

# Context Consistency between Training and Inference in Simultaneous Machine Translation

Anonymous ACL submission

## Abstract

Simultaneous Machine Translation (SiMT) aims to yield a real-time partial translation with a monotonically growing source-side context. However, there is a counterintuitive phenomenon about the context usage between training and inference: *e.g.*, in wait- $k$  inference, model consistently trained with wait- $k$  is much worse than that model inconsistently trained with wait- $k'$  ( $k' \neq k$ ) in terms of translation quality. To this end, we first investigate the underlying reasons behind this phenomenon and uncover the following two factors: 1) the limited correlation between translation quality and training (cross-entropy) loss; 2) exposure bias between training and inference. Based on both reasons, we then propose an effective training approach called context consistency training accordingly, which encourages consistent context usage between training and inference by optimizing translation quality and latency as bi-objectives and exposing the predictions to the model during the training. The experiments on three language pairs demonstrate our intuition: our system encouraging context consistency outperforms that existing systems with context inconsistency for the first time, with the help of our context consistency training approach.

## 1 Introduction

Simultaneous machine translation (SiMT) (Cho and Esipova, 2016; Gu et al., 2017; Ma et al., 2019) aims to generate a partial translation while incrementally receiving a prefix of a source sentence. A good SiMT system should not only have *low latency* in the generation process but also yield a complete translation with *high quality*. SiMT has been widely used in many real-world scenarios such as multilateral organizations and international summits (Ma et al., 2019). Hence, there has recently been witnessed a surge of interest in the research about SiMT (Elbayad et al., 2020; Ma et al., 2020; Zhang and Feng, 2021, 2022).

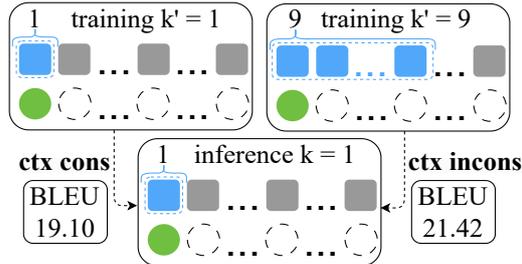


Figure 1: Counterintuitive phenomenon on the context usage between training and inference: in wait-1 inference ( $k = 1$ ), model trained with  $k'=9$  (denoted by “ctx incons”) outperforms the model trained with  $k'=1$  (denoted by “ctx cons”) in terms of BLEU, even though the former model (trained by  $k'=9$ ) induces a mismatch on context usage between training and inference.

In this paper, we shed light on a *counterintuitive phenomenon* on the context usage between training and inference in SiMT: in wait- $k$  inference, model consistently trained with wait- $k$  is worse than that model inconsistently trained with wait- $k'$  ( $k' \neq k$ ) in terms of the evaluation metrics of SiMT, as shown in Figure 1. This phenomenon was first observed by Ma et al. (2019) yet without explanations. Subsequently, such context inconsistency training becomes a standard practice (Elbayad et al., 2020; Zhang and Feng, 2021, 2022), even if this phenomenon is counterintuitive due to the mismatch between training and inference on the usage of partial source-side context.

To investigate the reasons behind the above counterintuitive phenomenon, we conduct experiments from two perspectives: calculating the correlation between translation quality and training (cross-entropy) loss, as well as evaluating the translation quality under the prefix-constrained decoding setting. Our empirical experiments demonstrate two reasons that are responsible for the phenomenon: 1) the limited correlation between translation quality and training loss; 2)

exposure bias between training and inference (§2). Moreover, based on our findings, this paper proposes an effective training approach called context consistency training accordingly, and breaks through the standard practice of inconsistent training. Its key idea is to make the context usage consistent between training and inference by optimizing translation quality and latency as bi-objectives and exposing the model to its own predictions to during the training stage. Particularly, this approach is general to be applied to most SiMT systems (§3).

