

与えた外部情報の再予測モデルを組み込んだ ニューラル文生成モデルの検討

隆辻 秀和^{1,a)} 吉野 幸一郎^{1,b)} 須藤 克仁^{1,c)} 中村 哲^{1,d)}

概要：言語生成は、与えられた外部情報のセットに対して、自然言語文をドメインに適当な形で生成するタスクである。近年、言語生成に用いられるニューラルネットワークを用いた手法は、より自然で柔軟な応答生成が実現できることが知られている。一方で、入力となる外部情報に対応する文生成を単語予測のモデルで行うため、モデルがどの情報を利用し文を生成したかを説明することが難しい。そこで本研究では、与えた外部情報を生成文に反映することを保証するため、与えた外部情報を再予測するモデルと再予測の結果に対する損失を利用した。アノテーション済みのコーパスを用いた実験を行い、生成された文の評価と、生成文に含まれる情報の精度評価を行った。

1. はじめに

言語生成は自然言語処理における重要なタスクの一つである。言語生成器は生成で考慮すべき外部情報を入力として受け取り、ドメインや文脈に合わせて適切な文生成を行う。機械翻訳であれば、翻訳元の言語で記述された文が与えられ、意味を損なわないように翻訳先の言語で記述した文を生成する。また、質問応答システムではユーザーから投げかけられた質問文の他に、システムが保持している知識などの追加情報も用いて回答を生成する。

言語生成の手法として、テンプレートを用いるものやルールに基づく生成など複数の手法が考えられてきた [1], [2], [3] が、近年ではニューラルネットワークを用いた手法が一般的になっている。ニューラルネットワークを用いた手法は、従来の手法に比べて柔軟かつ自然な文生成が行えることが、さまざまな研究結果より知られている。最もよく知られているモデルとして Sequence-to-Sequence (Seq2Seq) [4] がある。入力の単語系列と出力の単語系列が与えられた時にモデルが入出力間の対応関係を学習することで、学習時に存在しない未知の系列に対しても従来手法に比べて自然な生成を可能にしている。

入力として文以外の外部情報を考慮するような言語生成モデルはさまざまに提案されている。Li ら [5] は出力文に

反映すべき言語的特徴として個人性を取り扱う手法を提案している。この手法では、デコーダの各ステップで反映すべき特徴を表現したベクトルも入力として与えることで条件付き生成を行う。また、Eric ら [6] は Seq2Seq モデルのエンコーダに、質問文と与えられている外部情報のスロットの値を入力し、注意機構やコピーメカニズムを利用して外部情報を出力系列に反映する手法を提案している。

これらのモデルは、生成する単語系列に対する予測誤りを元に学習を行うため、追加で与えた情報が適切に考慮されていることが保証されないという問題がある。この問題は、特に、先述したような質問応答のシステムにおける応答生成のように、クエリの他に与えられる外部情報が出力に含まれることが期待される状況では大きな課題となる。

本研究では、モデルに入力文の他に与えた外部情報が生成文に含まれることを保証するような言語生成のモデルを提案する。具体的には出力となる単語系列の予測に加えて、与えた外部情報の再予測を行うモデルを用いることで、与えた外部情報についても損失を計算する。これによってモデルが生成する文が与えた外部情報を含むことを期待した。提案するモデルの有効性を確認するために、先行研究 [6] と同じコーパスを用いて実験を行い、結果を確認した。

2. 関連研究

ニューラルネットワークを用いた言語生成において、入力となる文以外の情報を考慮して生成を行うという研究は広く行われている。Seq2Seq によるモデルの他に、メモリネットワークと呼ばれる Key-Value 構造を行列で表現したネットワークによって外部情報を考慮するモデルが広く

¹ 奈良先端科学技術大学院大学
NARA Institute of Science and Technology
a) takatsuji.hidekazu.sx1@is.naist.jp
b) koichiro@is.naist.jp
c) sudoh@is.naist.jp
d) s-nakamura@is.naist.jp

