

A Teacher-Student Framework for Zero-Resource Neural Machine Translation

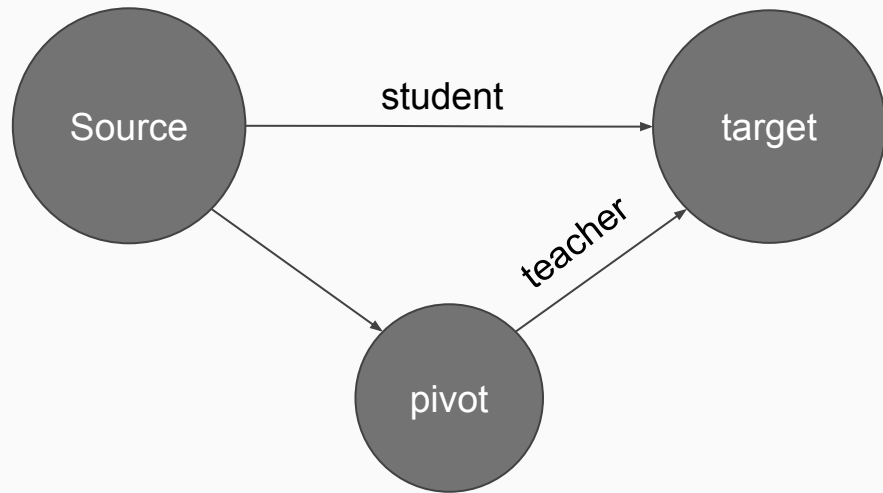
Authors: Yun Chen, Yang Liu, Yong Cheng, Victor O.K. Li

Conference: ACL - 2017

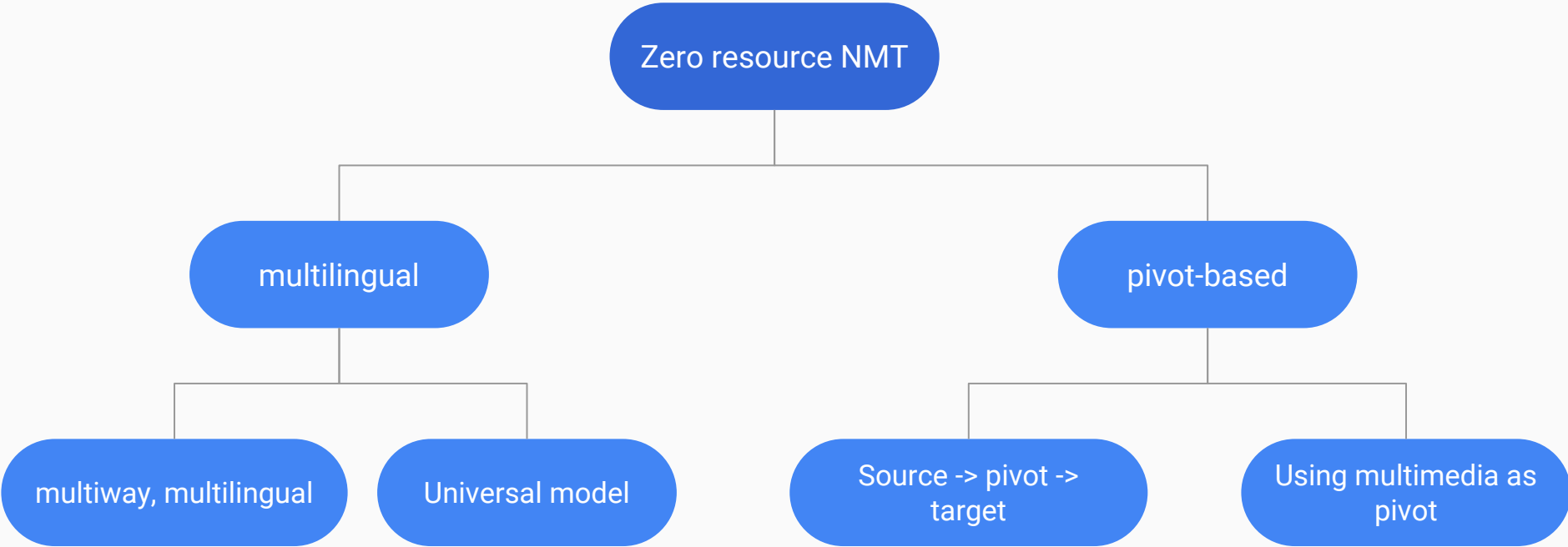
Link: http://nlp.csai.tsinghua.edu.cn/~ly/papers/acl2017_cy.pdf

Abstract

- Despite of recent remarkable progress in end-to-end NMT, the field still struggles to produce good results for languages with less data.
- This paper proposes a method for zero-resource (without parallel corpora for source-target) NMT by assuming that parallel sentences have close probabilities of generating a sentence in a third language.
- Based on this assumption, their method is able to train a source-to-target NMT model (“student”) without parallel corpora available, guided by an existing Pivot-to-target NMT model (“teacher”) on a source-pivot parallel corpus.
- Proposed method significantly improves over a baseline pivot-based model by +3.0 BLEU points across various language pairs.



