

Chapter N: Experimenting Interlinear Glossing Text

Yuan-Lu Chen

March 21, 2018

(
Assuming that in the previous chapters the following points are addressed already:

- The nature of glosses has been well-explained (Target audience: CS people without any formal linguistics background):
 - What glosses are: A basic intro of interlinear gloss for non-linguists
 - The golden nature of glosses (encodes NON-LINEAR syntax (i.e. structure parse) and semantics information)
 - The potential of gloss:
 - * potential: providing disambiguation, labeling important grammar morphemes in the source language, providing morphological analysis, providing one-to-many and many-to-one relations of source tokens and target tokens.
- A history of machine translation, and a non-mathy description of the methods of doing machine translation. (Target reader: theoretical linguists)

)

1 Introduction

The Innovation is to incorporate the gloss information of Interlinear Glossed Text data into machine translation.

In supervised machine learning models, two factors effects the performance of the trained systems (Kotsiantis et al., 2007): a.) the quality of the training data and b.) the choices of the features. The properties of the gloss data as described in *CHAPTERXYZ* make it a better training data than natural language data (Scottish Gaelic in the current case) for the following reasons. First, glosses are more purified than natural language words. The most ideal meaning representation system should be built with one-meaning-to-one-representation mappings; in other words, a meaning is mapped to one and only one representation. Natural languages fail to do so, given that synonyms and ambiguous words/phrases are ubiquitous in natural languages. Glosses provide this one-to-one mapping. Second, the gloss data provides hierarchical (non-linear) syntactic

parsing information to some degree. To determine what the gloss of a word is, linguists have to look for hierarchical (non-linear) context information.

Therefore, theoretically incorporation of the gloss data should improve the translation systems. Specifically, I propose the following hypothesis:

- (1) **Gloss-helps-hypothesis: the translation systems trained with the gloss data incorporated should outperform the systems trained with only Gaelic and English sentences pairs (i.e. without gloss data).**

The hypothesis can have two versions, strong and weak:

- a. Strong version: Gloss may replace the source natural language totally, and the system outperforms the system trained with source natural language to target language sentence pairs (i.e. the baseline systems).
- b. Weak version: Gloss only increases the performance of the baseline systems, but cannot replace the source language.

The experiments reveal that replacing Gaelic words with glosses doesn't boost up the performance of the translation systems. Thus, the strong version (replacing-Gaelic-with gloss) of the gloss-helps-hypothesis is not attested. However, it is found that if the Gaelic data and the gloss data are combined in a specific way as the training data, the performance of the systems is improved significantly.

This chapter describes the experiments conducted to test the gloss-helps-hypothesis and the results attest the weak version. The rest of the chapter is organized as follows: Section 2 describes the constant parameter settings across all the experiments, section 3 tests the hypothesis in (1a), section 4 tests the hypothesis in (1b), and section 5 concludes the chapter.

2 Technical Settings of the Machine Translation Experiments

The experiments are conducted by using OpenNMT (Klein et al., 2017), which implements the state-of-the-art neural net machine translation algorithms (Cho et al., 2014a,b; Bahdanau et al., 2014). The following default parameters settings of OpenNMT are used across all models so that the only independent variable is the type of the training data:

- Word vector size: 500
- Type of recurrent cell: Long Short Term Memory
- Number of recurrent layers of the encoder and decoder: 2
- Number of epochs: 13
- Size of mini batches: 64

The data and the scripts will be accessible on GitHub¹, so that the results can be reproduced.

¹https://github.com/lucien0410/Scottish_Gaelic

3 Gloss Representation Solely Does NOT Outperform Gaelic Sentences

This section tests the strong version of Gloss-helps-hypothesis in (1a). Given the assumption that gloss may be better than any natural language in terms of representing meanings, it is expected that for neural net machine translation systems it is easier to learn how to translate from the glosses of Scottish Gaelic to English than to learn how to translate from Scottish Gaelic to English. However, the results show that there is no significance difference between the two types of data (i.e. GLOSS \rightarrow English and Gaelic \rightarrow English).

3.1 Procedure of the Experiments

I use repeated random sub-sampling validation to compare the performances of the two type of models.