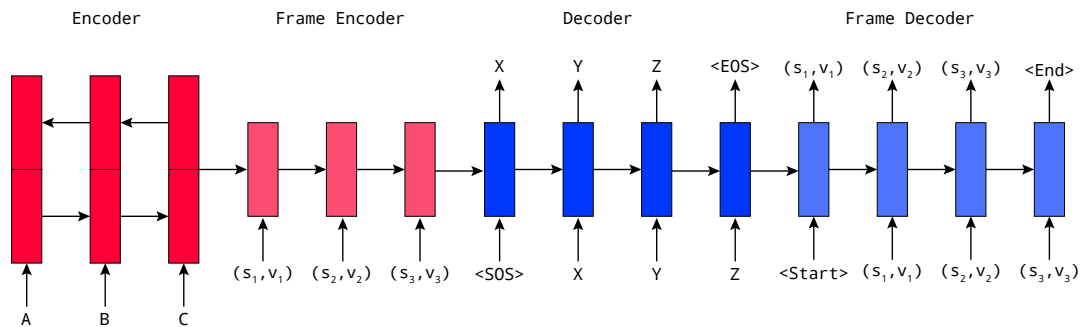


図 1 提案モデルの概要

Fig. 1 A blueprint of proposed model

知られている。Sukhbaatar ら [7] は外部情報を Key-Value からなる辞書オブジェクトとみなし、入力として与えられるクエリとの間の内積注意を計算することで、外部情報を参照した応答生成を行うモデルを提案している。このモデルでは回答が一つの単語に定まるような質問応答タスクにおいて有効性が確認されている。また、Madotto ら [8] は Sukhbaatar らのモデルの拡張を行い、入力として与えられるクエリもまたメモリに入力することで、出力として単語の系列を予測するモデルを提案している。

Seq2Seq の枠組みの上で、生成文の中にあらかじめ定められた辞書オブジェクトの値を埋め込むという研究として、Qian ら [9] によるものが知られている。この研究では、システムに与えられた辞書オブジェクトの値を利用するかどうか決定する予測器を用意し、予測期の出力に応じてデコード方法を変化させるという手法を取っている。本研究における手法と異なり、この手法では与えられた辞書オブジェクトのある Key-Value の組しか考慮できないという問題があるが、一方で辞書オブジェクト内の値が必ず出力されるようにデコードを実行しているため、出力に Key-Value の値を反映することを従来モデルに比べてより良い形で保証することができる。

これらの先行研究に対して提案モデルでは、目的の異なる複数のモデルを同時に学習させることで、生成タスクの精度を向上させることを意図している。このような学習はマルチタスク学習と呼ばれ、系列モデルに対してマルチタスク学習を適用する手法は既に Luong ら [10] によって研究されている。Luong らの研究ではエンコーダとデコーダに対して 1 対多や多対多のモデルについて取り上げており、これらのケースでは個々のエンコーダあるいはデコーダについては並列の関係となっている。一方で、本研究では複数のエンコーダ及びデコーダを直列に取り扱っており、この点において先行研究と異なる。

3. 条件付き応答生成モデル

本研究で提案するニューラル文生成モデルについて図 1 に提案モデルの概要を示す。

Seq2Seq における枠組みと同様にクエリ \mathbf{q} を受け取り、応答文 \mathbf{r} を生成する。提案モデルでは、外部情報 \mathcal{F} を与えるためにフレームエンコーダを導入し、これをクエリに対するエンコーダの後段に接続する。また、与えた外部情報が文 \mathbf{r} に含まれることを保証するために、学習時のみ与えた外部情報 \mathcal{F} を再予測するフレームデコーダを通常のデコーダの後段に接続する。これらの改良によって、与えた外部情報が出力に含まれることを保証することを目的としている。

3.1 エンコーダ

応答生成におけるコンテキストの一つであるクエリ $\mathbf{q} = \{q_1, q_2, \dots, q_N\}$ を一般的な Seq2Seq モデルと同様に Recurrent Neural Network (RNN) を用いて単一のベクトル表現に変換する。RNN のユニットには LSTM[11] を利用し、双方向 RNN を利用した。

$$\vec{h}_t^{enc} = \text{RNN}(q_t, \vec{h}_{t-1}^{enc}) \quad (1)$$

$$\overleftarrow{h}_t^{enc} = \text{RNN}(q_t, \overleftarrow{h}_{t+1}^{enc}) \quad (2)$$