Introduction



Multilingual

- Firat et al. (2016b) present a multi-way, multilingual model with shared attention to achieve zero-resource translation. They fine-tune the attention part using pseudo bilingual sentences for the zero-resource language pair.
- Another direction is to develop a universal NMT model in multilingual scenarios (Johnson et al., 2016; Ha et al., 2016). Use parallel corpora of multiple languages to train one single model, which is then able to translate a language pair without parallel corpora available
- Although these approaches prove to be effective, the combination of multiple languages in modeling and training leads to increased complexity compared with standard NMT.

Pivot-based

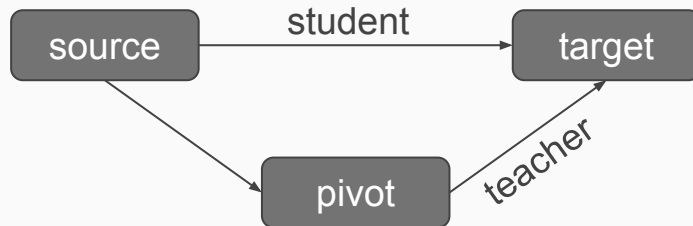
- Achieve source-to-target NMT without parallel data via a *pivot*, which is either text or image.
- Cheng et al. (2016a) propose a pivot-based method for zero-resource NMT:



- Use of multimedia information as pivot also benefits zero-resource translation (Nakayama and Nishida, 2016)
- However, pivot-based approaches usually need to divide the decoding process into two steps, which is not only more computationally expensive, but also potentially suffers from the error propagation problem (Zhu et al., 2013).

Proposed System

- The proposed system assumes that parallel sentences should have close probabilities of generating a sentence in a third language.
- An existing “teacher” model is leveraged to guide the learning of “student” model on a source-pivot parallel corpus.
- Compared with pivot based approaches (Cheng et al., 2016a), this method allows direct parameter estimation of the intended NMT model, without the need to divide decoding into two steps.
- This strategy improves efficiency and avoids error propagation in decoding.
- Experiments on the Europarl and WMT datasets show that this approach achieves significant improvements in terms of both translation quality and decoding efficiency over a baseline pivot-based approach to zero-resource NMT on Spanish-French and German-French translation tasks.



Background

- Let x be a source-language sentence and y be a target-language sentence. We use $P(y|x; \theta_{x \rightarrow y})$ to denote a source-to-target neural translation model, where $\theta_{x \rightarrow y}$ is a set of model parameters.
- Given a source-target parallel corpus $D_{x,y}$, which is a set of parallel source-target sentences, the model parameters can be learned by maximizing the log-likelihood of the parallel corpus:

$$\hat{\theta}_{x \rightarrow y} = \operatorname{argmax}_{\theta_{x \rightarrow y}} \left\{ \sum_{\langle \mathbf{x}, \mathbf{y} \rangle \in D_{x,y}} \log P(\mathbf{y} | \mathbf{x}; \theta_{x \rightarrow y}) \right\}$$

- Given learned model parameters $\hat{\theta}_{x \rightarrow y}$, the decision rule for finding the translation with the highest probability for a source sentence x is given by

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \left\{ P(\mathbf{y} | \mathbf{x}; \hat{\theta}_{x \rightarrow y}) \right\}. \quad (1)$$

