# COMP4801 Interim Report

## Machine translation models for Cantonese-English pair

Liu Hey Wing (3035474016)

Supervisor: Dr. Yip Beta

24 January, 2021

# Abstract

Although the studies on Cantonese language are conducted extensively and Machine Translation (MT) had been developed in the past decades, the study on Cantonese-based Machine Translations are under researched. This project aims to provide a foundation and preliminary research in this field, and investigate the difficulties and explanation of the absence. We will implement different MT models, including historically and up-to-date technologies, evaluate the results and post the trained models online for future development.

Data collection of bilingual corpus is finished and MT models are successfully installed and tested. 2 of the models had completed the pilot run. Discussion on results and evaluation the results are discussed including limitation and difficulties we encountered. In the following months, we will keep on improving data, train the remaining MT models and evaluate their results.

# Acknowledgement

I would like to thank you Dr. Yip Beta, the supervisor of this project, for his generous help for obtaining different resources and advice on this project. As a student having double major in Computer Science and Linguistics, the opportunity of working on natural language processing and computational linguistics is very precious and cherishable, to which I am very thankful.

Secondly, I would like to thank you the CAES English lecturer, Grace Chang. She had provided sufficient support on my final year project and outlined detailed improvement for my presentation and report. This report would not be done without her.

# Contents

# List of Figures

# List of Tables

# Abbreviations

**BLEU**  Bilingual Evaluation Understudy. 6, 15, 20–22, 25

**EBMT**  Example-based machine translation. 4–6, 9, 10

**FOS**  free/open source. 8–11, 24

**GRU**  Gated Recurrent Unit. vi, 12, 15, 17, 18, 20, 23, 24

**LM**  Language Model. 10, 11

**LSTM**  Long Short Term Memory. 6, 12

**MSC**  Modern Standard Chinese. 5, 23

**MT**  Machine Translation. 1–10, 12, 14–17, 20–26

**NMT**  Neural machine translation. 4, 6, 12, 18, 21, 22

**PBSMT**  Phrase-based Statistical Machine Translation. 6, 11, 15, 23

**RBMT**  Rule-based machine translation. 4–6, 9, 15, 23

**RNN**  Recurrent Neutral Network. 12, 13

**SL**  Source Language. 2, 9, 11, 23

**SMT**  Statistical Machine Translation. 4, 6, 9–11

**TL**  Target Language. 2, 9, 10, 23

**WMT**  The Confrence on Machine Translation. 6, 12, 21

# 1  Introduction

## 1.1  Background

Machine Translation (MT) has been developing since 1950s [9, Chapter 1.3], and started using in practice in recent years. However, numerous limitations has restricted the usage and functionalities of MT systems, such as translation direction (uni- or bi- directional), number of languages (multi- or bi- lingual), limited linguistic resources (unpopular languages), computational force and evaluating method (automatic or human).

An abstract concept of MT consist of 2 parts, 'declarative' and 'procedural' information, which can be simply understood as data and algorithms [9]. These information construct the basis of an MT model, and will be elaborated in the following content.

### 1.1.1 Linguistic aspect

Linguistic information are the 'declarative' in MT systems, which are based on facts, definitions and relationships of natural languages. These information may include grammatical description (dependency or generative), a corpus of orthography[1] or corpus of SL and TL.

Cantonese is a language ranked 18th in the world by the number of speakers [10]. The study of Cantonese linguistics are traced back to 1900s with Ball [11], since then, linguists had provided different aspects about Cantonese extensively and continuously, such as phonology, grammar [12] [13] and sociologist studies. However, there is a lack of a unified/standard writing system had created issues of choice in orthography. The written form is unified through internet and forums in recent years [14].

Bernard Vauquois' pyramid (Figure 1) is a very famous framework for preliminary MT, which focus on how linguistic information can be integrated. The pyramid is the ideological model for translation models.



Figure 1: Bernard Vauquois' pyramid [1].

Different levels and layers of translation are described, such as direct and syntactic transfer approaches. The concept is to decompose text into smaller units, such as morphemes, and reconstruct them into another language.

Some linguistic challenges are described in Figure 2 considering different aspects, including orthography, morphology, lexis etc. These problems had struggled with computational linguist for years, and some are still hurdles in the development of MT.

---

[1]Not all languages in the world uses Latin alphabet, and some are lacking a unified or standard writing system like Cantonese.

2

| Linguistic level | Challenge | Main related works |
|---|---|---|
| ORTHOGRAPHY | Spelling | Bertoldi et al. [2010], Farrús et al. [2011] |
| | Truecasing/Capitalization | Lita et al. [2003], Wang et al. [2006] |
| | Normalization | Riesa et al. [2006], Aw et al. [2006], Diab et al. [2007], Kobus et al. [2008] |
| | Tokenization | Farrús et al. [2011], El Kholy and Habash [2012] |
| | Transliteration | Boas [2002], Virga and Khudanpur [2003], Kondrak et al. [2003], Zhang et al. [2004], Kondrak [2005], Mulloni and Pekar [2006], Kumaran and Kellner [2007], Mitkov et al. [2007], Istvan and Shoichi [2009], Nakov and Ng [2009] |
| MORPHOLOGY | Inflections | Brants [2000], Ueffing and Ney [2003], Creutz and Lagus [2005], Minkov et al. [2007], Koehn and Hoang [2007], Virpioja et al. [2007], Avramidis and Koehn [2008], de Gispert et al. [2009] El-Kahlout and Oflazer [2010], Bojar and Tamchyna [2011], Green and DeNero [2012], Formiga et al. [2012], Rosa et al. [2012] |
| LEXIS | Unknown words | Knight and Graehl [1998], Al-Onaizan and Knight [2002], Koehn and Knight [2003], Fung and Cheung [2004], Shao and Ng [2004], Langlais and Patry [2007], Mirkin et al. [2009], Marton et al. [2009], Li et al. [2010], Huang et al. [2011], Zhang et al. [2012] |
| | Spurious words | Fraser and Marcu [2007], Li and Yarowsky [2008], Menezes and Quirk [2008] |
| SYNTAX | Word reordering | Wu [1997], Alshawi et al. [2000], Menezes and Richardson [2001], Yamada and Knight [2002], Aue et al. [2004], Galley et al. [2004], Ringger et al. [2004], Xia and McCord [2004], Chiang [2005], Collins et al. [2005], Ding and Palmer [2005], Quirk et al. [2005], Simard et al. [2005], Zhang and Gildea [2005], Galley et al. [2006], Liu et al. [2006], Huang et al. [2006], Langlais and Gotti [2006], Smith and Eisner [2006], Turian et al. [2006], Birch et al. [2007], Li et al. [2007], Zhang et al. [2007], Wang et al. [2007], Cowan [2008], Elming [2008], Graehl et al. [2008], Li and Yarowsky [2008], Badr et al. [2009], Genzel [2010], Shen et al. [2010], Khalilov and Fonollosa [2011], Bach [2012], Germann [2012] |
| SEMANTICS | Sense disambiguation | García-Varea et al. [2001], Chiang [2005], Bangalore et al. [2007], Carpuat and Wu [2007], Chan et al. [2007], Carpuat and Wu [2008], Shen et al. [2009], Wu and Fung [2009], España-Bonet et al. [2009], Haque [2011], Banchs and Costa-jussà [2011], Banarescu et al. [2013] |