3.2 フレームエンコーダ

生成時に外部情報 $\mathcal{F} = \{(s_1, v_1), (s_2, v_2), \dots, (s_K, v_K)\}$ によって生成文を制約するために、外部情報 \mathcal{F} をエンコーダによって得られるベクトル表現と同様のベクトル表現に変換する。外部情報はスロット s とスロットに紐づく値 v からなる組のセットとして表現されるが、これを任意の順序を持つ (s, v) の組の系列とみなして単方向 RNN を用いてエンコードする。フレームエンコーダは通常のエンコーダの後段に接続する形で利用するため、隠れ状態の初期状態はエンコーダの最終ステップにおける隠れ状態となる。

スロットと値の組み合わせ (s, v) をエンコードする手法はいくつか考えられるが、提案モデルでは (s, v) の組をひとつのトークンとしてエンコードする。

$$h_0^{fenc} = \text{concat}(\vec{h}_N^{enc}, \overleftarrow{h}_N^{enc}) \quad (3)$$

$$h_k^{fenc} = \text{RNN}((s_k, v_k), h_{k-1}^{fenc}) \quad (4)$$

3.3 デコーダ

エンコードされたクエリおよび外部情報を元にして応答文の単語列 $\mathbf{r} = \{r_1, r_2, \dots, r_M\}$ を予測する。生成には単方向のRNNを用いて系列の予測を行う。

各ステップで直前の予測単語 r'_{t-1} と隠れ状態 h_{t-1}^{dec} から現在の隠れ状態を計算し、単語予測を行う。学習時にはTeacher Forcingを行うため入力単語は直前の予測単語 r'_{t-1} ではなく正解系列における単語 r_{t-1} を用いる。

$$h_t^{dec} = \text{RNN}(r'_{t-1}, h_{t-1}^{dec}) \quad (5)$$

$$r'_t = \text{softmax}(\mathbf{W}_o^{dec} h_t^{dec}) \quad (6)$$

3.4 フレームデコーダ

入力として与えた外部情報 \mathcal{F} をデコーダが生成する単語列 \mathbf{r}' に対する制約とすることが、本研究における目的である。しかし、既存のデコーダでは単語予測誤りのみに基づいて損失を計算するため、与えた外部情報について陽に保証する仕組みが存在していない。提案モデルでは、与えた外部情報がデコーダが生成する単語列に考慮されていることを保証するために、与えた外部情報そのものを再予測するフレームデコーダを導入する。フレームデコーダはデコーダと同様に単方向のRNNを用いて与えた外部情報の再予測を行う。この時、フレームデコーダが再予測する外部情報の順序はフレームエンコーダに外部情報が与えられた際の順序と同一となるように学習を行う。

$$h_t^{fdec} = \text{RNN}((s'_{t-1}, v'_{t-1}), h_{t-1}^{fdec}) \quad (7)$$

$$(s'_t, v'_t) = \text{softmax}(\mathbf{W}_o^{fdec} h_t^{fdec}) \quad (8)$$

3.5 損失関数

提案モデルは、単語の予測誤りに基づく損失と外部情報の予測誤りに基づく損失の二つを学習に用いる。単語予測についての交差エントロピー損失 \mathcal{L}_w および外部情報予測についての交差エントロピー損失 \mathcal{L}_{fr} の重みつき線形和 $\mathcal{L} = \mathcal{L}_w + \alpha \mathcal{L}_{fr}$ を損失関数として利用する。

4. 実験設定

提案モデルについて二種類の実験を行った。一つは、提案モデルの有効性に関する実験。もう一つは損失関数もたらす影響についての実験である。はじめに、二つの実験に共通する設定について述べ、その後個々の実験で用いた設定について述べる。

4.1 共通設定

実験では生成文が含むべき外部情報についてアノテーションが行われたコーパスとして、DSTC2[12]において配布されたデータセットを利用した。コーパスはレストラン情報案内についての対話を収録しており、情報を案内するシステムとシステム利用者の発話が収録されている。各発話には対話状態として対話行為及び発話文が参照するフレームがアノテートされている。コーパス全体で16,551件の発話ペアが含まれており、そのうち15,723件の発話ペアを学習用に、残りの発話ペアを評価のために利用した。DSTC2コーパスには255種類の異なるフレームが含まれており、これらを外部情報として利用した。実験では、システムの発話をクエリとして、システム利用者の各発話にアノテートされた対話状態のフレームを外部情報として取り扱い、システム利用者の発話を予測した。