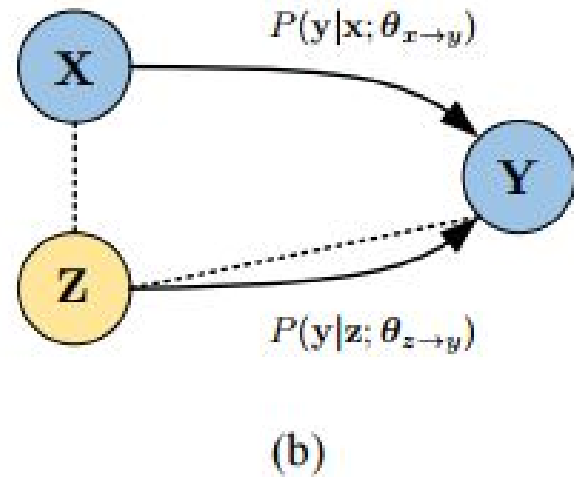
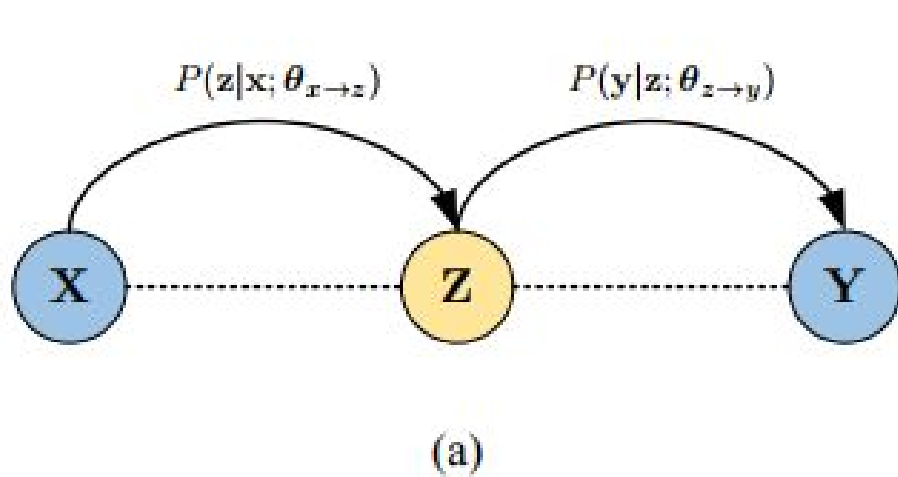


Figure 1: (a) The pivot-based approach and (b) the teacher-student approach to zero-resource neural machine translation. **X**, **Y**, and **Z** denote source, target, and pivot languages, respectively. We use a dashed line to denote that there is a parallel corpus available for the connected language pair. Solid lines with arrows represent translation directions. The pivot-based approach leverages a pivot to achieve indirect source-to-target translation: it first translates **x** into **z**, which is then translated into **y**. Our training algorithm is based on the translation equivalence assumption: if **x** is a translation of **z**, then $P(y|x; \theta_{x \rightarrow y})$ should be close to $P(y|z; \theta_{z \rightarrow y})$. Our approach directly trains the intended source-to-target model $P(y|x; \theta_{x \rightarrow y})$ (“student”) on a source-pivot parallel corpus, with the guidance of an existing pivot-to-target model $P(y|z; \hat{\theta}_{z \rightarrow y})$ (“teacher”).

What's the problem with this pivot based system?

- Due to the exponential search space of pivot sentences, the decoding process of translating an unseen source sentence x has to be divided into two steps:

$$\hat{z} = \operatorname{argmax}_{\mathbf{z}} \left\{ P(\mathbf{z}|\mathbf{x}; \hat{\boldsymbol{\theta}}_{x \rightarrow z}) \right\}, \quad (3)$$

$$\hat{y} = \operatorname{argmax}_{\mathbf{y}} \left\{ P(\mathbf{y}|\hat{z}; \hat{\boldsymbol{\theta}}_{z \rightarrow y}) \right\}. \quad (4)$$

- The above two-step decoding process potentially suffers from the error propagation problem (Zhu et al., 2013): the translation errors made in the first step (i.e., source-to-pivot translation) will affect the second step (i.e., pivot-to-target translation).
- Therefore, it is necessary to explore methods to directly model source-to-target translation without parallel corpora available.

So why the pivot?

- Simple and easy-to-implement, **pivot-based methods have been widely used in SMT for translating zero-resource language pairs** (de Gispert and Marino~ , 2006; Cohn and Lapata, 2007; Utiyama and Isahara, 2007; Wu and Wang, 2007; Bertoldi et al., 2008; Wu and Wang, 2009; Zahabi et al., 2013; Kholy et al., 2013).
- As pivot based methods are **agnostic to model structures**, they have been adapted to NMT recently (Cheng et al., 2016a; Johnson et al., 2016).



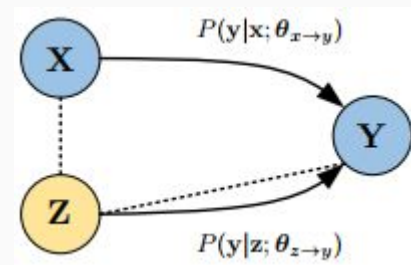
Approach

- The basic idea is to use a pre-trained pivot-to-target model (“teacher”) to guide the learning process of a source-to-target model (“student”) **without training data available on a source-pivot parallel corpus.**

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \left\{ P(\mathbf{y}|\mathbf{x}; \hat{\boldsymbol{\theta}}_{x \rightarrow y}) \right\}. \quad (1)$$

- One advantage of this approach is that Equation (1) can be used as the decision rule for decoding, which avoids the error propagation problem faced by two-step decoding in pivot-based approaches.

