

# Agreement-based Learning of Parallel Lexicons and Phrases from Non-Parallel Corpora

Chunyang Liu<sup>†</sup>, Yang Liu<sup>†#\*</sup>, Huanbo Luan<sup>†</sup>, Maosong Sun<sup>†#</sup>, and Heng Yu<sup>‡</sup>

<sup>†</sup> State Key Laboratory of Intelligent Technology and Systems

Tsinghua National Laboratory for Information Science and Technology

Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China

<sup>#</sup> Jiangsu Collaborative Innovation Center for Language Competence, Jiangsu, China

<sup>‡</sup> Samsung R&D Institute of China, Beijing 100028, China

{liuchunyang2012, liuyang.china, luanhuanbo}@gmail.com, sms@tsinghua.edu.cn

h0517.yu@samsung.com

## Abstract

We introduce an agreement-based approach to learning parallel lexicons and phrases from non-parallel corpora. The basic idea is to encourage two asymmetric latent-variable translation models (i.e., source-to-target and target-to-source) to agree on identifying latent phrase and word alignments. The agreement is defined at both word and phrase levels. We develop a Viterbi EM algorithm for jointly training the two unidirectional models efficiently. Experiments on the Chinese-English dataset show that agreement-based learning significantly improves both alignment and translation performance.

## 1 Introduction

Parallel corpora, which are large collections of parallel texts, serve as an important resource for inducing translation correspondences, either at the level of words (Brown et al., 1993; Smadja and McKeown, 1994; Wu and Xia, 1994) or phrases (Kupiec, 1993; Melamed, 1997; Marcu and Wong, 2002; Koehn et al., 2003). However, the availability of large-scale, wide-coverage corpora still remains a challenge even in the era of big data: parallel corpora are usually only existent for resource-rich languages and restricted to limited domains such as government documents and news articles.

Therefore, intensive attention has been drawn to exploiting non-parallel corpora for acquiring translation correspondences. Most previous efforts have concentrated on learning parallel lexicons from non-parallel corpora, including parallel sentence and lexicon extraction via bootstrapping (Fung and Cheung, 2004), inducing parallel lexicons via canonical correlation analysis (Haghighi

et al., 2008), training IBM models on monolingual corpora as decipherment (Ravi and Knight, 2011; Nuhn et al., 2012; Dou et al., 2014), and deriving parallel lexicons from bilingual word embeddings (Vulić and Moens, 2013; Mikolov et al., 2013; Vulić and Moens, 2015).

Recently, a number of authors have turned to a more challenging task: learning parallel phrases from non-parallel corpora (Zhang and Zong, 2013; Dong et al., 2015). Zhang and Zong (2013) present a method for retrieving parallel phrases from non-parallel corpora using a seed parallel lexicon. Dong et al. (2015) continue this line of research to further introduce an iterative approach to joint learning of parallel lexicons and phrases. They introduce a corpus-level latent-variable translation model in a non-parallel scenario and develop a training algorithm that alternates between (1) using a parallel lexicon to extract parallel phrases from non-parallel corpora and (2) using the extracted parallel phrases to enlarge the parallel lexicon. They show that starting from a small seed lexicon, their approach is capable of learning both new words and phrases gradually over time.

However, due to the structural divergence between natural languages as well as the presence of noisy data, only using asymmetric translation models might be insufficient to accurately identify parallel lexicons and phrases from non-parallel corpora. Dong et al. (2015) report that the accuracy on Chinese-English dataset is only around 40% after running for 70 iterations. In addition, their approach seems prone to be affected by noisy data in non-parallel corpora as the accuracy drops significantly with the increase of noise.

Since asymmetric word alignment and phrase alignment models are usually complementary, it is natural to combine them to make more accurate predictions. In this work, we propose to in-

\*Corresponding author: Yang Liu.

roduce agreement-based learning (Liang et al., 2006; Liang et al., 2008) into extracting parallel lexicons and phrases from non-parallel corpora. Based on the latent-variable model proposed by Dong et al. (2015), we propose two kinds of loss functions to take into account the agreement between both phrase alignment and word alignment in two directions. As the inference is intractable, we resort to a Viterbi EM algorithm to train the two models efficiently. Experiments on the Chinese-English dataset show that agreement-based learning is more robust to noisy data and leads to substantial improvements in phrase alignment and machine translation evaluations.

## 2 Background

Given a monolingual corpus of source language phrases  $E = \{\mathbf{e}^{(s)}\}_{s=1}^S$  and a monolingual corpus of target language phrases  $F = \{\mathbf{f}^{(t)}\}_{t=1}^T$ , we assume there exists a parallel corpus  $D = \{(\mathbf{e}^{(s)}, \mathbf{f}^{(t)}) | \mathbf{e}^{(s)} \leftrightarrow \mathbf{f}^{(t)}\}$ , where  $\mathbf{e}^{(s)} \leftrightarrow \mathbf{f}^{(t)}$  denotes that  $\mathbf{e}^{(s)}$  and  $\mathbf{f}^{(t)}$  are translations of each other.

As a long sentence in  $E$  is usually unlikely to have an translation in  $F$  and vice versa, most previous efforts build on the assumption that *phrases* are more likely to have translational equivalents on the other side (Munteanu and Marcu, 2006; Cettolo et al., 2010; Zhang and Zong, 2013; Dong et al., 2015). Such a set of phrases can be constructed by collecting either constituents of parsed sentences or strings with hyperlinks on webpages (e.g., Wikipedia). Therefore, we assume the two monolingual corpora are readily available and focus on how to extract  $D$  from  $E$  and  $F$ .

To address this problem, Dong et al. (2015) introduce a corpus-level latent-variable translation model in a non-parallel scenario:

$$P(F|E; \theta) = \sum_{\mathbf{m}} \underbrace{P(F, \mathbf{m}|E; \theta)}_{\text{phrase alignment}}, \quad (1)$$

where  $\mathbf{m}$  is *phrase alignment* and  $\theta$  is a set of model parameters. Each target phrase  $\mathbf{f}^{(t)}$  is restricted to connect to exactly one source phrase:  $\mathbf{m} = (\mathbf{m}_1, \dots, \mathbf{m}_t, \dots, \mathbf{m}_T)$ , where  $\mathbf{m}_t \in \{0, 1, \dots, S\}$ . For example,  $\mathbf{m}_t = s$  denotes that  $\mathbf{f}^{(t)}$  is aligned to  $\mathbf{e}^{(s)}$ . Note that  $\mathbf{e}^{(0)}$  represents an empty source phrase.