Totally we have 8,388 indexed 3-tuples of Gaelic sentence, a gloss line and an English translation. In the interlinear glossed text example below, each line is an argument of a 3-tuple sample.

- (2) Tha a athair nas sine na a mhàthair.
 be.pres 3sm.poss father comp old.cmpr comp 3sm.poss mother
 ‘His father is older than his mother.’

The 3-tuple above is:

- (3) <“Tha a athair nas sine na a mhàthair.”, “be.pres 3sm.poss father comp old.cmpr comp 3sm.poss mother”, “His father is older than his mother.”>

First, the samples (i.e. the 3-tuples) are randomly split into three datasets: training set (N=6,388), validation set (N=1,000), and test set (N=1,000).

- (4) Definitions of datasets:

Let:

- a. $Index_{Train}$, $Index_{Validation}$, and $Index_{Test}$ be sets of random indexes from 0 to 8,387.
- b. $Index_{Train} \cap Index_{Validation} \cap Index_{Test} = \emptyset$
- c. $|Index_{Train}| = 6,388$; $|Index_{Validation}| = 1,000$; $|Index_{Test}| = 1,000$.

The step above just randomly splits the indexes of the 3-tuples into three distinct sets: $Index_{Train}$, $Index_{Validation}$, and $Index_{Test}$. Based on the indexes, we generate the sets of samples. For each index, the 3-tuple is split into two pairs: <gloss, English>, <Gaelic, English>, so that later we can compare the different effects of gloss lines and Gaelic sentences. For each pair, the first item is the source sequence, and the second item is the target sequence. The systems learns how to map the source sequence to the target sequence.

- (5) Gloss to English

- a. $GLOSToEN_{Train} = \{ \langle gloss_i, En_i \rangle \mid i \in Index_{Train} \}$
- b. $GLOSToEN_{Validation} = \{ \langle gloss_i, En_i \rangle \mid i \in Index_{Validation} \}$

$$c. \text{ GLOSStoEN}_{\text{Test}} = \{ \langle gloss_i, En_i \rangle \mid i \in Index_{\text{Test}} \}$$

(6) Gaelic to English

$$a. \text{ GDtoEN}_{\text{Train}} = \{ \langle GD_i, En_i \rangle \mid i \in Index_{\text{Train}} \}$$

$$b. \text{ GDtoEN}_{\text{Validation}} = \{ \langle GD_i, En_i \rangle \mid i \in Index_{\text{Validation}} \}$$

$$c. \text{ GDtoEN}_{\text{Test}} = \{ \langle GD_i, En_i \rangle \mid i \in Index_{\text{Test}} \}$$

The models are trained with the training set and validation set (i.e. the model learns how to map the source sequence to the target sequence). Both training set and validation set are known information for the models². Specifically, the neural net system learns how to maps gloss lines to English sentences from samples in (5a) and (5b), and another neural net system learns how to maps Gaelic sentences to English sentences from from samples in (6a) and (6b).

(7) Models:

- a. $\text{Model}_{\text{GLOSStoEN}}$ = Model trained with $\text{GLOSStoEN}_{\text{Train}}$ in (5a) and $\text{GLOSStoEN}_{\text{Validation}}$ in (5b)
- b. $\text{Model}_{\text{GDtoEN}}$ = Model trained with $\text{GDtoEN}_{\text{Train}}$ in (6a) and $\text{GDtoEN}_{\text{Validation}}$ in (6b)

The two trained models (gloss-to-English and Gaelic-to-English) then take the right source sequences of the test sets (i.e. glossing lines and Gaelic sentences for $\text{Model}_{\text{GLOSStoEN}}$ and $\text{Model}_{\text{GDtoEN}}$ respectively) as inputs and then generate the predicted target sequences (i.e. English sentences).