Figure 2: Some linguistic challenges in MT and related work [2]

Linguistics considerations and factors in an MT model are overwhelming, where each new language pair requires professional linguists to describe particular features and construct relevant data. With all the challenges described, translation frameworks and models switched from knowledge-based approach to corpus-based approach.

### 1.1.2 Computational aspect

Algorithms, softwares and models are the 'procedural' information in MT, which means how the system interpret the 'declarative' information, represented as a form of program or procedures. There are different approaches for performing MT, each requires different 'declarative' language information and perform different actions [15].

Figure 3 shows the development and variants of MT since 1960s [16]. The inheriting tree structure is describing the relationship between different approaches.



Figure 3: MT approaches and models. Based on [3]

The first MT system introduced was Rule-based machine translation (RBMT). Example-based machine translation (EBMT) was introduced in the 1980s which focus on comparison and extraction [17]. In 1990s, IBM had performed a series of experiment for a new Machine Translation (MT) system based on sentence alignment which is later known as Statistical Machine Translation (SMT). Neural machine translation (NMT), which is a kind of deep learning, started around mid-2010s, together with the rise of AI and machine learning, and the relevant research is still ongoing.

From the history of the development of MT, the approaches switched from relying linguistic knowledge to statistics, data science and AI. In this project, several historical and advanced MT models are implemented, and the details will be discussed.

## 1.2   Previous works on MTs in Cantonese-English pair

To my knowledge, PoluU-MT-99 is the first MT system implemented for Cantonese-English pair in 1999. Both RBMT and EBMT approaches are included in the system, which described the details such as Cantonese segmentation, bilingual algorithm and target construction algorithm [18].

Later in 2005, one of the author of PoluU-MT-99, Yan Wu, presented another MT called LangCompMT05 [4], where the architecture (Figure 4) is nearly identical to PoluU-MT-99 with extended algorithms.



Figure 4: Structure of LangCompMT05 [4].

However, the corpus data used in both PoluU-MT-99[2] and LangComp05[3] is not reliable. Some non-Cantonese words, sentences, aspect markers are used in the study, such as 有些, 的, 正在 originated from Modern Standard Chinese (MSC)[4]. Thus, we have certain reservation towards the result of this study and leave it out of our consideration.

Unfortunately, the source code is not publicly available but the algorithms and structure is presented in the study, thus we are unable to test the model. Re-implementing the models from scratch is over-complicated and out of the scope of this project.

---

[2]The source of its example base is not described, but the testing data was retrieved from Mingpao, a Hong Kong newpaper uses MSC rather than written form of Cantonese in their print.

[3]The example based was constructed by Shiwen Yu, but his/her corpus is not cited in [4]. The source of testing data was created by the authors, which are targeted to perform specific sentence structure translation rather than real data.

[4]Or written form of Mandarin. Note that although Mandarin and Cantonese are both considered as the same language under the name of 'Chinese' and share some degree of written intelligibility of Hanzi and cognates, they are mutually unintelligible and are considered as different languages linguistically.

## 1.3 Objectives

Despite the abundant resources, barely any MT systems include Cantonese, with only Microsoft Bing translator and Baidu App. The lack of Cantonese MT systems (regardless of direction) is unusual. The purpose of this project is to investigate and discover such obstacles, provide possibly improvement for existing MT systems, and outline the future development for Cantonese MT.

## 1.4 Scope and Deliverable

MT is evolving rapidly, with numerous proposed variant and commercial implementation. Some of them are outdated and some have outstanding performance, so our scope is to narrow down and focus on some categories of MT systems listed above, including cutting-edge and historically popular models.

In The Confrence on Machine Translation (WMT) 2019 [19], the most advanced and newly developed MT were presented, which are mainly SMT (PBSMT) and NMT (LSTM and Transformer). Some historically popular MT will be discussed in this project, including RBMT and EBMT. Different MT systems will be implemented with training and testing set of data. We will evaluate the results and discover the reason for the absence of Cantonese MT.

At the end of the project, we want to deliver at least 1 MT model with a BLEU score of 10. With such system, we hope that it will increase the interest in research of Cantonese-based MT and potentially develop into a fully functional translation system.

## 1.5 Outline

In Introduction (section 1), the overall background is provided with objectives and scope. In Methodology (section 2), we will investigate MT systems from data, algorithm and evaluation perspectives. Details of the implementation will be discussed in Experimental Setup(Section 3). The preliminary results and discussion will be presented in (Section 4). Finally, the limitations and future work will be discussed in Project status (Section 5).

# 2 Methodology

In this section, the desciptions of each MT models are listed out, including data, algorithms and evaluation. Particular MT models are chosen, which cover historical models and advanced technologies.

## 2.1 Data

### 2.1.1 Linguistics choice

There are around 15,000 frequently used vocabularies in Hong Kong, whereas 9706 vocabularies are considered to be sufficient for primary school children, and those vocabularies are composited by 3171 Chinese characters [20]. The number of frequently used characters and vocabularies are comparable, thus, both formats are implemented to compare performances in MT models.

Chinese Character is chosen to represent Cantonese due to the high number of lexicons/cognates shared [12]. The meaning of characters are related in different vocabularies. Using romanized text would cause difficulties in data preprocessing also.

Considering the vocabularies without consensus on Chinese character and especially the polysemous, we constructed principles based on the finding in [14] to transformer them, (1) preferably choose the mostly recognized Chinese character, (2) choose the most disambiguated Chinese character, (3) use the most widely adopted romanized form and (4) use Jyutping to transcribe[5]. All forms are chosen based on the usage in Hong Kong, but not the forms in Mainland China or overseas Cantonese community.

### 2.1.2 Source

Bilingual corpus are preferably used and sometimes required as training and testing data in MT. To my knowledge, the only available Cantonese-English bilingual corpus is ShefCE [21]. Despite that the corpus exist, the transcription only includes Jyutping romanization without tone markers[6], which makes it not applicable for MT.