They follow IBM Model 1 (Brown et al., 1993) to further decompose the model as

$$P(F, \mathbf{m}|E; \theta) =$$

$$\frac{p(T|S)}{(S+1)^T} \prod_{t=1}^T P(\mathbf{f}^{(t)} | \mathbf{e}^{(\mathbf{m}_t)}; \theta), \quad (2)$$

where  $P(\mathbf{f}^{(t)} | \mathbf{e}^{(\mathbf{m}_t)}; \theta)$  is a *phrase translation* model that can be further defined as

$$\begin{aligned} & P(\mathbf{f}^{(t)} | \mathbf{e}^{(\mathbf{m}_t)}; \theta) \\ &= \delta(\mathbf{m}_t, 0)\epsilon + \\ & (1 - \delta(\mathbf{m}_t, 0)) \sum_{\mathbf{a}} \underbrace{P(\mathbf{f}^{(t)}, \mathbf{a} | \mathbf{e}^{(\mathbf{m}_t)}; \theta)}_{\text{word alignment}}. \end{aligned} \quad (3)$$

Dong et al. (2015) distinguish between *empty* and *non-empty* phrase translations. If a target phrase  $\mathbf{f}^{(t)}$  is aligned to the empty source phrase  $\mathbf{e}^{(0)}$  (i.e.,  $\mathbf{m}_t = 0$ ), they set the phrase translation probability to a fixed number  $\epsilon$ . Otherwise, conventional word alignment models such as IBM Model 1 can be used for non-empty phrase translation:

$$\begin{aligned} & P(\mathbf{f}^{(t)}, \mathbf{a} | \mathbf{e}^{(\mathbf{m}_t)}; \theta) \\ &= \frac{p(J^{(t)} | I^{(\mathbf{m}_t)})}{(I^{(\mathbf{m}_t)} + 1)^{J^{(t)}}} \prod_{j=1}^{J^{(t)}} p(\mathbf{f}_j^{(t)} | \mathbf{e}_{\mathbf{a}_j}^{(\mathbf{m}_t)}), \end{aligned} \quad (4)$$

where  $p(J|I)$  is a *length* model and  $p(f|e)$  is a *translation* model. We use  $J^{(t)}$  to denote the length of  $\mathbf{f}^{(t)}$ .

Therefore, the latent-variable model involves two kinds of latent structures: (1) *phrase alignment*  $\mathbf{m}$  between source and target phrases, (2) *word alignment*  $\mathbf{a}$  between source and target words within phrases.

Given the two monolingual corpora  $E$  and  $F$ , the training objective is to maximize the likelihood of the training data:

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \{ \mathcal{L}(\theta) \}, \quad (5)$$

where

$$\begin{aligned} \mathcal{L}(\theta) &= \log P(F|E; \theta) - \\ & \sum_I \lambda_I \left( \sum_J p(J|I) - 1 \right) - \\ & \sum_e \gamma_e \left( \sum_f p(f|e) - 1 \right) - \\ & \sum_f \sum_e \sigma(f, e, \mathbf{d}) \log \frac{\sigma(f, e, \mathbf{d})}{p(f|e)}. \end{aligned} \quad (6)$$

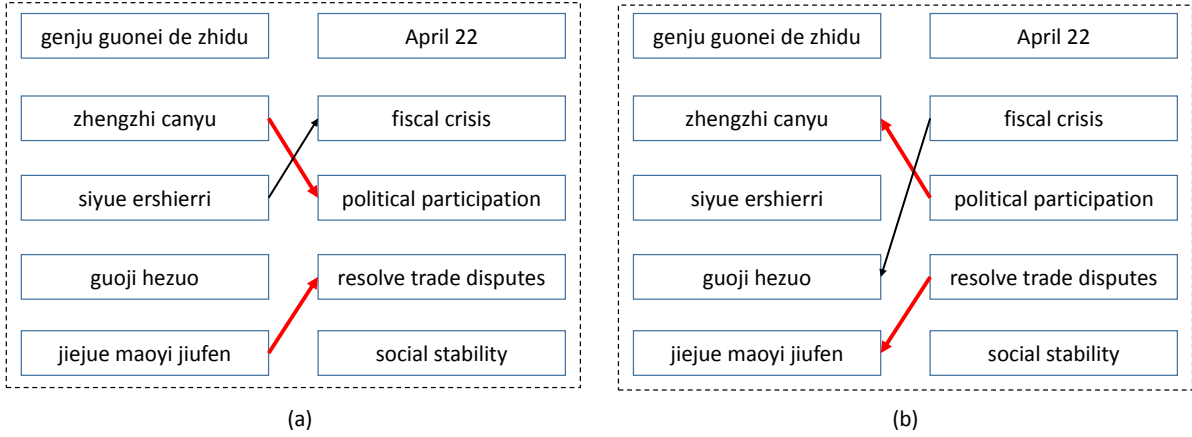


Figure 1: Agreement between (a) Chinese-to-English and (b) English-to-Chinese phrase alignments. The arrows indicate translation directions. The links on which two models agree are highlighted in bold red. The *outer agreement* loss function (see Eq. (14)) aims to encourage the agreement at the phrase level.

Note that  $\mathbf{d}$  is a small seed parallel lexicon for initializing training<sup>1</sup> and  $\sigma(f, e, \mathbf{d})$  checks whether an entry  $\langle f, e \rangle$  exists in  $\mathbf{d}$ .