Experiments conducted across various benchmarks demonstrate that the proposed context consistency training towards bi-objectives achieves substantial gains over the original consistency training based on cross-entropy. In particular, with the help of our training approach, our SiMT systems encouraging context consistency outperform the existing systems with context inconsistency in terms of translation quality and latency (§4).

Our main contributions are:

- This paper sheds light on a counterintuitive phenomenon about context usage between training and inference in SiMT and provides comprehensive explanations for this phenomenon.
- Based on our findings, this paper proposes a simple yet effective approach, known as context consistency training, which encourages consistent context usage between training and inference in SiMT.
- Experiments conducted on three benchmarks and several SiMT systems demonstrate that our system encouraging context consistency outperforms the existing systems with context inconsistency for the first time.

## 2 Rethinking Counterintuitive Phenomenon on Context Usage

### 2.1 Counterintuitive Phenomenon

**Counterintuitive Phenomenon on Valid Set** In wait- $k$  systems, the counterintuitive phenomenon of the context usage between training and inference was first observed by Ma et al. (2019) yet without explanations: *in wait- $k$  inference, model trained consistently with the same wait- $k$  setting is worse than the model trained with the wait- $k'$  setting*

Train \ Inference	Inference				
	$k=1$	$k=3$	$k=5$	$k=7$	$k=9$
$k'=1$	<u>19.10</u>	18.06	17.42	16.94	16.80
$k'=3$	19.29	<u>23.76</u>	24.97	25.00	24.40
$k'=5$	20.33	<b>24.89</b>	26.36	26.93	27.27
$k'=7$	20.48	24.60	26.46	<u>27.26</u>	27.81
$k'=9$	<b>21.42</b>	24.82	<b>26.92</b>	<b>27.84</b>	<u>28.63</u>

Table 1: Evaluation by BLEU score on the valid set of the WMT15 De-En task for wait- $k$  policy. Bold: best in a column. Underline: training context is consistent with inference context. (§4 provides detailed settings.)

( $k' \neq k$ ) in terms of translation quality, as illustrated in Table 1.<sup>1</sup> For example, the BLEU score obtained by the model trained with wait-9 surpasses the model trained with wait-1 by a large margin with wait-1 inference. As a result, it has become a standard practice to utilize inconsistent context for training, and this practice is widely followed by (Elbayad et al., 2020; Zhang and Feng, 2021, 2022; Zhang et al., 2022; Guo et al., 2022, 2023), even if this phenomenon is counterintuitive due to the mismatch between training and inference on the usage of source-side context.

**Counterintuitive Phenomenon on Training Subset** One might hypothesize that this phenomenon is attributed to the generation issue from training data to valid data. To verify this hypothesis, we conduct similar experiments on a subset of the training data. We sample examples from the training data as a training subset with the same size as the valid set. Table 2 depicts that the situation

Train \ Inference	Inference				
	$k=1$	$k=3$	$k=5$	$k=7$	$k=9$
$k'=1$	<u>21.42</u>	21.21	21.00	20.25	19.67
$k'=3$	22.07	<u>25.51</u>	26.73	26.69	26.33
$k'=5$	22.53	25.55	<u>27.27</u>	28.06	28.07
$k'=7$	23.15	25.73	27.20	<u>28.34</u>	28.63
$k'=9$	<b>23.22</b>	<b>26.21</b>	<b>27.52</b>	<b>28.66</b>	<u>29.33</u>

Table 2: Evaluation by BLEU score on the training subset of the WMT15 De-En task for wait- $k$  policy.

on the training subset is almost similar to that on the valid set except for  $k = 3$ , where the optimal  $k' = 9$  for the training subset rather than  $k' = 5$  as for the valid set. This shows that generalization from training data to valid data is not the main reason for this counterintuitive phenomenon and it is non-trivial to analyze its reasons. Therefore, in the next subsection, we plan to investigate the reason for this phenomenon in depth.

<sup>1</sup>Actually, this observation focuses on the lower triangle.