---

[5]For instance, '係' is the mostly recognized form to represent *haai2* 'at', but it is overlapped with *haai6* 'to be', so '喺' is chosen instead. Another example is the expressions 'hea' and 'chur' without Chinese character, and they remains this form to be recognizable.

[6]The study used a refined lexicon derived by linguistic knowledge and sub-syllable unit combinatorial constraints, and get rid of the lexical tones [22]. However, in the precise environment like translation, the

Cantonese consists of high frequency of code-switching in daily speech or text. While Audio originated corpus encourage code-switched text, translation originated corpus discourage the phenomenon. This is an important language phenomenon in Cantonese used in Hong Kong, and possibly contribute to the absence of Cantonese-based MT, so both audio and translation originated data are included.

TED [23] is a nonprofit organization which invite speakers to give talks, namely TED talks. It is open-source for translation and subtitles in more than 100 languages, including Cantonese and English. Following web crawling/scraping method proposed in [24] with Beautiful Soup [25], a total of 193 talks are available to construct the corpus. The timestamp is used to subdivide the whole talk into sentences, however, some subtitles are not divided into pieces due to formatting issue.

Netflix movies and TV shows is chosen as another data source, where they support both Cantonese and English audios and subtitles. There are 300 videos with Cantonese audio and English subtitle, while there are only 6 videos vice versa. For fairness between the language pair, it is better to include both direction. However, the Cantonese subtitles are encrypted and incomprehensible. Therefore, pyTranscriber [26], a FOS software, performs Cantonese transcription based on Google Speech-to-Text API [27].

The corpus is not exactly parallel, because of the sentence structure between Cantonese and English and the difficulties in aligning sentences. The noisiness issue will be addressed again in Section 4.2.

### 2.1.3 Preprocessing

The data collected are unified into the same format described above. Some of the data in TED talks are in simplified character and Mainland version of the controversial forms, which are further transformed into our standardized format.

Punctuation marks[7] are removed because full text translation and long sentence translation is not the focus in this project [28]. Word order are relatively consistent in

---

disappearance of tones would cause loss in information.

[7]In addition, the output of pyTranscriber do not include any punctuation marks, and the movie subtitles contain punctuation marks seldomly.

English, which reduce ambiguities in context. While word order in Cantonese varies with right-dislocation [12], but still understandable without punctuation marks.

## 2.2 Description of MT models

### 2.2.1 RBMT

Rule-based machine translation (RBMT) is designed based on linguistic knowledge, in which decompose the Source Language (SL) into parse tree, restructure the lexicons, and reformat into the Target Language (TL). This was a popular model before 1980s, but it had faded out in use because of the improvement of computational power, which later give rise to SMT and other systems.

Apertium [5] is a free/open source (FOS) platform and the most popular framework for RBMT. The structure is shown below (Figure 5).



Figure 5: Structure of Apertium [5].

We adopt Apertium as an example of RBMT to perform our training. It is heavily based on parallel (Bilingual) data, dictionary and explicit linguistic data such as grammars and structural transfer rules.

### 2.2.2 EBMT

Example-based machine translation (EBMT) is the first MT system get rid of the 'knowledge' of linguistics and based on 'pure' data, which was developed by Makoto Nagao [17] in 1984. The input sentence would be searched in the trained model, substitute the phrase or character, and generate the output sentence, as shown in example (1)-(2) [29].

(1) (a) He buys a notebook.

Kare wa noto o kau.

(b) I read a book on international politics.

Watashi wa kokusai seiji nitsuite kakareta hon o yomu.

(2)   (a)  He buys a book on international politics.

       (b)  Kare wa kokusai seiji nitsuite kakareta hon o kau.

On the other hand, Somers [30] had a thorough review about EBMT and addressed several limitations and difficulties[8] on further development and improvement. As a result, it introduced the idea of corpus-based MT and give rise of SMT.

CMU-EMBT [6] in Figure 6 uses a nearly 'pure' EMBT, with optional components and integration to other MT models available.



Figure 6: Structure of CMU-EBMT [6].

Input sentences are converted into lattice where character/phrase segmentation would be performed, and convert into the Language Model (LM). The decoder would perform lexicon lookup for decomposed TL sentence in the corpus, and substitute the sentence. The model is already not functional and cannot be implemented in this project, the structure is provided to increase the understanding of EBMT.

The rise and fall of EBMT was fast, which induced low effort or resources[9] on the topic. Thus, EBMT will not be used in this project.

---

[8]Such as the size of example, lack of parallel corpora of language pairs, suitability of examples etc.

[9]Some more known EBMT FOS platforms are already unavailable now, such as Cunei [31], OpenMa-TrEx [32] and uploaded in 2010s. Even if they are available, the coding is not compatible with modern technology anymore.

### 2.2.3 PBSMT

Statistical Machine Translation (SMT) uses a probabilistic approach to perform translations, which involves much less linguistic knowledge. The most popular variant is Phrase-based Statistical Machine Translation (PBSMT).

Bayes Decision Rule (Figure 7) describes $Pr(e_1^I|f_1^F)$ as the translation model (SMT), where the source sentence $f_1^J = f^1, ..., f^j, ..., f^J$ is translated into target sentence $e_1^I = e^1, ..., e^i, ..., e^I$.



Figure 7: Structure of Bayes Decision Rule applying
Source: [33]

PBSMT segment SL into phrases and translate phrase-to-phrase with the probability, and the mathematical description is below [33].

$$\hat{e}_1^I \approx \underset{e_i^I, B}{argmax} \left\{ \prod_{i=1}^{I} p(e_i|e_{i-1}) \cdot \prod_{k=1}^{K} p(\tilde{f}_k|\tilde{e}_k)^\lambda \right\}$$

We will implement Moses [34] in our project, one of the most popular FOS toolkit for SMT, consisting 2 main components, *training pipeline* (LM) and *decoder* which obeys Bayes' rule.

### 2.2.4 GRU

Neural machine translation (NMT) uses 2 approaches, attention mechanism (section 2.2.5) and Recurrent Neutral Network (RNN) [35]. NMT has been proven with significant improvement comparing with other models, especially in WMT'16 [36] in different fields. The limitation for NMT is the ambiguity within the hidden layers, which is not understandable by human and thus increase the difficulty for further improvement [37].