Given the monolingual corpora and the optimized model parameters, the Viterbi phrase alignment is calculated as

$$\mathbf{m}^* = \operatorname{argmax}_{\mathbf{m}} \left\{ P(F, \mathbf{m} | E; \theta^*) \right\} \quad (7)$$

$$= \operatorname{argmax}_{\mathbf{m}} \left\{ \prod_{t=1}^T P(\mathbf{f}^{(t)} | \mathbf{e}^{(m_t)}; \theta^*) \right\}. \quad (8)$$

Finally, parallel lexicons can be derived from the translation probability table of IBM model 1  $\theta^*$  and parallel phrases can be collected from the Viterbi phrase alignment  $\mathbf{m}^*$ . This process iterates and enlarges parallel lexicons and phrases gradually over time.

As it is very challenging to extract parallel phrases from non-parallel corpora, unidirectional models might only capture partial aspects of translation modeling on non-parallel corpora. Indeed, Dong et al. (2015) find that the accuracy of phrase alignment is only around 50% on the Chinese-English dataset. More importantly, their approach seems to be vulnerable to noise as the accuracy drops significantly with the increase of noise. As source-to-target and target-to-source translation models are usually complementary (Och and Ney, 2003; Koehn et al., 2003; Liang et al., 2006),

<sup>1</sup>Due to the difficulty of learning translation correspondences from non-parallel corpora, many authors have assumed that a small seed lexicon is readily available (Gaussier et al., 2004; Zhang and Zong, 2013; Vulić and Moens, 2013; Mikolov et al., 2013; Dong et al., 2015).

it is appealing to combine them to improve alignment accuracy.

### 3 Approach

#### 3.1 Agreement-based Learning

The basic idea of our work is to encourage the source-to-target and target-to-source translation models to agree on both phrase and word alignments.

For example, Figure 1 shows two example Chinese-to-English and English-to-Chinese phrase alignments on the same non-parallel data. As each model only captures partial aspects of translation modeling, our intuition is that the links on which two models agree (highlighted in red) are more likely to be correct.

More formally, let  $P(F|E; \vec{\theta})$  be a source-to-target translation model and  $P(E|F; \overleftarrow{\theta})$  be a target-to-source model, where  $\vec{\theta}$  and  $\overleftarrow{\theta}$  are corresponding model parameters. We use  $\vec{\mathbf{m}} = (\vec{\mathbf{m}}_1, \dots, \vec{\mathbf{m}}_t, \dots, \vec{\mathbf{m}}_T)$  to denote source-to-target phrase alignment. Likewise, the target-to-source phrase alignment is denoted by  $\overleftarrow{\mathbf{m}} = (\overleftarrow{\mathbf{m}}_1, \dots, \overleftarrow{\mathbf{m}}_s, \dots, \overleftarrow{\mathbf{m}}_S)$ .

To ease the comparison between  $\vec{\mathbf{m}}$  and  $\overleftarrow{\mathbf{m}}$ , we represent them as sets of non-empty links equivalently:

$$\vec{\mathbf{m}} = \left\{ \langle \vec{\mathbf{m}}_t, t \rangle \mid \vec{\mathbf{m}}_t \neq 0 \right\} \quad (9)$$

$$\overleftarrow{\mathbf{m}} = \left\{ \langle s, \overleftarrow{\mathbf{m}}_s \rangle \mid \overleftarrow{\mathbf{m}}_s \neq 0 \right\}. \quad (10)$$

For example, suppose the source-to-target and target-to-source phrase alignments are  $\vec{\mathbf{m}} =$

```

1: procedure VITERBIEM( $E, F, \mathbf{d}$ )
2:   Initialize  $\Theta^{(0)}$ 
3:   for all  $k = 1, \dots, K$  do
4:      $\hat{\mathbf{m}}^{(k)} \leftarrow \text{SEARCH}(E, F, \Theta^{(k-1)})$ 
5:      $\Theta^{(k)} \leftarrow \text{UPDATE}(E, F, \mathbf{d}, \hat{\mathbf{m}}^{(k)})$ 
6:   end for
7:   return  $\hat{\mathbf{m}}^{(K)}, \Theta^{(K)}$ 
8: end procedure

```

Figure 2: A Viterbi EM algorithm for agreement-based learning of parallel lexicons and phrases from non-parallel corpora.  $F$  and  $E$  are non-parallel corpora,  $\mathbf{d}$  is a seed parallel lexicon,  $\Theta^{(k)}$  is the set of model parameters at the  $k$ -th iteration,  $\hat{\mathbf{m}}^{(k)}$  is the Viterbi phrase alignment on which two models agree at the  $k$ -th iteration.

(2, 3, 0, 0) and  $\hat{\mathbf{m}} = (0, 1, 2)$ . The equivalent link sets are  $\bar{\mathbf{m}} = \{\langle 2, 1 \rangle, \langle 3, 2 \rangle\}$  and  $\overleftarrow{\mathbf{m}} = \{\langle 2, 1 \rangle, \langle 3, 2 \rangle\}$ . Therefore,  $\bar{\mathbf{m}}$  is said to be *equal* to  $\overleftarrow{\mathbf{m}}$  (i.e.,  $\delta(\bar{\mathbf{m}}, \overleftarrow{\mathbf{m}}) = 1$ ).

Following Liang et al. (2006), we introduce a new training objective that favors the agreement between two unidirectional models:

$$\begin{aligned} & \mathcal{J}(\vec{\theta}, \overleftarrow{\theta}) \\ &= \log P(F|E; \vec{\theta}) + \log P(E|F; \overleftarrow{\theta}) - \\ & \log \sum_{\bar{\mathbf{m}}, \overleftarrow{\mathbf{m}}} P(\bar{\mathbf{m}}|E, F; \vec{\theta}) P(\overleftarrow{\mathbf{m}}|F, E; \overleftarrow{\theta}) \\ & \quad \times \Delta(E, F, \bar{\mathbf{m}}, \overleftarrow{\mathbf{m}}, \vec{\theta}, \overleftarrow{\theta}), \end{aligned} \quad (11)$$

where the posterior probabilities in two directions are defined as

$$P(\bar{\mathbf{m}}|E, F; \vec{\theta}) = \prod_{t=1}^T \frac{P(\mathbf{f}^{(t)} | \mathbf{e}^{(\bar{\mathbf{m}}_t)}; \vec{\theta})}{\sum_{s=0}^S P(\mathbf{f}^{(t)} | \mathbf{e}^{(s)}; \vec{\theta})} \quad (12)$$

$$P(\overleftarrow{\mathbf{m}}|F, E; \overleftarrow{\theta}) = \prod_{s=1}^S \frac{P(\mathbf{e}^{(s)} | \mathbf{f}^{(\overleftarrow{\mathbf{m}}_s)}; \overleftarrow{\theta})}{\sum_{t=0}^T P(\mathbf{e}^{(s)} | \mathbf{f}^{(t)}; \overleftarrow{\theta})}. \quad (13)$$

The *loss function*  $\Delta(E, F, \bar{\mathbf{m}}, \overleftarrow{\mathbf{m}}, \vec{\theta}, \overleftarrow{\theta})$  measures the disagreement between the two models.