## 2.2 Reasons of Counterintuitive Phenomenon

### Correlation between BLEU and Cross-entropy Loss in SiMT

Firstly, we explore the correlation between translation quality and training loss. To investigate correlation, we measure both the training loss and translation quality of each sample and calculate their Absolute Pearson Correlation in the training subset. In the majority of SiMT systems, the training objective is based on the cross-entropy objective. Therefore, we assess the training loss using cross-entropy loss score in our experiments. However, training loss is measured at the word level, while translation quality (BLEU score) is measured at the sentence level. To bridge this disparity, we compute the average training loss for each word within a sentence, thus representing it as sentence-level training loss.

$k$	1	3	5	7	9	$\infty$
<b>Entire</b>	0.62	0.70	0.73	0.74	0.75	0.75
<b>Low</b>	0.68	0.73	0.74	0.75	0.76	0.75
<b>High</b>	0.27	0.44	0.51	0.56	0.60	0.64

Table 3: Correlation between BLEU score and training (cross-entropy) loss on three subsets from the training subset of the WMT15 De-En task for wait- $k$  policy, where  $k = \infty$  means Full-sentence MT. **Entire** denotes the entire training subset, **Low** consists of those samples whose cross-entropy loss is lower than the averaged loss, **High** consists of those samples whose loss is higher than the averaged loss.

Table 3 presents the results of the correlation between BLEU and training (cross-entropy) loss in the wait- $k$  policy. We reveal the following insights. 1) In wait- $k$  systems, especially when  $k$  is smaller, the correlation is lower than that in Full-sentence MT. 2) When evaluating samples with high training (cross-entropy) loss, we observe a weaker correlation (between training loss and BLEU) compared to that with low training loss. This observation is not difficult to understand: taking a two-class classification task as an example, if the cross-entropy loss of an example is very high (e.g., the loss is  $-\log 0.2$ ), then the model can not predict the correct label for this example even if its loss is improved to  $-\log 0.4$ , because the probability of the ground-truth label is 0.4, which is less than 0.5. *This suggests the reason for the counterintuitive phenomenon on context usage is attributed to the relatively high cross-*

*entropy loss for SiMT,*<sup>2</sup> *leading to the weak correlation between training (cross-entropy) loss and translation quality.*

### Effects of Exposure Bias on the Models Trained Consistently and Inconsistently

Since the SiMT model is typically trained by cross-entropy loss, it suffers from the well-known exposure bias, i.e., during training, the model is only exposed to the training data distribution, instead of its predictions. Therefore, we propose to study the effects of exposure bias on the model trained with consistent context as well as the model trained with inconsistent context. To control the extent of exposure bias during inference stage, we measure translation quality by BLEU for both models (e.g., the former wait-1 inference model is trained with wait-1 setting and the later wait-1 inference model is trained with wait-9 setting) under the prefix-constrained decoding setting (Wuebker et al., 2016), where each model requires to predict the suffix for a given gold prefix. Under this setting, as the gold prefix gets shorter, more predicted tokens are used as the context during the prefix-decoding stage and the exposure bias is more severe.

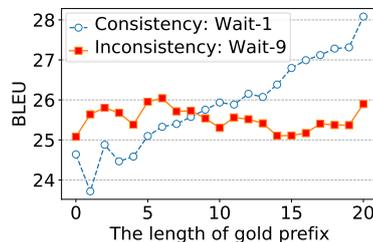


Figure 2: BLEU score comparison between context consistency and context inconsistency under the prefix-constrained decoding setting. The x-axis denotes the number of tokens for the gold prefix.

The results as presented in Figure 2 are averaged from a subset of 400 sentence pairs in the train set, all having the same number of tokens in the target (20 target tokens). It is evident that as the gold prefix becomes shorter (i.e., exposure bias is more severe) the performance of the consistent model significantly deteriorates, while the inconsistent model’s performance remains relatively better; however, when the number of tokens in the gold prefix is larger than 10 (i.e., exposure bias is less severe), the consistent model performs better. *This*

<sup>2</sup>Compared with full-sentence translation, SiMT uses less source-side context, which essentially results in a higher cross-entropy loss.