(8) Predictions:

- a. $\text{Predictions}_{\text{GLOSStoEN}}$ = A list of English sequences that $\text{Model}_{\text{GLOSStoEN}}$ maps to from the gloss sequences in (5c)
- b. $\text{Predictions}_{\text{GDtoEN}}$ = A list of English sequences that $\text{Model}_{\text{GDtoEN}}$ maps to from the Gaelic sentences in (6c)

To evaluate the model, the predicted target sequences are checked against the target sequences of the test set (i.e. the gold standard/human-translated English sentences). Specifically, the BLEU (bilingual evaluation understudy)³ score metric (Papineni et al., 2002) of each prediction is calculated using the `multi-bleu.perl`⁴ script, a public implementation of Moses (Koehn et al., 2007). The BLEU score calculation is an automatic evaluation of how similar two

²Technically speaking, the validation set is part of the training data in terms of machine learning. The presence of the validation set is a special requirement of neural net machine learning, which uses the validation set to further tune the parameters learned from the training set.

³There are other automatic machine translation evaluation algorithms available, such as translation edit rate (Snover et al., 2006) and Damerau–Levenshtein distance (Damerau, 1964; Levenshtein, 1966). BLEU is chosen for the current experiments because it is the most widely used evaluation algorithm, and the correlation between the BLEU score evaluation and human judgment evaluation is also well-acknowledged.

⁴The script can be downloaded from: <https://github.com/moses-smt/mosesdecoder/blob/master/scripts/generic/multi-bleu.perl>

copora are. In the current experiments we are comparing the predicted target sequences with the gold standard. The BLEU score of 100 means the two copora are identical, and the BLEU score of 0 means the two copora are completely distinct from each other.

$$(9) \quad \text{Gold-Standard} = \text{English sentences in (5c)} = \text{English sentences in (6c)}$$

Note that the gold-standard is the same because they are the same English sentences in the 3-tuples samples. Then the two sets of predicted English sentences are evaluated, yielding two BLEU scores.

(10) Scores:

$$\text{a. } \text{Score}_{\text{GLOStoEN}} = \text{BLEU}(\text{Gold-Standard}, \text{Predictions}_{\text{GLOStoEN}})$$

$$\text{b. } \text{Score}_{\text{GDtoEN}} = \text{BLEU}(\text{Gold-Standard}, \text{Predictions}_{\text{GDtoEN}})$$

This procedure of splitting the data into three sub-sets, training the models, and evaluating the models is executed for ten times.

3.2 Result

After ten rounds of repeated random sub-sampling validation, ten pairs of scores of the two models are generated, as shown in the following table.

Round	GLOSS	Gaelic
1	18.39	17.29
2	18.00	16.42
3	16.02	15.29
4	20.22	15.97
5	19.02	17.79
6	15.53	16.73
7	18.00	17.11
8	20.08	16.37
9	15.82	15.93
10	15.93	16.99
Mean	17.70	16.59

Table 1: BLEU scores of Model_{GLOStoEN} and Model_{GDtoEN}

The average score of the Models_{GLOStoEN} is only sightly higher than the average score of the Models_{GDtoEN}. Also, after doing a paired T-test, the difference between the two types of models is not attested ($M_{\text{GlossToEn}}=17.70$, $SD_{\text{GlossToEn}}=1.78$, $M_{\text{GDToEn}}=16.59$, $SD_{\text{GDToEn}}=0.74$); $t(9)=1.97$, $p=0.08$).

3.3 Summary

The ultimate practical goal of the dissertation is to use glossing data to develop better machine translation systems. Here *better* means to be better than a

baseline system, which is the machine translation system trained with Gaelic-to-English translation samples. The models in (7b) are the baseline systems, and their scores are in the Gaelic column of table (1). These are the target scores that we aim to outperform. The experiment above is the first attempt to improve that scores by using the *gloss treatment*, in which the Gaelic sentences are replaced with gloss lines. However, the result shows that this *gloss treatment* is not effective as the scores of the gloss models are not statistically higher than the baseline Gaelic-to-English models.

3.4 Discussion

It is assumed that the performances of the machine translation systems are correlated with the quality of the representation of meanings in the source sequences. Better representations of meanings yield better machine translation systems. Given the results in (3.2) that the two types of models have the same performance, it is concluded that glosses and natural languages are about the same in terms of representing meanings. The strong version of the Gloss-helps-hypothesis does not hold.

We may now combine Gaelic and Gloss sentences as the training data to test the weak version of the Gloss-helps-hypothesis. The experiments and results are reported in the next section.