## 3.2 Outer Agreement

### 3.2.1 Definition

A straightforward loss function is to force the two models to generate identical phrase alignments:

$$\Delta_{\text{outer}}(E, F, \bar{\mathbf{m}}, \overleftarrow{\mathbf{m}}, \vec{\theta}, \overleftarrow{\theta}) = 1 - \delta(\bar{\mathbf{m}}, \overleftarrow{\mathbf{m}}). \quad (14)$$

We refer to Eq. (14) as *outer agreement* since it only considers phrase alignment and ignores the word alignment within aligned phrases.

### 3.2.2 Training Objective

Since the outer agreement forces two models to generate identical phrase alignments, the training objective can be written as

$$\begin{aligned} & \mathcal{J}_{\text{outer}}(\vec{\theta}, \overleftarrow{\theta}) \\ &= \log P(F|E; \vec{\theta}) + \log P(E|F; \overleftarrow{\theta}) + \\ & \log \sum_{\mathbf{m}} P(\mathbf{m}|E, F; \vec{\theta}) P(\mathbf{m}|F, E; \overleftarrow{\theta}), \end{aligned} \quad (15)$$

where  $\mathbf{m}$  is a phrase alignment on which two models agree.

The partial derivatives of the training objective with respect to source-to-target model parameters  $\vec{\theta}$  are given by

$$\begin{aligned} & \frac{\partial \mathcal{J}_{\text{outer}}(\vec{\theta}, \overleftarrow{\theta})}{\partial \vec{\theta}} \\ &= \frac{\partial P(F|E; \vec{\theta}) / \partial \vec{\theta}}{P(F|E; \vec{\theta})} + \\ & \frac{\mathbb{E}_{\mathbf{m}|F, E; \overleftarrow{\theta}} [\partial P(F|E; \vec{\theta}) / \partial \vec{\theta}]}{\sum_{\mathbf{m}} P(\mathbf{m}|E, F; \vec{\theta}) P(\mathbf{m}|F, E; \overleftarrow{\theta})}. \end{aligned} \quad (16)$$

The partial derivatives with respect to  $\overleftarrow{\theta}$  are defined likewise.

### 3.2.3 Training Algorithm

As the expectation in Eq. (16) is usually intractable to calculate due to the exponential search space of phrase alignment, we follow Dong et al. (2015) to use a Viterbi EM algorithm instead.

As shown in Figure 2, the algorithm takes a set of source phrases  $E$ , a set of target phrases  $F$ , and a seed parallel lexicon  $\mathbf{d}$  as input (line 1). After initializing model parameters  $\Theta = \{\vec{\theta}, \overleftarrow{\theta}\}$  (line 2), the algorithm calls the procedure  $\text{ALIGN}(F, E, \Theta)$  to compute the Viterbi phrase alignment between  $E$  and  $F$  on which two models agree. Then, the algorithm updates the two models by normalizing counts collected from the Viterbi phrase alignment. The process iterates for  $K$  iterations and returns the final Viterbi phrase alignment and model parameters.

### 3.2.4 Computing Viterbi Phrase Alignments

The procedure  $\text{ALIGN}(F, E, \Theta)$  computes the Viterbi phrase alignment  $\hat{\mathbf{m}}$  between  $E$  and  $F$  on which two models agree as follows:

$$\hat{\mathbf{m}} = \underset{\mathbf{m}}{\text{argmax}} \left\{ P(\mathbf{m}|E, F; \vec{\theta}) \times P(\mathbf{m}|F, E; \overleftarrow{\theta}) \right\}. \quad (17)$$

Unfortunately, due to the exponential search space of phrase alignment, computing  $\hat{\mathbf{m}}$  is also intractable. As a result, we approximate it as the intersection of two unidirectional Viterbi phrase alignments:

$$\hat{\mathbf{m}} \approx \overrightarrow{\mathbf{m}}^* \cap \overleftarrow{\mathbf{m}}^*, \quad (18)$$

where the unidirectional Viterbi phrase alignments are calculated as

$$\overrightarrow{\mathbf{m}}^* = \operatorname{argmax}_{\overrightarrow{\mathbf{m}}} \left\{ \prod_{t=1}^T P(\mathbf{f}^{(t)} | \mathbf{e}^{(\overrightarrow{\mathbf{m}}_t)}; \overrightarrow{\theta}) \right\} \quad (19)$$

$$\overleftarrow{\mathbf{m}}^* = \operatorname{argmax}_{\overleftarrow{\mathbf{m}}} \left\{ \prod_{s=1}^S P(\mathbf{e}^{(s)} | \mathbf{f}^{(\overleftarrow{\mathbf{m}}_s)}; \overleftarrow{\theta}) \right\}. \quad (20)$$

The source-to-target Viterbi phrase alignment is calculated as

$$\overrightarrow{\mathbf{m}}^* = \operatorname{argmax}_{\overrightarrow{\mathbf{m}}} \left\{ P(\overrightarrow{\mathbf{m}} | E, F; \overrightarrow{\theta}) \right\} \quad (21)$$

$$= \operatorname{argmax}_{\overrightarrow{\mathbf{m}}} \left\{ \prod_{t=1}^T P(\mathbf{f}^{(t)} | \mathbf{e}^{(\overrightarrow{\mathbf{m}}_t)}; \overrightarrow{\theta}) \right\}. \quad (22)$$

Dong et al. (2015) indicate that computing the Viterbi alignment for individual target phrases is *independent* and only need to focus on finding the most probable source phrase for each target phrase:

$$\overrightarrow{\mathbf{m}}_t^* = \operatorname{argmax}_{s \in \{0, 1, \dots, S\}} \left\{ P(\mathbf{f}^{(t)} | \mathbf{e}^{(s)}; \overrightarrow{\theta}) \right\}. \quad (23)$$

This can be cast as a translation retrieval problem (Zhang and Zong, 2013; Dong et al., 2014). Please refer to (Dong et al., 2015) for more details. The target-to-source Viterbi phrase alignment can be calculated similarly.