DSTC2におけるコーパスは含まれる発話ペアの数が少なく、自然な応答を学習するためには不足していると考えられる。そのため、モデルの事前学習を行うためにReddit*¹より収集した発話ペア、約50万件を利用した。

それぞれのコーパスに対する前処理としてSentencePiece[13]を用いたトークナイズを行った。トークナイザの学習にはRedditより収集したコーパスを利用し、単語サイズは16,000とした。また、学習時には最適化手法はAdam[14]を利用した。

4.2 提案モデルの有効性検証

提案モデルの有効性を確認するために、ベースラインモデルとしてフレームデコーダを含まないモデルを用いて、結果の比較を行った。この実験では、損失関数の重み $\alpha = 1.0$ に固定した。

4.3 損失関数の重みによる差異

提案モデルでは単語の予測誤りに基づく損失と外部情報の再予測誤りに基づく損失の二つの損失を考慮している。学習において、外部情報の再予測誤りに基づく損失をより重視することが、生成文に対して与えた外部情報が反映されるという保証をより強く与えるものとなるかについて比較を行った。具体的には損失関数の重み α を変更とした時に、生成文に含まれる外部情報に変化が発生するのかについて実験を行う。実験では $\alpha = \{2.0, 5.0, 10.0\}$ の場合について生成結果の比較を行った。

5. 実験結果

それぞれの実験における結果を判断するために複数の評価尺度による結果と生成例を示す。評価尺度として、BLEU[15]、パープレキシティ、Entity.F1、および提案モ

*¹ <https://reddit.com>

表 1 生成例の比較

Table 1 A comparison of generation example

query	Hello, welcome to the Cambridge restaurant system? You can ask for restaurants by area, price range or food type. How may I help you?
response	restaurant in the south part of town that serves indian food
frame	{{(food, indian) (area, south)}}
baseline	im looking for a restaurant in the south part of town that serves indian food
proposed	im looking for a restaurant in the south part of town that serves indian food
query	I'm sorry but there is no restaurant serving kosher food
response	im looking for a restaurant in the north part of town serving french type of food
frame	{{(food, french),(area, north)}}
baseline	i want a restaurant in the north part of town that serves french food
proposed	im looking for a restaurant in the north part of town that serves korean food

表 2 ベースラインと提案モデルの評価尺度による比較

Table 2 A comparison of metric evaluation result between baseline and proposed model

Model	BLEU	Ppl.	Ent.F1	Accuracy
baseline	55.42	2.1890	69.15	-
proposed	56.06	2.1922	67.67	0.1181

デルのみを対象として外部情報の再予測精度を用いた。

BLEU は主に機械翻訳における翻訳の評価を行うために用いられるが、複数の先行研究 [5], [6] においてモデルが生成した応答を評価する手法として用いられている。モデルが生成した文ごとに BLEU スコアの計算を行い、その平均値をモデルの BLEU として利用した。

Entity.F1 は先行研究 [6] においてモデルが与えた外部情報を考慮しているか確認するために用いられた指標である。テストデータには正解となる外部情報が与えられているため、これを利用してモデルが生成した応答内に正解となる外部情報が含まれているかマイクロ平均した F1 スコアを使って評価する。

5.1 提案モデルの有効性検証

表 1 にベースラインと提案モデルにおける生成結果の比較を、表 2 に各評価尺度における評価結果を示す。まず、定量的な評価に着目すると、表 2 より BLEU およびパープレキシティを用いた評価では、ベースラインと提案モデルの間に大きな違いは見られなかった。Entity.F1 に着目すると、提案モデルはベースラインに比べて 1.5 ポイント程度低いスコアとなった。また、提案モデルにおける外部情報再予測の精度は 11.8% にとどまった。テストコーパスにおける外部情報再予測のチャンスレートは 12.7% であるため、提案モデルは外部情報を適切に予測できていないことが示されている。

評価尺度を用いた比較ではベースラインと提案モデルにおいて大きな違いが見られなかったが、生成結果による定性的な比較ではどのような差異が見られるのかを確認する。

表 4 再予測に関する損失の重みを変更した際の評価結果

Table 4 A metric evaluation result changing coefficient that multiplies frame prediction loss

weight	BLEU	ppl.	Ent.F1	Accuracy
$\alpha = 2.0$	53.71	2.1702	68.39	0.1162
$\alpha = 5.0$	55.76	2.1965	68.26	0.1124
$\alpha = 10.0$	54.79	2.3078	66.50	0.1272