RNN approach consists of numerous layers of feed-forward computation and RNN between the encoder-decoder and send results backwards. The encoder transforms $x$ into a vector space of hidden layers[10], and the decoder transforms those vectors into $y$, using the likelyhood equation below.

$$p(y|x;\theta) = \prod_{j=1}^{m+1} p(y_j|y_{0:j-1}, x; \theta)$$

LSTM cell is an extension of RNN proposed back in 1997 [38], and implemented for MT in 2014 [39]. Cho et al. [40] proposed a modified version called Gated Recurrent Unit (GRU) as in Figure 8 with higher computational speed and low memory space.



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$
$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$
$$\tilde{h}_t = tanh(W \cdot [r_t * h_{t-1}, x_t])$$
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Figure 8: GRU cell [7].

GRU combines the input and forget gate into the *update* gate $z_t$, and introduced the *reset* gate $r_t$, but still unable to perform multi-threading and have low training speed.

The pyTorch framework provided by D2L [41] is adopted for GRU model.

---

[10]Hidden Markov model (HMM) is a complicated topic, which is out of the scope of this project

12

### 2.2.5 Transformer

In Attention mechanism, each segment has a vector of non-fixed size in contrast to RNN, which represent the attention driven in a sentence. The decoder then derives output once at a time based on the attention towards the input, and passes the results backwards . This avoids RNN's problem of mapping sentences into fixed-length vector, which cause lost in input information and over-matching [8] [42] [43].

Transformer model was invented in 2017 [42], with the innovation of multi-head attention in the encoder. Figure 9 presents the structure of Transformer model.

Figure 9: Structure of Transformer model [8].

The left side is the encoding of input sentence, which is a attention block generating vector based on sentence length. The decoder on the right side are separated into 2 blocks, where the first attention block encodes the output sentence, and the second attention block computes the attention between the input vector and output vector. The details of attention blocks and multi-layer is as below (Figure 10), where $Q$ is the attention ma-

trices, $K$ is the key index and $V$ is the value.



Figure 10: (left) Scaled Dot-Product Attention. (right) Multi-Head Attention. [8]

Multi-head attention uses self-attention mechanism, which increase the quality of the system. In addition, the structure allows multi-threading in training on several GPU, and highly increase the training speed. Finally, The result would be converted into probability functions to predict output.

The Transformer model in Tensorflow [44] with the configurations and library provided by Vaswasni [45] will be implemented in this project.

## 2.3 Evaluation of translation output

Each time the MT translate a sentence, it is curious about how it performed, whether it translate correctly and how can we improve its result. The performance of MT systems are the key for further improvement, and evaluate methods will be discussed.

### 2.3.1 Manual evaluation

Manual evaluation is always considered time consuming and high cost, and the consistency in a large-sclae rating is difficult to maintain [46]. In current advancement in technology, machine can evaluate translation results in hours where human take days or weeks to complete, thus, manual evaluation are not preferred in most circumstances.

[47] had provided another perspective using Amazon's Mechanical Turk to recruit human evaluators for numerous languages, which showed evaluating result are actually near to experts. Considering the status of Cantonese [10], the evaluation is feasible, however, the cost and time are not preferred in this project.

### 2.3.2   Automatic evaluation

MT systems often uses automated metrics to evaluate translation results. Different metrics are proposed for better performance.

Bilingual Evaluation Understudy (BLEU) [48] is widely adopted to evaluate translation results [1, 2, 43, 19, 35, 33, 49, 36, 39, 50, 42, 8, 4], which is well-established and preforms quality evaluation close to human. Since then, more metrics are proposed according to different dependencies, such as NIST, RED and WER. In this project, BLEU is adopted to conduct comparison with previous studies.

## 2.4   Summary

To conclude, our project will implement the softwares / models as below.

- *Data*: TED talks + Netflix in Chinese character
- *MT models*:
    - *RBMT*: Apertium
    - *PBSMT*: Moses
    - *GRU*: Tensorflow with GRU cell
    - *Transformer*: Tensorflow with Transformer model
- *Evaluation*: BLEU

# 3 Experimental Setup

In this section, We will describe the experimental setup, including the details of generated corpus and implemented MT models.

## 3.1 Current progress

The progress is keeping to the schedule. However, the workload of bilingual corpus is higher than expected, so we will estimate and propose to put more effort on the data collection. Testing on MT models are required to assure quality and performance.

### 3.1.1 Data

The corpus is generated from TED talks and Netflix with over 9000 parallel sentences. All talks in TED are extracted, while only 2 movies from Netflix are included in the corpus. The transcription accuracy of pyTranscriber for Netflix is relatively low, due to musics and background noises. Thus, manual proofreading is required to preprocess the data. In Figure 11, some sample sentences from the corpus are listed out.

```
而且 解決 咗 呢個 問題   and solved it
呢個 係 我哋 嘅 醫療 保 體系 統 中        It s the one great preventive health success
一個 非常 成功 嘅 疾病 預防 嘅 例子      we have in our health care system
```

Figure 11: Example sentences from corpus.

The sentences are in the format proposed in Section 2.1.3. Some observation is that the sentence are not completely identically parallel due to the reasons mentioned in Section 2.1.2.

Overly-long sentences of over 40 units[11] in both languages are excluded in the data, as long-range attention and documental translation are not the focus in this project. The input data is split into fixed training and testing set, with a 80:20 ratio. Although the testing results would not be absolutely fair in different splittings, the computational cost and time are limited for other validation methods.

---

[11]Different subword encoder have different approach for dividing word units. In English, the morphology can be decomposed like '-ing' and '-ed' tense morphemes. In Cantonese, we adopt both vocabulary and character based approach, with the minimum subword unit is restrict into 1 Chinese character, as the input are parsed as UTF-8 coding.

### 3.1.2  MT

All 4 MT models mentioned above are implemented and tested for functionalities. GRU and Transformer model had completed the pilot run and have some preliminary results to be discussed below. The structure of GRU and Transformer model are shown in Figure 12 and Figure 13 respectively.

```
EncoderDecoder(
  (encoder): Seq2SeqEncoder(
    (embedding): Embedding(2862, 32)
    (rnn): GRU(32, 32, num_layers=2, dropout=0.1)
  )
  (decoder): Seq2SeqDecoder(
    (embedding): Embedding(3559, 32)
    (rnn): GRU(64, 32, num_layers=2, dropout=0.1)
    (dense): Linear(in_features=32, out_features=3559, bias=True)
  )
)
```

Figure 12: Structure summary of GRU model.

```
Model: "transformer_1"
_____
Layer (type)               Output Shape          Param #
=================================================================
encoder_2 (Encoder)        multiple              1790592
_____
decoder_2 (Decoder)        multiple              1489024
_____
dense_151 (Dense)          multiple              434085
=================================================================
Total params: 3,713,701
Trainable params: 3,713,701
Non-trainable params: 0
_____
```

Figure 13: Structure summary of Transformer model.