### 3.2.5 Updating Model Parameters

Following Liang et al. (2006), we collect counts of model parameters only from the agreement term.<sup>2</sup>

Given the agreed Viterbi phrase alignment  $\hat{\mathbf{m}}$ , the count of the source-to-target length model  $p(J|I)$  is given by

$$c(J|I; E, F) = \sum_{\langle s, t \rangle \in \hat{\mathbf{m}}} \delta(J^{(t)}, J) \delta(I^{(s)}, I). \quad (24)$$

The new length probabilities can be obtained by

$$p(J|I) = \frac{c(J|I; E, F)}{\sum_{J'} c(J'|I; E, F)}. \quad (25)$$

<sup>2</sup>We experimented with collecting counts from both the unidirectional and agreement terms but obtained much worse results than counting only from the agreement term.

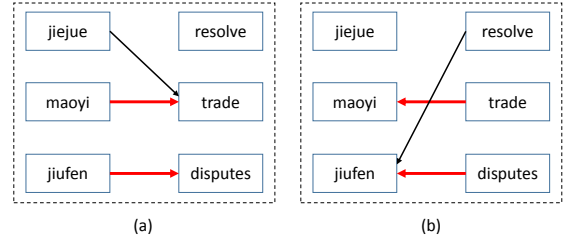


Figure 3: Agreement between (a) Chinese-to-English and (b) English-to-Chinese word alignments. The links on which two models agree are highlighted in red. The *inner agreement* loss function (see Eq. (28)) aims to encourage the agreement at both the phrase and word levels.

The count of the source-to-target translation model  $p(f|e)$  is given by

$$\begin{aligned} c(f|e; E, F) &= \sum_{\langle s, t \rangle \in \hat{\mathbf{m}}} \frac{p(f|e)}{\sum_{i=0}^{I^{(s)}} p(f|e_i^{(s)})} \times \\ &\quad \sum_{j=1}^{J^{(t)}} \delta(f, \mathbf{f}_j^{(t)}) \sum_{i=0}^{I^{(s)}} \delta(e, \mathbf{e}_i^{(s)}) \\ &\quad + \sigma(f, e, \mathbf{d}). \end{aligned} \quad (26)$$

The new translation probabilities can be obtained by

$$p(f|e) = \frac{c(f|e; E, F)}{\sum_{f'} c(f'|e; E, F)}. \quad (27)$$

Counts of target-to-source length and translation models can be calculated in a similar way.

## 3.3 Inner Agreement

### 3.3.1 Definition

As the outer agreement only considers the phrase alignment, the *inner agreement* takes both phrase alignment and word alignment into consideration:

$$\begin{aligned} \Delta_{\text{inner}}(E, F, \overrightarrow{\mathbf{m}}, \overleftarrow{\mathbf{m}}, \overrightarrow{\theta}, \overleftarrow{\theta}) &= -\delta(\overrightarrow{\mathbf{m}}, \overleftarrow{\mathbf{m}}) \times \\ &\quad \sum_{\langle s, t \rangle \in \overrightarrow{\mathbf{m}}} \sum_{\overrightarrow{\mathbf{a}}, \overleftarrow{\mathbf{a}}} P(\overrightarrow{\mathbf{a}} | \mathbf{e}^{(s)}, \mathbf{f}^{(t)}; \overrightarrow{\theta}) \times \\ &\quad P(\overleftarrow{\mathbf{a}} | \mathbf{f}^{(t)}, \mathbf{e}^{(s)}; \overleftarrow{\theta}) \times \\ &\quad \delta(\overrightarrow{\mathbf{a}}, \overleftarrow{\mathbf{a}}). \end{aligned} \quad (28)$$

For example, Figure 3 shows two examples of Chinese-to-English and English-to-Chinese word alignments. The shared links are highlighted in

red. Our intuition is that a source phrase and a target phrase are more likely to be translations of each other if the two translation models also agree on word alignment within aligned phrases.

### 3.3.2 Training Objective and Algorithm

The training objective for inner agreement is given by

$$\begin{aligned} & \mathcal{J}_{\text{inner}}(\vec{\theta}, \overleftarrow{\theta}) \\ = & \log P(F|E; \vec{\theta}) + \log P(E|F; \overleftarrow{\theta}) + \\ & \log \sum_{\mathbf{m}} P(\mathbf{m}|E, F; \vec{\theta}) P(\mathbf{m}|F, E; \overleftarrow{\theta}) \times \\ & \sum_{\langle s,t \rangle \in \mathbf{m}} \sum_{\mathbf{a}} P(\mathbf{a}|\mathbf{e}^{(s)}, \mathbf{f}^{(t)}; \vec{\theta}) \times \\ & P(\mathbf{a}|\mathbf{f}^{(t)}, \mathbf{e}^{(s)}; \overleftarrow{\theta}). \end{aligned} \quad (29)$$

We still use the Viterbi EM algorithm as shown in Figure 2 for training the two models.