表 1 を見れば、一つめの生成例ではベースラインと提案モデルでまったく同じ文が生成されており、どちらの生成結果も外部情報として与えられる情報を過不足なく含んでいる。一方、二つめの生成例ではベースラインでは正しく与えられた外部情報を生成文に含んでいるが、提案モデルでは与えられた情報とは異なる内容についての文を生成している。文中斜体とした単語では、ベースラインは正しく与えられた外部情報を反映して french と生成できているが、提案モデルは与えられた情報ではなく korean を生成している。他の生成例においても、ベースラインの生成結果に比べて提案モデルが与えられた外部情報をよりよく含んでいるような生成例は発見できなかった。

5.2 損失関数の重みによる差異

損失関数の重みを変更した時に、生成される文に対してどのように影響が発生するかについて、評価尺度および生成例を用いて確認した。表 3 は表 1 と同様の入力を与えたときの $\alpha = 2.0, 5.0, 10.0$ における生成結果を示す表であり、表 4 は α を変更した際のそれぞれの評価尺度を用いた評価結果を示す表である。まず、定量的評価において比較すると、表 2 から $\alpha = 5.0$ の条件において、ベースラインから BLEU がわずかに向上しているが、差としては非常に小さいものであり、ほぼ誤差と言える。Entity.F1 について比較すると、 $\alpha = 10.0$ の場合が最も低いスコアとなっており、 $\alpha = 2.0$ および 5.0 ではベースラインより低く、しかし提案モデルよりは高いスコアを示している。この結果より、 α を変更したことによって生成結果に外部情

表 3 α を変化したときの生成文の変化
 Table 3 Difference of generated sentence when alpha increases

weight	sentence
frame	{(food, indian) (area, south)}
$\alpha = 2.0$	im looking for a restaurant in the south part of town that serves indian food
$\alpha = 5.0$	im looking for a restaurant in the south part of town that serves indian food
$\alpha = 10.0$	im looking for a restaurant in the south part of town that serves
frame	{(food, french),(area, north)}
$\alpha = 2.0$	im looking for a restaurant in the north part of town serving <i>spanish</i> food
$\alpha = 5.0$	i want a restaurant in the north part of town that serves <i>french</i> food
$\alpha = 10.0$	i want a restaurant in the north part of town that serves <i>korean</i> food

報が反映されやすくなったわけではないということがわかる。また、再予測精度は $\alpha = 5.0$ の時が最も小さく、 α を変更した際の再予測精度の上昇幅に関しては微々たるもので、 $\alpha = 1.0$ と $\alpha = 10.0$ を比較しても精度は 1% 向上したのみにとどまっている。このことから、 α の大きさと再予測精度の関係はおおむね誤差の範囲に収まるともと考えられる。また、重みを変更してもなお、再予測精度がチャンスレートを上回らないことから、提案モデルが再予測を行えるものでないことが裏付けられている。

次に α を変更した際の生成文の変化について比較する。表 3 より、一つめの生成例では $\alpha = 10.0$ を除く全ての生成結果において同じ文章が生成されている。したがって、どの α においても与えた外部情報を過不足なく含んだ生成結果となっていることがわかる。一方、二つめの生成例では、 α を変化した時に生成される文章はどれも細部が異なっており、具体的には表内斜体で示した単語がどの生成結果でも変化している。 $\alpha = 5.0$ における生成例では、与えられた外部情報 (food, french) を正しく反映した生成結果となっているが、それ以外の事例ではそれぞれ *spanish*, *korea* とコーパス内に登場する同一スロットを持つ与えられていない値を生成してしまっている。この結果より、今回提案したアーキテクチャが外部情報を保証するには不十分であるということがわかった。

6. おわりに

本研究では、ニューラルネットワークを用いた応答生成において、生成された応答に外部情報が含まれることを保証するために新しいモデルを提案した。提案モデルでは、与えた外部情報を再予測するモデルを応答生成後に行うことで、与えた外部情報が正しく再予測されることを応答に対する制約として機能することを期待した。しかし、実験の結果より、提案モデルは期待に反し、与えた外部情報が生成結果に含まれることを保証するには不十分であるということがわかった。Entity F1 を用いた評価結果においてスコアが下がっていることから、生成結果に対して与えた外部情報を含むことを保証する形になっていないことが