The figures are another representation of the models identical to structures discussed in Section 2.2. The Transformer model is too complicated to be shown as the form of summary, so the above format is used to summarize the structure. The term 'params' in Transformer is dependent on the size of data input, which is subject to change.

Considering the fact that many Cantonese speakers in Hong Kong, Macau and overseas are bilingual speakers, the demand for a Cantonese-English system is lower than vice versa. The ultimate goal for this project is to implement bi-directional MT system, but for current stage, only English-Cantonese direction is implemented.

## 3.2 Parameters

The hyperparameter of GRU and Transformer models are shown in Table 1.

| Transformer | GRU |
|---|---|
| num_layers = 4 | embed_size = 32 |
| d_model = 128 | num_hiddens = 32 |
| dff = 512 | num_layers = 2 |
| num_heads = 8 | dropout = 0.1 |

Table 1: Hyperparameter setting of Transformer and GRU

The default setting is used and fine-tuning is not preformed. Each model had trained on 100 epoches. These hyperparameter are subject to change based on numerous reasons, including the size of dataset, number of epochs, quality of data etc.

## 3.3 Training interface

The NMT models learn by evaluating themselves in the training process, where either the training accuracy or training loss are used as the evaluating metric. As mentioned in Section 2.1.1, both character and vocabulary based Cantonese input are used.

### 3.3.1 GRU

The loss function defined in GRU model is cross-entropy loss. Using the loss function, the GRU models adjusts the weighting and biases of artificial neurons. The training loss through epoches, i.e. the learning progress, is shown in Figure 14.



Figure 14: Training loss of GRU. (left) vocabulary-based. (right) character-based.

18

The vocabulary-based approach demonstrate a training curve with steep gradient. On the other hand, The training curve of character-based approach shows gentle improvement throughout epoches.

### 3.3.2 Transformer

The loss function used in Transformer model is cross-entropy loss. Using the loss function, the Transformer model learns to pay attention towards different embedding space. Along with training loss, the training accuracy is also recorded during the epoches using the accuracy function in Figure

Using the two functions, we can observe the learning curve of models and contribute to fine-tune. The training process is recorded throughout epoches (Figure 15).



Figure 15: Comparison of vocabulary and character base in Transformer. (left) Training loss. (right) Training accuracy.

All the curves increase or decrease gently. For training loss, different based input started differently and converge towards zero. For training accuracy, both base approach started at zero, but diverge to different result. The character-based input has accuracy higher than vocabulary-based approach.

# 4 Results and Discussion

## 4.1 Results

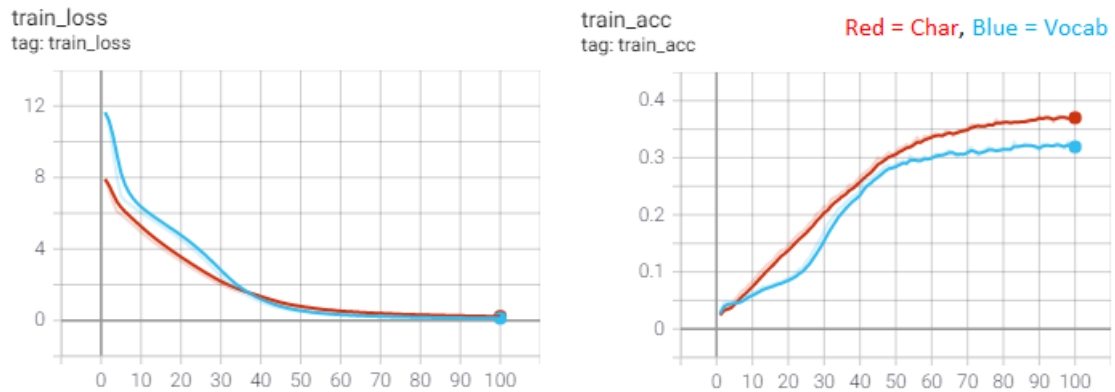The testing dataset is used to evaluate the performance of MT models towards unseen data. In Figure 16, an example test sentence of the multi-head attention in Transformer model is shown.



Figure 16: Sample output of multihead attention in Transformer.

The attention paid in all heads are messy and without pattern. The matching between input and output sentence is not learnt by the model, and there is no direct linkage shown in the different attention heads.

The BLEU scores of Transformer and GRU model of both vocabulary-based are shown in Table 2. The results of the original paper are listed out for comparison.

| Model | | Vocab | Character |
|---|---|---|---|
| Implemented models | Eng-Yue Transformer | 2.46 | 3.75 |
| | Eng-Yue GRU | 1.16 | 0 |
| Original Transformer | En-De | 28.4 | |
| | En-Fr | 41.8 | |

Table 2: BLEU scores of different models

The scores of the implemented models are much lower than the original transformer. The character-based approach performs better than vocabulary-based approach in Transformer. The BLEU scores of different approaches in GRU models is nearly zero and considered negligible.

## 4.2 Discussions

### 4.2.1 Overview

In the Transformer model, we can clearly visualise how the system maps input tokens into output tokens. Thus, attention graph similar to Figure 16 is the generalization of overall performance. Some output sentence have similar meaning comparing to the input sentence, but the overall accuracy is considerably low. BLEU score lower than 10 are considered as useless [51], therefore the implemented systems are not functional at all at the current moment. The results obtained show the distance to our goal in this project, and improvement is required in the future.

In various NMT competition like WMT, human translators still outperform the most advanced models, especially in English-Chinese pair. Another observation is that contextual information have higher translation quality, such as news translation task, laws and documentations.