### 3.3.3 Computing Viterbi Phrase Alignments

The agreed Viterbi phrase alignment is defined as

$$\begin{aligned} \hat{\mathbf{m}} = \operatorname{argmax}_{\mathbf{m}} & \left\{ P(\mathbf{m}|E, F; \vec{\theta}) P(\mathbf{m}|F, E; \overleftarrow{\theta}) \right. \\ & \times \sum_{\langle s,t \rangle \in \mathbf{m}} \sum_{\mathbf{a}} P(\mathbf{a}|\mathbf{e}^{(s)}, \mathbf{f}^{(t)}; \vec{\theta}) \\ & \left. \times P(\mathbf{a}|\mathbf{f}^{(t)}, \mathbf{e}^{(s)}; \overleftarrow{\theta}) \right\}. \end{aligned} \quad (30)$$

As computing  $\hat{\mathbf{m}}$  is intractable, we still approximate it using the intersection of two unidirectional Viterbi phrase alignments (see Eq. (18)). The source-to-target Viterbi phrase alignment is calculated as

$$\begin{aligned} \vec{\mathbf{m}}^* = \operatorname{argmax}_{\vec{\mathbf{m}}} & \left\{ P(\vec{\mathbf{m}}|E, F; \vec{\theta}) \times \right. \\ & \sum_{\langle s,t \rangle \in \vec{\mathbf{m}}} \sum_{j=1}^{J^{(t)}} \sum_{i=1}^{I^{(s)}} P(\langle i, j \rangle | \mathbf{e}^{(s)}, \mathbf{f}^{(t)}; \vec{\theta}) \times \\ & \left. P(\langle i, j \rangle | \mathbf{f}^{(t)}, \mathbf{e}^{(s)}; \overleftarrow{\theta}) \right\}, \end{aligned} \quad (31)$$

where  $P(\langle i, j \rangle | \mathbf{e}^{(s)}, \mathbf{f}^{(t)}; \vec{\theta})$  is source-to-target link posterior probability of the link  $\langle i, j \rangle$  being present (or absent) in the word alignment according to the source-to-target model,  $P(\langle i, j \rangle | \mathbf{f}^{(t)}, \mathbf{e}^{(s)}; \overleftarrow{\theta})$  is target-to-source link posterior probability. We follow Liang et al. (2006) to use the product of link posteriors to encourage the agreement at the level of word alignment.

We use a coarse-to-fine approach (Dong et al., 2015) to compute the Viterbi alignment: first retrieving a coarse set of candidate source phrases using translation probabilities and then selecting the candidate with the highest score according to Eq. (31). The target-to-source Viterbi phrase alignment can be calculated similarly.

### 3.3.4 Updating Model Parameters

Given the agreed Viterbi phrase alignment  $\hat{\mathbf{m}}$ , the count of the source-to-target length model  $p(J|I)$  is still given by Eq. (24). The count of the translation model  $p(f|e)$  is calculated as

$$\begin{aligned} & c(f|e; E, F) \\ = & \sum_{\langle s,t \rangle \in \hat{\mathbf{m}}} \sum_{i=1}^{I^{(s)}} \sum_{j=1}^{J^{(t)}} P(\langle i, j \rangle | \mathbf{e}^{(s)}, \mathbf{f}^{(t)}; \vec{\theta}) \times \\ & P(\langle i, j \rangle | \mathbf{f}^{(t)}, \mathbf{e}^{(s)}; \overleftarrow{\theta}) \times \\ & \delta(f, \mathbf{f}^{(t)}) \delta(e, \mathbf{e}^{(s)}) \\ & + \sigma(f, e, \mathbf{d}). \end{aligned} \quad (32)$$

Counts of target-to-source length and translation models can be calculated in a similar way.

## 4 Experiments

In this section, we evaluate our approach in two tasks: phrase alignment (Section 4.1) and machine translation (Section 4.2).

### 4.1 Alignment Evaluation

#### 4.1.1 Evaluation Metrics

Given two monolingual corpora  $E$  and  $F$ , we suppose there exists a ground truth parallel corpus  $G$  and denote an extracted parallel corpus as  $D$ . The quality of an extracted parallel corpus can be measured by  $F1 = 2|D \cap G| / (|D| + |G|)$ .

#### 4.1.2 Data Preparation

Although it is appealing to apply our approach to dealing with real-world non-parallel corpora, it is time-consuming and labor-intensive to manually construct a ground truth parallel corpus. Therefore, we follow Dong et al. (2015) to build synthetic  $E$ ,  $F$ , and  $G$  to facilitate the evaluation.

We first extract a set of parallel phrases from a sentence-level parallel corpus using the state-of-the-art phrase-based translation system Moses (Koehn et al., 2007) and discard low-probability parallel phrases. Then,  $E$  and  $F$  can be constructed by corrupting the parallel phrase set by

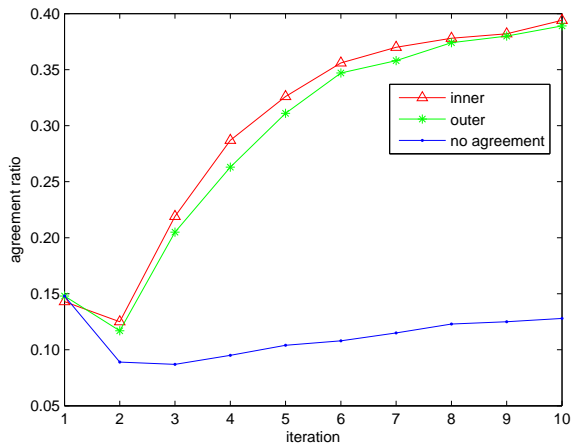


Figure 4: Comparison of agreement ratios on the development set.

seed	C → E	E → C	Outer	Inner
50	4.1	4.8	60.8	66.2
100	5.1	5.5	65.6	69.8
500	7.5	8.4	70.4	72.5
1,000	22.4	23.1	73.6	74.3

Table 1: Effect of seed lexicon size in terms of F1 on the development set.

adding irrelevant source and target phrases randomly. Note that the parallel phrase set can serve as the ground truth parallel corpus  $G$ . We refer to the non-parallel phrases in  $E$  and  $F$  as *noise*.

From LDC Chinese-English parallel corpora, we constructed a *development set* and a *test set*. The development set contains 20K parallel phrases, 20K noisy Chinese phrases, and 20K noisy English phrases. The test set contains 20K parallel phrases, 180K noisy Chinese phrases, and 180K noisy English phrases. The seed parallel lexicon contains 1K entries.