4 Section: Combining Gaelic words with Glosses

Round	GLOSS	Gaelic	ConcatGLOSSGaelic	GLOSS.Gaelic	hyGD	hyGW
0	18.39	17.29	15.42	13.67	9.44	15.95
1	18.00	16.42	14.31	12.49	9.07	15.60
2	16.02	15.29	15.38	11.01	7.69	14.15
3	20.22	15.97	14.18	12.33	9.12	14.72
4	19.02	17.79	18.63	12.56	9.08	15.74
5	15.53	16.73	14.89	12.13	10.45	14.88
6	18.00	17.11	15.16	11.55	8.62	14.45
7	20.08	16.37	15.20	12.78	10.00	16.41
8	15.82	15.93	15.50	12.43	10.52	15.15
9	15.93	16.99	15.72	11.65	8.46	17.61
Mean	17.70	16.59	15.44	12.26	9.24	10.24

Table 2: More BLEU scores of the Models, where ‘ConcatGLOSSGaelic’ is ‘GLOSS_Gaelic’, ‘GLOSS.Gaelic’ is ‘GLOSS Gaelic’, ‘hyGD’ is using the less ambiguous one from either Gaelic word or gloss with Gaelic word as the default, ‘hyGW’ is the same as ‘hyGD’ but with the default as gloss.

References

Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio (2014), “Neural machine translation by jointly learning to align and translate.” *arXiv preprint arXiv:1409.0473*.

Round	Para	ParaPart	ParaPartHalf	ParaPartHalfOver	Over.Gaelic	google
0	25.42	32.64	24.07	26.31	29.05	22.09
1	25.32	32.28	18.58	24.85	28.61	25.38
2	20.72	29.94	18.00	22.96	23.78	23.72
3	22.22	31.18	22.25	25.48	27.50	23.21
4	24.27	32.83	23.79	25.33	25.51	22.31
5	24.55	31.11	21.40	24.39	27.88	23.41
6	27.03	32.19	23.61	26.29	25.72	24.53
7	25.34	33.52	23.73	25.61	27.12	22.78
8	24.24	30.93	23.16	25.59	25.20	25.67
9	25.96	34.35	24.49	26.32	26.39	23.42
Mean	24.51	32.10	22.31	25.31	26.68	23.65

Table 3: ‘Para’ is ‘GLOSS -> Gaelic; Gloss -> En; Gaelic -> En’, ‘ParaPart’ is ‘Para plus Gaelic word token -> Gloss token’, ‘Over’ means oversampling, ‘Half’ means using half of the training data.

Cho, Kyunghyun, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio (2014a), “On the properties of neural machine translation: Encoder-decoder approaches.” *arXiv preprint arXiv:1409.1259*.

Cho, Kyunghyun, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio (2014b), “Learning phrase representations using rnn encoder-decoder for statistical machine translation.” *arXiv preprint arXiv:1406.1078*.

Damerau, Fred J (1964), “A technique for computer detection and correction of spelling errors.” *Communications of the ACM*, 7, 171–176.

Klein, Guillaume, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander M. Rush (2017), “Opennmt: Open-source toolkit for neural machine translation.” *CoRR*, abs/1701.02810, URL <http://arxiv.org/abs/1701.02810>.

Koehn, Philipp, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, et al. (2007), “Moses: Open source toolkit for statistical machine translation.” In *Proceedings of the 45th annual meeting of the ACL on interactive poster and demonstration sessions*, 177–180, Association for Computational Linguistics.

Kotsiantis, Sotiris B, I Zaharakis, and P Pintelas (2007), “Supervised machine learning: A review of classification techniques.” *Emerging artificial intelligence applications in computer engineering*, 160, 3–24.

Levenshtein, Vladimir I (1966), “Binary codes capable of correcting deletions, insertions, and reversals.” In *Soviet physics doklady*, volume 10, 707–710.

Papineni, Kishore, Salim Roukos, Todd Ward, and Wei-Jing Zhu (2002), “Bleu: a method for automatic evaluation of machine translation.” In *Proceedings of the 40th annual meeting on association for computational linguistics*, 311–318, Association for Computational Linguistics.

Snover, Matthew, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul (2006), “A study of translation edit rate with targeted human annotation.” In *In Proceedings of Association for Machine Translation in the Americas*, 223–231.