There was outstanding performance in the MT system. Specifically, the output result was unexpected, as revealed by the mismatch between input sentence and targeted output (Figure 17).
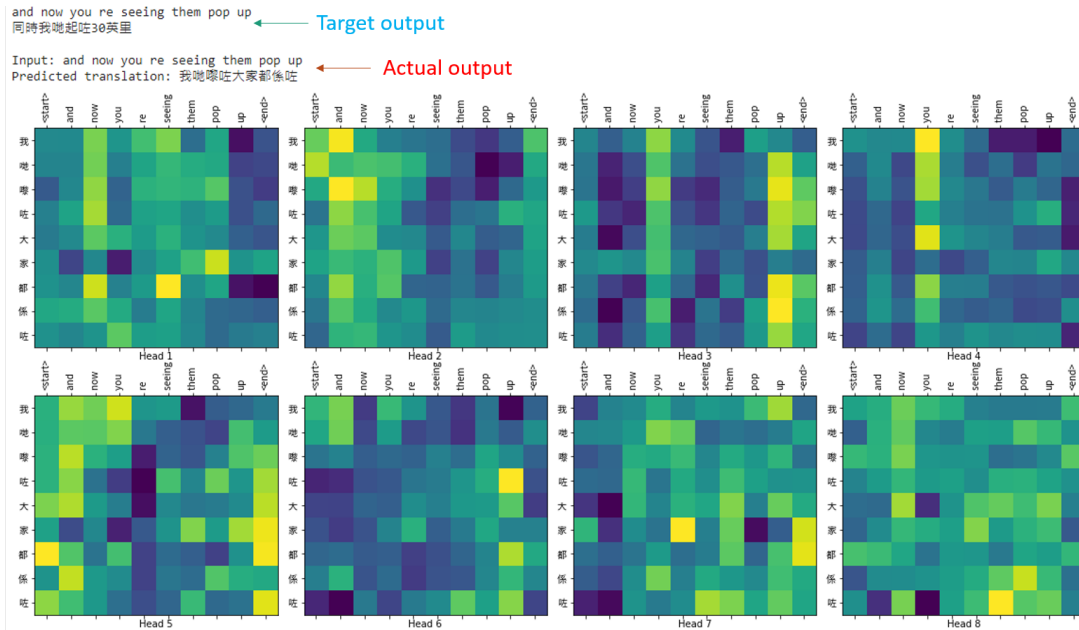


Figure 17: Sample output of multi-head attention in Transformer.

The actual input and targeted output have different meanings. However, the predicted sentence (actual output) have related meaning to the input sentence, the attention

matrix does not show correlations to explain the output. This result is incredible and considered as unexpected, showing that the MT models are learning from the training.

### 4.2.2 Corpus

In this project, the size of corpus data obtained is around 9000 pairs. While MT models generally uses > 100k pairs of data to perform training. Since the dataset is comparatively small, the statistical/probabilistic output is unable to obtain enough training data, which contributes to low accuracy.

As shown in Figure 17, the corpus is not exactly parallel and not completely accurate. Even when the predicted output has similar meaning to the input sentence, the score is still low because of inaccurate targeted output. Noisy data affect the training process in MT models, and thereafter affect the evaluation result.

### 4.2.3 MT

Overfit and underfit is one of the largest challenge in training NMT models [52], where the hyperparameter needs to match with the size of input data. Without the fine-tune of the hyperparameter described in Table 1, the performance of the MT systems would not be the best. Whereas the fine-tuning process is time consuming requires previous experience to perform efficiently, this process is not included at the current stage.

As shown in left side of Figure 14, One of the possibility is that these models show inability in Cantonese. This is the least likely situation, because the models adopted in this project are carefully selected based on their popularity and functionality.

### 4.2.4 Evaluation

BLEU is used as the only metric to evaluate the MT models. Although it had been used as the core evaluation metric in the past 20 years, different critiques had arose issues around it. One of them is [53], which had shown that BLEU cannot correctly evaluate languages without word boundaries, and Cantonese is one of those language. This problem can contribute to the current results.

## 4.3 Expected results

### 4.3.1 RBMT

RBMT is not highly dependent on the amount of training data. The core principle is to use grammatical rules to transform the SL into TL. Training data are used to providing the transferring lexicons and different sentence structure. Grammatical rules and lexicons are considerably reliable because of the efforts in research for Cantonese. Thus, the models is expected perform better than the 2 trained models..

### 4.3.2 PBSMT

PBSMT is a probabilistic model requires clean and large amount of data, and possibly face the same issue mentioned in Section 4.2.2. On the other hand. fine-tune are not required in the model which reduce the difficulty to obtain maximized result. Therefore, the expected outcome would be similar to the results of GRU and Transformer models or can perform slightly better.

## 4.4 Possible obstacle for Cantonese MT

Several possible explanations of the current results are discussed above (Section 4.2). These discussions are likely to correlate with the lack of Cantonese-MT.

### 4.4.1 Lack of data source

Despite both languages are the official languages in Hong Kong, Cantonese-English bilingual corpus and materials barely exist. Most documents are written in MSC instead of Cantonese. In Section 1.2, the study even used MSC as like 'Cantonese'. The lack in written Cantonese text has indeed discouraged researchers to construct an corpus from scratch, and thus discourages the implementation of MT systems.

### 4.4.2 Lack of unified representation

During data collection and preprocessing process (Section 2.1.1, 2.1.3), various assumptions are made to handle the data. The conflict in borrowing phonemic characters, transferring from simplified to traditional Chinese and resolving polysemous, are obstructing scholars to make consensus on on Character choice.

# 5 Project status

## 5.1 Limitations and recommendations

### 5.1.1 Vocabularies and genres of corpus

The vocabulary of movies and talks are heavily based on the genre. For instance, third person narratives are more likely to appear in speeches and talks, and the content are full of technical jargon of specific topic.

The solution is to include as many genre as possible by manual effort. Data retrieving are limited in speech-to-text method, since the quality of transcription are required to be assure by human. So, we hope to discover other data sources to create a corpus larger in scale with ensured quality.

A brand new Cantonese-English bilingual corpus, SpiCE [54], is going to be released in early 2021. The transcription are done by profession linguist and certainly have reliable and assured quality comparing to our current corpus. Thus, we will include it as the data source once it release.

### 5.1.2 Choice of MT

Only 4 MT models are chosen to be implemented in this project. The coverage is smaller than the paradigm of MT, so thorough analysis and comparison cannot be done. Some models are not exist with any FOS platform as discussed in Section 2.2.2. We hope to include more alternatives in this project, however due to the resources and cost, we had covered enough variants and had to give up less popular models.

### 5.1.3 Cost

Each training requires computational power and a lot of time in hours or even days. Restricted by the architecture, some MT models are unable to perform multi-thread training like GRU. With the amount of models to be implemented and the ever growing size of corpus in the future, it is difficult to obtain the best performance of the MT models.

Therefore, the refinement in MT models are expected to be limited and not to be perfect. Exhaustive training are not recommended, and the hyperparameter will be

adjusted based on experience rather than blind test.

## 5.2   Future planning

Remaining works of this project include enhancing corpus data, training the remaining MT models and evaluating their results. Base on the results, we will attempt to derive the reasons for the lack of Cantonese-based MT.

### 5.2.1   Improvement of corpus

The corpus size is not enough for training a functional MT model, and improvement in quality are also required. This process are considered to be continuous and without a fix date, whenever new usable data exist.