### 4.1.3 Comparison of Agreement Ratios

We introduce *agreement ratio* to measure to what extent two unidirectional models agree on phrase alignment:

$$\text{ratio} = \frac{2|\vec{\mathbf{m}}^* \cap \overleftarrow{\mathbf{m}}^*|}{|\vec{\mathbf{m}}^*| + |\overleftarrow{\mathbf{m}}^*|}. \quad (33)$$

Figure 4 shows the agreement ratios of independent training (“no agreement”), joint training with the outer agreement (“outer”), and joint training with the inner agreement (“inner”). We find that independently trained unidirectional models

noise		C → E	E → C	Outer	Inner
C	E				
0	0	58.5	61.2	86.5	86.1
0	10K	41.0	54.4	83.6	83.8
0	20K	28.3	48.3	80.1	81.2
10K	0	54.7	43.1	84.9	84.3
20K	0	50.4	31.4	83.8	83.6
10K	10K	34.9	34.4	80.0	79.7
20K	20K	22.4	23.1	73.6	74.3

Table 2: Effect of noise in terms of F1 on the development set.

hardly agree on phrase alignment, suggesting that each model can only capture partial aspects of translation modeling on non-parallel corpora. In contrast, imposing the agreement term significantly increases the agreement ratios: after 10 iterations, about 40% of phrase alignment links are shared by two models.

### 4.1.4 Effect of Seed Lexicon Size

Table 1 shows the F1 scores of the Chinese-to-English model (“C → E”), the English-to-Chinese model (“E → C”), joint learning based on the outer agreement (“outer”), and joint learning based on the inner agreement (“inner”) over various sizes of seed lexicons on the development set.

We find that agreement-based learning obtains substantial improvements over independent learning across all sizes. More importantly, even with a seed lexicon containing only 50 entries, agreement-based learning is able to achieve F1 scores above 60%. The inner agreement performs better than the outer agreement by taking the consensus at the word level into account.

### 4.1.5 Effect of Noise

Table 2 demonstrates the effect of noise on the development set. In row 1, “0+0” denotes there is no noise, which can be seen as an upper bound. Adding noise, either on the Chinese side or on the English side, deteriorates the F1 scores for all methods. Adding noise on the English side makes predicting phrase alignment in the C → E direction more challenging due to the enlarged search space. The situation is similar in the reverse direction. It is clear that agreement-based learning is more robust to noise: while independent training suffers from a reduction of 40% in terms of F1 for the “20K + 20K” setting, agreement-based learning still achieves F1 scores over 70%.

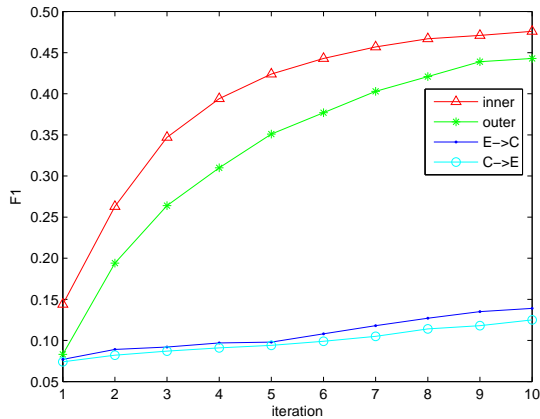


Figure 5: Comparison of F1 scores on the test set.

Chinese	jingji
English	<i>economy</i>
Chinese	jialebi
English	<i>caribbean</i>
Chinese	zhengzhi huanjing
English	political environment
Chinese	jiaoyisuo shichang jiage zhishu
English	<i>exchange market price index</i>
Chinese	qianding bianjing maoyi xieding
English	<i>signed border trade agreements</i>

Table 3: Example learned parallel lexicons and phrases. New words that are not included in the seed lexicon are highlighted in italic.

#### 4.1.6 Results

Figure 5 gives the final results on the test set. We find that agreement-based training achieves significant improvements over independent training. By considering the consensus on both phrase and word alignments, the inner agreement significantly outperforms the outer agreement. Notice that Dong et al. (2015) only add noise on one side while we add noisy phrases on both sides, which makes phrase alignment more challenging.

Table 3 shows example learned parallel words and phrases. The lexicon is built from the translation table by retaining high-probability word pairs. Therefore, our approach is capable of learning both new words and new phrases unseen in the seed lexicon.

## 4.2 Translation Evaluation

Following Zhang and Zong (2013) and Dong et al. (2015), we evaluate our approach on domain

adaptation for machine translation.

The data set consists of two in-domain non-parallel corpora and an out-domain parallel corpus. The in-domain non-parallel corpora consists of 2.65M Chinese phrases and 3.67M English phrases extracted from LDC news articles. We use a small out-domain parallel corpus extracted from financial news of FTChina which contains 10K phrase pairs. The task is to extract a parallel corpus from in-domain non-parallel corpora starting from a small out-domain parallel corpus.

We use the state-of-the-art translation system Moses (Koehn et al., 2007) and evaluate the performance on Chinese-English NIST datasets. The development set is NIST 2006 and the test set is NIST 2005. The evaluation metric is case-insensitive BLEU4 (Papineni et al., 2002). We use the SRILM toolkit (Stolcke, 2002) to train a 4-gram English language model on a monolingual corpus with 399M English words.

Table 4 shows the results. At iteration 0, only the out-domain corpus is used and the BLEU score is 5.61. All methods iteratively extract parallel phrases from non-parallel corpora and enlarge the extracted parallel corpus. We find that agreement-based learning achieves much higher BLEU scores while obtains a smaller parallel corpus as compared with independent learning. One possible reason is that the agreement-based learning rules out most unlikely phrase pairs by encouraging consensus between two models.

## 5 Conclusion

We have presented agreement-based training for learning parallel lexicons and phrases from non-parallel corpora. By modeling the agreement on both phrase alignment and word alignment, our approach achieves significant improvements in both alignment and translation evaluations.