- As shown in Figure, we still assume that a source-pivot parallel corpus $D_{x,z}$ and a pivot-target parallel corpus $D_{z,y}$ are available. Unlike pivot-based approaches, we first use the pivot-target parallel corpus $D_{z,y}$ to obtain a teacher model $P(y|z; \hat{\boldsymbol{\theta}}_{z \rightarrow y})$, where $\hat{\boldsymbol{\theta}}_{z \rightarrow y}$ is a set of learned model parameters. Then, the teacher model “teaches” the student model $P(y|x; \boldsymbol{\theta}_{x \rightarrow y})$ on the source-pivot parallel corpus $D_{x,z}$ based on the following assumptions.



Assumption 1

If a source sentence x is a translation of a pivot sentence z , then the probability of generating a target sentence y from x should be close to that from its counterpart z .

We can further introduce a word-level assumption:

Assumption 2

If a source sentence x is a translation of a pivot sentence z , then the probability of generating a target word y from x should be close to that from its counterpart z , given the already obtained partial translation $y_{<j}$.

These two assumptions are empirically verified in our experiments

Sentence level teaching

Given a source-pivot parallel corpus $D_{x,z}$, the training objective based on Assumption 1 is defined as follows:

$$\begin{aligned} & \mathcal{J}_{\text{SENT}}(\theta_{x \rightarrow y}) \\ &= \sum_{\langle \mathbf{x}, \mathbf{z} \rangle \in D_{x,z}} \text{KL} \left(P(\mathbf{y}|\mathbf{z}; \hat{\theta}_{z \rightarrow y}) \parallel P(\mathbf{y}|\mathbf{x}; \theta_{x \rightarrow y}) \right), \quad (5) \end{aligned}$$

where the KL divergence sums over all possible target sentences:

$$\text{KL} \left(P(\mathbf{y}|\mathbf{z}; \hat{\theta}_{z \rightarrow y}) \parallel P(\mathbf{y}|\mathbf{x}; \theta_{x \rightarrow y}) \right) = \sum_{\mathbf{y}} P(\mathbf{y}|\mathbf{z}; \hat{\theta}_{z \rightarrow y}) \log \frac{P(\mathbf{y}|\mathbf{z}; \hat{\theta}_{z \rightarrow y})}{P(\mathbf{y}|\mathbf{x}; \theta_{x \rightarrow y})}. \quad (6)$$

As the teacher model parameters are fixed, the training objective can be equivalently written as

$$\begin{aligned} & \mathcal{J}_{\text{SENT}}(\theta_{x \rightarrow y}) \\ &= - \sum_{\langle \mathbf{x}, \mathbf{z} \rangle \in D_{x,z}} \mathbb{E}_{\mathbf{y}|\mathbf{z}; \hat{\theta}_{z \rightarrow y}} \left[\log P(\mathbf{y}|\mathbf{x}; \theta_{x \rightarrow y}) \right]. \quad (7) \end{aligned}$$

Sentence level training

In training, our goal is to find a set of source-to-target model parameters that minimizes the training objective:

$$\hat{\theta}_{x \rightarrow y} = \underset{\theta_{x \rightarrow y}}{\operatorname{argmin}} \left\{ \mathcal{J}_{\text{SENT}}(\theta_{x \rightarrow y}) \right\}. \quad (8)$$

- With learned source-to-target model parameters $\hat{\theta}_{x \rightarrow y}$, they use the standard decision rule as shown in Equation (1) to find the translation \hat{y} for a source sentence x .
- However, a major difficulty faced by their approach is the intractability in calculating the gradients because of the exponential search space of target sentences.
- To address this problem, it is possible to construct a subspace by either sampling (Shen et al., 2016), generating a k-best list (Cheng et al., 2016b) or mode approximation (Kim and Rush, 2016). Then, standard stochastic gradient descent algorithms can be used to optimize model parameters.