示されている。

今後は再予測によらない方法で外部情報を反映することを保証する方法を検討する必要がある。

参考文献

- [1] Dale, R., Geldof, S. and Prost, J.-P.: CORAL: Using Natural Language Generation for Navigational Assistance, *Proceedings of the 26th Australasian Computer Science Conference - Volume 16*, ACSC '03, Darlinghurst, Australia, Australia, Australian Computer Society, Inc., pp. 35–44 (2003).
- [2] 山崎 健史, 吉野 幸一郎, 前田 浩邦, 笹田 鉄郎, 橋本 敦史, 船富卓哉, 山肩 洋子, 森信介: フローグラフからの手順書の生成, 情報処理学会論文誌, Vol. 57, No. 3, pp. 849–862 (2016).
- [3] Kondadadi, R., Howald, B. and Schilder, F.: A Statistical NLG Framework for Aggregated Planning and Realization, *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Sofia, Bulgaria, Association for Computational Linguistics, pp. 1406–1415 (2013).
- [4] Sutskever, I., Vinyals, O. and Le, Q. V.: Sequence to Sequence Learning with Neural Networks, *Advances in Neural Information Processing Systems 27* (Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N. D. and Weinberger, K. Q., eds.), Curran Associates, Inc., pp. 3104–3112 (2014).
- [5] Li, J., Galley, M., Brockett, C., Spithourakis, G., Gao, J. and Dolan, B.: A Persona-Based Neural Conversation Model, *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Berlin, Germany, Association for Computational Linguistics, pp. 994–1003 (2016).
- [6] Eric, M. and Manning, C.: A Copy-Augmented Sequence-to-Sequence Architecture Gives Good Performance on Task-Oriented Dialogue, *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, Valencia, Spain, Association for Computational Linguistics, pp. 468–473 (2017).
- [7] Sukhbaatar, S., szlam, a., Weston, J. and Fergus, R.: End-To-End Memory Networks, *Advances in Neural Information Processing Systems 28* (Cortes, C., Lawrence, N. D., Lee, D. D., Sugiyama, M. and Garnett, R., eds.), Curran Associates, Inc., pp. 2440–2448 (2015).
- [8] Madotto, A., Wu, C.-S. and Fung, P.: Mem2Seq: Effectively Incorporating Knowledge Bases into End-to-End Task-Oriented Dialog Systems, *Proceedings of the 56th*

Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Melbourne, Australia, Association for Computational Linguistics, pp. 1468–1478 (2018).

- [9] Qian, Q., Huang, M., Zhao, H., Xu, J. and Zhu, X.: Assigning Personality/Profile to a Chatting Machine for Coherent Conversation Generation, *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence*, Stockholm, Sweden, International Joint Conferences on Artificial Intelligence Organization, pp. 4279–4285 (2018).
- [10] Luong, M.-T., Le, Q. V., Sutskever, I., Vinyals, O. and Kaiser, L.: Multi-task Sequence to Sequence Learning, *the 4th International Conference on Learning Representations (ICLR)* (2016).
- [11] Gers, F.: Long short-term memory in recurrent neural networks, PhD Thesis (2001).
- [12] Henderson, M., Thomson, B. and Williams, J. D.: The Second Dialog State Tracking Challenge, *Proceedings of the 15th Annual Meeting of the Special Interest Group on Discourse and Dialogue (SIGDIAL)*, Philadelphia, PA, U.S.A., Association for Computational Linguistics, pp. 263–272 (2014).
- [13] Kudo, T. and Richardson, J.: SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing, *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, Brussels, Belgium, Association for Computational Linguistics, pp. 66–71 (2018).
- [14] Kingma, D. P. and Ba, J.: Adam: A Method for Stochastic Optimization, *the 3rd International Conference on Learning Representations (ICLR)* (2015).
- [15] Papineni, K., Roukos, S., Ward, T. and Zhu, W.-J.: BLEU: A Method for Automatic Evaluation of Machine Translation, *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, ACL '02*, Stroudsburg, PA, USA, Association for Computational Linguistics, pp. 311–318 (2002).