### 5.2.2   Training and testing MTs

The training of the remaining MT models are proposed to be done at the next stage. The fine-tune process will be performed during training. We will perform training in January to March.

Then, we will test, evaluate and discuss the results of all MT models in March to April. Different domains will be investigated such as translation direction, vocabulary and character based input, translation based and speech based input, and syntax structure like right dislocation.

### 5.2.3   Completion of project

At the end of this project, we hope to provide a foundation and preliminary study for further development on Cantonese-based MT, including the resource used and difficulties encountered. We hope to provide an functional system with BLEU score of 10. The trained models with corpus data will be available and open-source for future developers, and the code used will also be presented.

# 6    Conclusion

There is a lack in researches of Machine Translation between Cantonese and English. We discussed the linguistic and computational information for MT, and described details about data collection, MT systems and evaluating methods.

In this project, We implemented several Cantonese-based MT systems and investigated their performance. We noticed that Cantonese-English bilingual corpus do not exist, and chose to use TED talks and Netflix to construct our own corpus. Then, the well-known automated metric BLEU is used to evaluate translation outputs.

Currently, we had successfully installed the all chosen MTs and constructed the bilingual corpus. Two of the models had completed pilot run and tested for their functionalities. The input data are evaluated for the performance in vocabulary and character based approach. The MT models use default hyperparameter to perform training.

The current results are unsatisfactory with the higher score of 3.75, and do not meeting our goal of 10 in performance. However, some systems had unexpected output and incredible results. Possible reasons behind malfunctioning includes corpus size, corpus quality, lack of fine-tune, and metric that cannot evaluate correctly.

There are several limitations and difficulties found, including the vocabulary base and the choice of MT and training cost. We hope to improve these areas in manageable time and cost by manual effort and avoid exhaustive training.

In the next stage, we will continue on improving the corpus, train and test MTs before April. At the end of this project, we aim to provide a foundation for further researches in Cantonese-based MT. We will also release the source code and trained models with corpus data for future researchers and developers.

# 7 Bibliography

[1] K. Shah, "Model adaptation techniques in machine translation," Ph.D. dissertation, Université du Maine, 06 2012.

[2] M. R. Costa-Jussà and M. Farrús, "Statistical machine translation enhancements through linguistic levels: A survey," *ACM Computing Surveys (CSUR)*, vol. 46, no. 3, pp. 1–28, 2014.

[3] S. Tripathi and J. K. Sarkhel, "Approaches to machine translation," *Annals of Library and Information Studies*, vol. 57, pp. 388–393, 2010.

[4] Y. Wu, X. Li, and S. C. Lun, "A structural-based approach to cantonese-english machine translation," in *International Journal of Computational Linguistics & Chinese Language Processing, Volume 11, Number 2, June 2006*, 2006, pp. 137–158.

[5] M. Forcada, M. Ginestí-Rosell, J. Nordfalk, J. O'Regan, S. Ortiz-Rojas, J. Pérez-Ortiz, F. Sánchez-Martínez, G. Ramírez-Sánchez, and F. Tyers, "Apertium: A free/open-source platform for rule-based machine translation," *Machine Translation*, vol. 25, pp. 127–144, 06 2011.

[6] R. D. Brown, "The cmu-ebmt machine translation system," *Machine translation*, vol. 25, no. 2, p. 179, 2011.

[7] C. Olah, "Understanding lstm networks," 2015. [Online]. Available: https://colah.github.io/posts/2015-08-Understanding-LSTMs/

[8] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.

[9] W. J. Hutchins and H. L. Somers, *An introduction to machine translation*. Academic Press London, 1992, vol. 362.

[10] D. M. Eberhand, G. F. Simons, and C. D. Fenning, Eds., *Ethnologue: Languages of the World*, twenty-third ed. Dallas, TX: SIL International, 2020, https://www.ethnologue.com/.

[11] J. Ball, *Cantonese Made Easy: A Book of Simple Sentences in the Cantonese Dialect, with Free and Literal Translations, and Directions for the Rendering of*

*English Grammatical Forms in Chinese.* Kelly & Walsh, Limited, 1907. [Online]. Available: https://books.google.com.hk/books?id=3msuAAAAYAAJ

[12] S. Matthews and V. Yip, *Cantonese: A comprehensive grammar.* Routledge, 2013.

[13] T.-s. Wong, K. Gerdes, H. Leung, and J. Lee, "Quantitative comparative syntax on the cantonese-mandarin parallel dependency treebank," in *Proceedings of the Fourth International Conference on Dependency Linguistics (Depling 2017), September 18-20, 2017, Università di Pisa, Italy*, no. 139. Linköping University Electronic Press, 2017, pp. 266–275.

[14] D. Snow, *Cantonese as written language: The growth of a written Chinese vernacular.* Hong Kong University Press, 2004, vol. 1.

[15] J. Oladosu, A. Esan, I. Adeyanju, B. Adegoke, O. Olaniyan, and B. Omodunbi, "Approaches to machine translation: a review," *FUOYE Journal of Engineering and Technology*, vol. 1, no. 1, 2016.

[16] T. Poibeau, *Machine Translation*, ser. The MIT Press Essential Knowledge series. MIT Press, 2017. [Online]. Available: https://books.google.com.hk/books?id=HYc3DwAAQBAJ

[17] M. Nagao, "A framework of a mechanical translation between japanese and english by analogy principle," *Artificial and human intelligence*, pp. 351–354, 1984.

[18] J. Liu and Y. Yu, "A cantonese-english machine translation system polyu-mt-99," in *MT SUMMIT VII: MT in the great translation era : proceedings of Machine Translation Summit VII.* Kent Ridge Digital Labs, Singapore: Asia-Pacific Association for Machine Translation, September 2019, pp. 481–486.

[19] L. Barrault, O. Bojar, M. R. Costa-Jussà, C. Federmann, M. Fishel, Y. Graham, B. Haddow, M. Huck, P. Koehn, S. Malmasi *et al.*, "Findings of the 2019 conference on machine translation (wmt19)," in *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, 2019, pp. 1–61.

[20] *Hong Kong Chinese Lexical Lists for Primary Learning.* Hong Kong: Chinese Language Education Section, Curriculum Development Institute, Education Bureau, The Government of the Hong Kong Special Administrative Region, 2007.