In the future, we plan to apply our approach to real-world non-parallel corpora to further verify its effectiveness. It is also interesting to extend the phrase translation model to more sophisticated models such as IBM models 2-5 (Brown et al., 1993) and HMM (Vogel and Ney, 1996).

## Acknowledgments

We sincerely thank the reviewers for their valuable suggestions. We also thank Meng Zhang, Yankai Lin, Shiqi Shen, Meiping Dong and Congyu Fu for their insightful discussions. Yang Liu is sup-



Iteration	Corpus Size				BLEU			
	E→C	C→E	Outer	Inner	E→C	C→E	Outer	Inner
0	10k				5.61			
1	145k	162k	59k	73k	8.65	8.90	13.53	13.74
2	195k	215k	69k	101k	8.82	9.47	15.26	15.61
3	209k	231k	88k	132k	8.42	9.29	16.88	16.94
4	214k	238k	106k	159k	8.46	9.27	17.15	17.83
5	217k	241k	123k	181k	8.87	9.40	17.94	18.89
6	219k	243k	137k	197k	8.52	9.30	18.56	19.47
7	222k	247k	140k	207k	8.81	9.22	18.72	19.46
8	224k	249k	153k	220k	8.71	9.26	18.84	19.50
9	227k	251k	159k	233k	8.92	9.35	19.05	19.63
10	229k	254k	163k	239k	8.33	9.06	19.39	19.78

Table 4: Results on domain adaptation for machine translation.

ported by the National Natural Science Foundation of China (No. 61522204), the 863 Program (2015AA011808), and Samsung R&D Institute of China. Huanbo Luan is supported by the National Natural Science Foundation of China (No. 61303075). Maosong Sun is supported by the Major Project of the National Social Science Foundation of China (13&ZD190).

## References

- Peter F. Brown, Stephen A. Della Pietra, Vincent J. Della Pietra, and Robert L. Mercer. 1993. The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*.
- Mauro Cettolo, Marcello Federico, and Nicola Bertoldi. 2010. Mining parallel fragments from comparable texts. In *Proceedings of IWSLT*.
- Meiping Dong, Yong Cheng, Yang Liu, Jia Xu, Maosong Sun, Tatsuya Izuha, and Jie Hao. 2014. Query lattice for translation retrieval. In *Proceedings of COLING*.
- Meiping Dong, Yang Liu, Huanbo Luan, Maosong Sun, Tatsuya Izuha, and Dakun Zhang. 2015. Iterative learning of parallel lexicons and phrases from non-parallel corpora. In *Proceedings of IJCAI*.
- Qing Dou, Ashish Vaswani, and Kevin Knight. 2014. Beyond parallel data: Joint word alignment and decipherment improves machine translation. In *Proceedings of EMNLP*.
- Pascale Fung and Percy Cheung. 2004. Mining very-non-parallel corpora: Parallel sentence and lexicon extraction via bootstrapping and em. In *Proceedings of EMNLP*.
- Eric Gaussier, J.M. Renders, I. Matveeva, C. Goutte, and H. Dejean. 2004. A geometric view on bilingual lexicon extraction from comparable corpora. In *Proceedings of ACL*.
- Aria Haghighi, Percy Liang, Taylor Berg-Kirkpatrick, and Dan Klein. 2008. Learning bilingual lexicons from monolingual corpora. In *Proceedings of ACL*.
- Philipp Koehn, Franz J. Och, and Daniel Marcu. 2003. Statistical phrase-based translation. In *Proceedings of NAACL*.
- Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, and Evan Herbst. 2007. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the Association for Computational Linguistics Companion Volume Proceedings of the Demo and Poster Sessions*.
- Julian Kupiec. 1993. An algorithm for finding noun phrase correspondences in bilingual corpora. In *Proceedings of ACL*.
- Percy Liang, Ben Taskar, and Dan Klein. 2006. Alignment by agreement. In *Proceedings of NAACL*.
- Percy Liang, Dan Klein, and I. Jordan, Michael. 2008. Alignment-based learning. In *Proceedings of NIPS*.
- Daniel Marcu and Daniel Wong. 2002. A phrase-based joint probability model for statistical machine translation. In *Proceedings of EMNLP*.
- I. Dan Melamed. 1997. Automatic discovery of non-compositional compounds in parallel data. In *Proceedings of EMNLP*.
- Tomas Mikolov, Quoc V Le, and Ilya Sutskever. 2013. Exploiting similarities among languages for machine translation. arXiv:1309.4168.

- Dragos Stefan Munteanu and Daniel Marcu. 2006. Extracting parallel sub-sentential fragments from non-parallel corpora. In *Proceedings of ACL*.
- Malte Nuhn, Arne Mauser, and Hermann Ney. 2012. Deciphering foreign language by combining language models and context vectors. In *Proceedings of ACL*.
- Franz J. Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of ACL*.
- Sujith Ravi and Kevin Knight. 2011. Deciphering foreign language. In *Proceedings of ACL*.
- Frank Smadja and Kathleen McKeown. 1994. Translating collocations for use in bilingual lexicons. In *Proceedings of the ARPA Human Language Technology Workshop*.
- Andreas Stolcke. 2002. Srilm - an extensible language modeling toolkit. In *Proceedings of ICSLP*.
- Stephan Vogel and Hermann Ney. 1996. Hhm-based word alignment in statistical translation. In *Proceedings of COLING*.
- Ivan Vulić and Marie-Francine Moens. 2013. A study on bootstrapping bilingual vector spaces from non-parallel data (and nothing else). In *Proceedings of EMNLP*.
- Ivan Vulić and Marie-Francine Moens. 2015. Bilingual word embeddings from non-parallel document-aligned data applied to bilingual lexicon induction. In *Proceedings of ACL*.
- Dekai Wu and Xuanyin Xia. 1994. Learning an english-chinese lexicon from a parallel corpus. In *Proceedings of the ARPA Human Language Technology Workshop*.
- Jiajun Zhang and Chengqing Zong. 2013. Learning a phrase-based translation model from monolingual data with application to domain adaptation. In *Proceedings of ACL*.