Word level training

Instead of minimizing the KL divergence between the teacher and student models at the sentence level, we further define a training objective at the word level based on Assumption 2:

$$\begin{aligned} \mathcal{J}_{\text{WORD}}(\boldsymbol{\theta}_{x \rightarrow y}) \\ = \sum_{\langle \mathbf{x}, \mathbf{z} \rangle \in D_{\mathbf{x}, \mathbf{z}}} \mathbb{E}_{\mathbf{y} | \mathbf{z}; \hat{\boldsymbol{\theta}}_{\mathbf{z} \rightarrow \mathbf{y}}} \left[J(\mathbf{x}, \mathbf{y}, \mathbf{z}, \hat{\boldsymbol{\theta}}_{\mathbf{z} \rightarrow \mathbf{y}}, \boldsymbol{\theta}_{x \rightarrow y}) \right], \quad (9) \end{aligned}$$

Where:

$$J(\mathbf{x}, \mathbf{y}, \mathbf{z}, \hat{\boldsymbol{\theta}}_{\mathbf{z} \rightarrow \mathbf{y}}, \boldsymbol{\theta}_{x \rightarrow y}) = \sum_{j=1}^{|\mathbf{y}|} \text{KL} \left(P(y | \mathbf{z}, \mathbf{y}_{< j}; \hat{\boldsymbol{\theta}}_{\mathbf{z} \rightarrow \mathbf{y}}) \parallel P(y | \mathbf{x}, \mathbf{y}_{< j}; \boldsymbol{\theta}_{x \rightarrow y}) \right). \quad (10)$$

Equation (9) suggests that the teacher model $P(y | \mathbf{z}, \mathbf{y}_{< j}; \hat{\boldsymbol{\theta}}_{\mathbf{z} \rightarrow \mathbf{y}})$ “teaches” the student model $P(y | \mathbf{x}, \mathbf{y}_{< j}; \boldsymbol{\theta}_{x \rightarrow y})$ in a word-by-word way. Note that the KL-divergence between the two models is defined at the word level:

$$\text{KL} \left(P(y | \mathbf{z}, \mathbf{y}_{< j}; \hat{\boldsymbol{\theta}}_{\mathbf{z} \rightarrow \mathbf{y}}) \parallel P(y | \mathbf{x}, \mathbf{y}_{< j}; \boldsymbol{\theta}_{x \rightarrow y}) \right) = \sum_{y \in \mathcal{V}_y} P(y | \mathbf{z}, \mathbf{y}_{< j}; \hat{\boldsymbol{\theta}}_{\mathbf{z} \rightarrow \mathbf{y}}) \log \frac{P(y | \mathbf{z}, \mathbf{y}_{< j}; \hat{\boldsymbol{\theta}}_{\mathbf{z} \rightarrow \mathbf{y}})}{P(y | \mathbf{x}, \mathbf{y}_{< j}; \boldsymbol{\theta}_{x \rightarrow y})},$$

where \mathcal{V}_y is the target vocabulary

Word level training

As the parameters of the teacher model are fixed, the training objective can be equivalently written as:

$$\begin{aligned} & \mathcal{J}_{\text{WORD}}(\boldsymbol{\theta}_{x \rightarrow y}) \\ &= - \sum_{\langle \mathbf{x}, \mathbf{z} \rangle \in D_{x, z}} \mathbb{E}_{y | \mathbf{z}; \hat{\boldsymbol{\theta}}_{z \rightarrow y}} \left[S(\mathbf{x}, y, \mathbf{z}, \hat{\boldsymbol{\theta}}_{z \rightarrow y}, \boldsymbol{\theta}_{x \rightarrow y}) \right], \quad (11) \end{aligned}$$

Where:

$$S(\mathbf{x}, y, \mathbf{z}, \hat{\boldsymbol{\theta}}_{z \rightarrow y}, \boldsymbol{\theta}_{x \rightarrow y}) = \sum_{j=1}^{|\mathbf{y}|} \sum_{y \in \mathcal{V}_y} P(y | \mathbf{z}, y_{< j}; \hat{\boldsymbol{\theta}}_{z \rightarrow y}) \times \log P(y | \mathbf{x}, y_{< j}; \boldsymbol{\theta}_{x \rightarrow y}). \quad (12)$$

They use similar approaches as described in sentence level training for approximating the full search space with sentence-level teaching. After obtaining $\hat{\boldsymbol{\theta}}_{x \rightarrow y}$, the same decision rule as shown in Equation (1) can be utilized to find the most probable target sentence \hat{y} for a source sentence x

Conclusion

Conclusion

- This paper proposes a novel framework to train the student model without parallel corpora available under the guidance of the pre-trained teacher model on a source-pivot parallel corpus.
- They introduce sentence-level and word-level teaching to guide the learning process of the student model.
- They also analyze zero-resource translation with small source-pivot data, and combine our word level sampling method with initialization and parameter freezing suggested by (Zoph et al., 2016). The experiments on the Europarl corpus show that our approach obtains a significant improvement over the pivot-based baseline.