[21] W. M. Ng, A. C. Kwan, T. Lee, and T. Hain, "Shefce: A cantonese-english bilingual speech corpus – speech recognition model sets and recording transcripts," Mar 2017. [Online]. Available: https://figshare.shef.ac.uk/articles/dataset/ShefCE_A_Cantonese-English_bilingual_speech_corpus_--_speech_recognition_model_sets_and_recording_transcripts/4522925/1

[22] R. W. Ng, A. C. Kwan, T. Lee, and T. Hain, "Shefce: A cantonese-english bilingual speech corpus for pronunciation assessment," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2017, pp. 5825–5829.

[23] "Translate." [Online]. Available: https://www.ted.com/participate/translate

[24] M. Cettolo, C. Girardi, and M. Federico, "Wit3: Web inventory of transcribed and translated talks," in *Conference of european association for machine translation*, 2012, pp. 261–268.

[25] L. Richardson, "Beautiful soup documentation," *April*, 2007.

[26] R. C. Souza, "raryelcostasouza/pytranscriber v1.4," Jan. 2020. [Online]. Available: https://github.com/raryelcostasouza/pyTranscriber

[27] "Speech-to-text documentation nbsp;|nbsp; cloud speech-to-text documentation." [Online]. Available: https://cloud.google.com/speech-to-text/docs/

[28] J. L. Lee, "Pycantonese: Cantonese linguistic research in the age of big data," *Talk at the Childhood Bilingualism Research Centre, the Chinese University of Hong Kong*, 2015.

[29] S. Sato and M. Nagao, "Toward memory-based translation," in *COLNG 1990 Volume 3: Papers presented to the 13th International Conference on Computational Linguistics*, 1990.

[30] H. Somers, "Example-based machine translation," *Machine translation*, vol. 14, no. 2, pp. 113–157, 1999.

[31] A. B. Phillips, "Cunei: open-source machine translation with relevance-based models of each translation instance," *Machine Translation*, vol. 25, no. 2, p. 161, 2011.

[32] S. Dandapat, M. L. Forcada, D. Groves, S. Penkale, J. Tinsley, and A. Way, "Open-matrex: a free/open-source marker-driven example-based machine translation system," in *International Conference on Natural Language Processing*. Springer, 2010, pp. 121–126.

[33] R. Zens, F. J. Och, and H. Ney, "Phrase-based statistical machine translation," in *Annual Conference on Artificial Intelligence*. Springer, 2002, pp. 18–32.

[34] P. Koehn, H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens *et al.*, "Moses: Open source toolkit for statistical machine translation," in *Proceedings of the 45th annual meeting of the ACL on interactive poster and demonstration sessions*. Association for Computational Linguistics, 2007, pp. 177–180.

[35] K. Cho, B. Van Merriënboer, D. Bahdanau, and Y. Bengio, "On the properties of neural machine translation: Encoder-decoder approaches," *arXiv preprint arXiv:1409.1259*, 2014.

[36] O. Bojar, R. Chatterjee, C. Federmann, Y. Graham, B. Haddow, M. Huck, A. Jimeno Yepes, P. Koehn, V. Logacheva, C. Monz, M. Negri, A. Névéol, M. Neves, M. Popel, M. Post, R. Rubino, C. Scarton, L. Specia, M. Turchi, K. Verspoor, and M. Zampieri, "Findings of the 2016 conference on machine translation," in *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*. Berlin, Germany: Association for Computational Linguistics, Aug. 2016, pp. 131–198. [Online]. Available: https://www.aclweb.org/anthology/W16-2301

[37] P. Koehn and R. Knowles, "Six challenges for neural machine translation," in *Proceedings of the First Workshop on Neural Machine Translation*, 2017, pp. 28–39.

[38] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[39] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.

[40] K. Cho, B. van Merrienboer, D. Bahdanau, and Y. Bengio, "On the properties of neural machine translation: Encoder-decoder approaches," *CoRR*, vol. abs/1409.1259, 2014. [Online]. Available: http://arxiv.org/abs/1409.1259

[41] A. Zhang, Z. C. Lipton, M. Li, and A. J. Smola, *Dive into Deep Learning*, 2020, https://d2l.ai.

[42] O. Firat, K. Cho, and Y. Bengio, "Multi-way, multilingual neural machine translation with a shared attention mechanism," *arXiv preprint arXiv:1601.01073*, 2016.

[43] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," 2014.

[44] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, M. Kudlur, J. Levenberg, R. Monga, S. Moore, D. G. Murray, B. Steiner, P. Tucker, V. Vasudevan, P. Warden, M. Wicke, Y. Yu, and X. Zheng, "Tensorflow: A system for large-scale machine learning," in *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*. Savannah, GA: USENIX Association, Nov. 2016, pp. 265–283. [Online]. Available: https://www.usenix.org/conference/osdi16/technical-sessions/presentation/abadi

[45] A. Vaswani, S. Bengio, E. Brevdo, F. Chollet, A. N. Gomez, S. Gouws, L. Jones, Łukasz Kaiser, N. Kalchbrenner, N. Parmar, R. Sepassi, N. Shazeer, and J. Uszkoreit, "Tensor2tensor for neural machine translation," 2018.

[46] P. Koehn and C. Monz, "Manual and automatic evaluation of machine translation between european languages," in *Proceedings on the Workshop on Statistical Machine Translation*, 2006, pp. 102–121.

[47] C. Callison-Burch, "Fast, cheap, and creative: Evaluating translation quality using amazon' s mechanical turk," in *Proceedings of the 2009 conference on empirical methods in natural language processing*, 2009, pp. 286–295.

[48] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 2002, pp. 311–318.

[49] J. Zhou, Y. Cao, X. Wang, P. Li, and W. Xu, "Deep recurrent models with fast-forward connections for neural machine translation," *Transactions of the Association for Computational Linguistics*, vol. 4, pp. 371–383, 2016.

[50] S. Varsamopoulos, K. Bertels, and C. G. Almudever, "Designing neural network based decoders for surface codes," *arXiv preprint arXiv:1811.12456*, 2018.

[51] A. Lavie, "Evaluating the output of machine translation systems," *AMTA Tutorial*, vol. 86, 2010.

[52] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An introduction to statistical learning*. Springer, 2013, vol. 112.

[53] C. Callison-Burch, M. Osborne, and P. Koehn, "Re-evaluation the role of bleu in machine translation research," in *11th Conference of the European Chapter of the Association for Computational Linguistics*, 2006.

[54] K. A. Johnson, M. Babel, I. Fong, and N. Yiu, "SpiCE: A new open-access corpus of conversational bilingual speech in Cantonese and English," in *Proceedings of the 12th Language Resources and Evaluation Conference*. Marseille, France: European Language Resources Association, May 2020, pp. 4089–4095. [Online]. Available: https://www.aclweb.org/anthology/2020.lrec-1.503