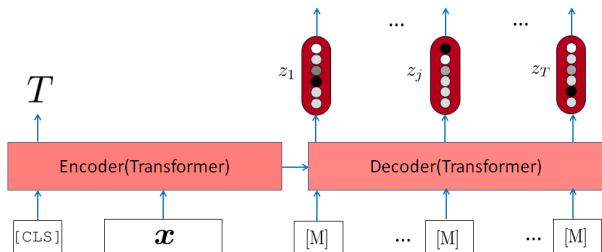
CMLM [Ghazvininejad et al., 2019]



Predicting Length

Inference Networks Architecture

CMLM [Ghazvininejad et al., 2019]

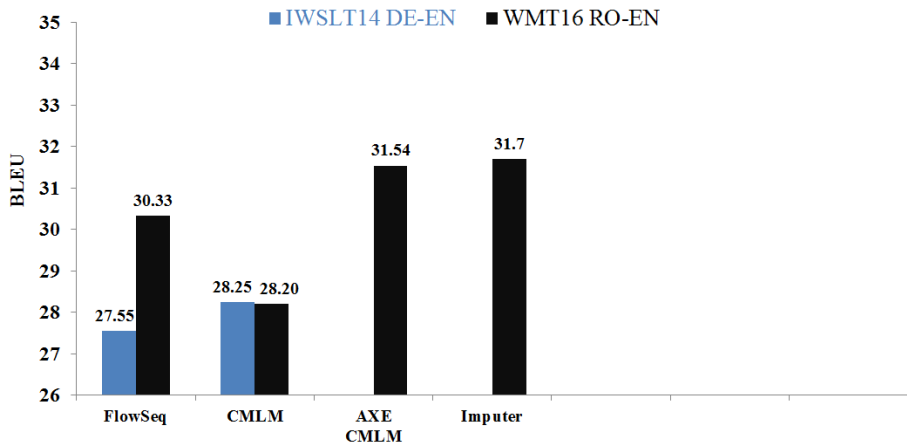


The decoder inputs are the special masked tokens $[M]$.

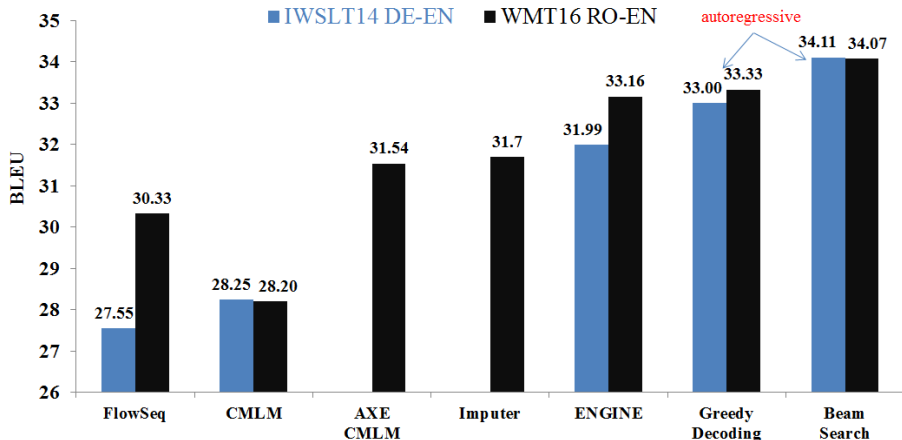
Experimental Setup

- Datasets: IWSLT2014 DE-EN and WMT2016 RO-EN
- Energy function: pretrained autoregressive model
- Inference network architecture: CMLM
- One decoding iteration is used (meaning they are purely non-autoregressive.)

Results



Results



- achieves state-of-the-art results for non-autoregressive translation
- approaching the performance of autoregressive models

Example 1

Ground Truth :

the u.n. chief again urged all parties , including the divided u.n. security council , to unite and support inclusive negotiations to find a political solution .

CMLM :

the un chief **again again** urged all parties , including the divided un security council to unify and support negotiations in order to find a political solution .

ENGINE (ours) :

the un chief has again urged all parties , including the divided un security council to unify and support negotiations in order to find a political solution .

An example from the WMT'16 RO-EN test set.

Example 2

Ground Truth :

the study 's conductors transmit that " romanians feel the need for a little more adventure in their lives (24 %) , followed by affection (21 %) , money (21 %) , safety (20 %) , new things (19 %) ...

CMLM :

survey survey makers say that ' romanians romanians some some-thing adventadventure ure their lives 24 24 %) followed followed by % % % % % , (21 % %) , safety (% % %) , new19% %) , , 19 % % %) , respect 18 % % % % % % % % % , , % % % % % % % % , , % , 14 % , 12 % %

ENGINE (ours) :

realisation of the survey say that ' romanians feel a slightly more adventure in their lives (24 %) followed by aff% (21 %) , money (21 %) , safety (20 %) , new 19 %) , sex (19 %) , respect 18 % , confidence 17 % , 17 % , connecting 17 % , knowledge % % , 14 % , 14 % , 12 % %

An longer example from the WMT'16 RO-EN test set.

Thanks!