		Valid set					Training subset				
		<i>k</i> =1	<i>k</i> =3	<i>k</i> =5	<i>k</i> =7	<i>k</i> =9	<i>k</i> =1	<i>k</i> =3	<i>k</i> =5	<i>k</i> =7	<i>k</i> =9
Train	Inference										
	<i>k</i> '=1	<b>5.78</b>	5.26	5.00	4.87	4.81	<b>5.43</b>	5.11	4.95	4.87	4.83
	<i>k</i> '=3	5.78	<u>5.12</u>	4.79	4.61	4.53	5.48	<b>5.03</b>	4.83	4.73	4.67
	<i>k</i> '=5	5.81	<b>5.10</b>	<b>4.73</b>	4.53	4.42	5.54	5.06	<b>4.81</b>	4.69	4.61
	<i>k</i> '=7	5.86	5.12	4.72	<u>4.50</u>	4.38	5.60	5.09	4.82	<b>4.67</b>	4.59
	<i>k</i> '=9	5.91	5.14	4.72	<b>4.49</b>	<b>4.36</b>	5.65	5.12	4.84	4.68	<b>4.58</b>

Table 4: Evaluation by cross-entropy loss on valid set and training subset of WMT15 De-En task for wait-*k* policy.

finding reveals that one of the underlying causes of the counterintuitive phenomenon is attributed to exposure bias (Ranzato et al., 2016; Bengio et al., 2015; Zhang et al., 2019).

### 2.3 Counterintuitive Phenomenon Depends on Evaluation Metrics

The above reasons motivate us to study the counterintuitive phenomenon by using the cross-entropy loss for evaluation, in addition to BLEU as before, because training and inference criteria are the same, and there is no exposure bias issue in this case. We evaluate cross-entropy loss for the wait-*k* inference while models trained with wait-*k*' settings on the valid set and training subset, as illustrated in Table 4. On the valid set, we almost notice a diagonal trend, indicating the superiority of the consistent model. On the training subset, we observe a similar diagonal trend, indicating the counterintuitive phenomenon disappears in terms of cross-entropy loss as the evaluation metric. *These observation suggests that the counterintuitive phenomenon of context usage between training and inference depends on evaluation metrics, and it might be helpful to address this phenomenon by encouraging the consistent criterion between training and inference.*

## 3 Context Consistency Training for SiMT

Previous findings have shown that: 1) it is helpful to address the counterintuitive phenomenon by encouraging the consistent criterion between training and inference; 2) exposure bias is a reason for the counterintuitive phenomenon. To address the counterintuitive phenomenon and make the consistent model successful, we propose a simple and effective training approach, called context consistency training for SiMT, which not only incorporates the evaluation metrics for SiMT as training objectives (§3.1) but also allows the model

to expose its predictions during training (§3.2).

### 3.1 Bi-Objectives Optimization for SiMT

In SiMT, the evaluation metrics of models are translation quality and latency. Therefore, we intend to leverage both of these metrics as bi-objectives in our proposed method.

**Translation Quality** To measure the translation quality of SiMT models, we employ BLEU score (Papineni et al., 2002).

**Latency** Latency measurement is conducted using Average Lagging (AL) (Ma et al., 2019). AL quantifies the number of tokens of hypotheses that fall behind the ideal policy and is calculated as:

$$AL_g(\mathbf{x}, \mathbf{u}) = \frac{1}{\tau} \sum_{i=1}^{\tau} g(i, \mathbf{u}) - \frac{i-1}{|\mathbf{u}|/|\mathbf{x}|}, \quad (1)$$

where  $\tau = \operatorname{argmax}_i \{i \mid g(i) = |\mathbf{x}|\}$ ,  $\mathbf{x}$  is the source sentence,  $\mathbf{u}$  is the hypothesis sentence, and  $g(i)$  is the number of waited source tokens before translating  $\mathbf{u}_i$  and thus it is dependent on  $\mathbf{u}_{<i}$ , and its detailed definition depends on different read/write policies.

Formally, the SiMT model parametrized by  $\theta$  can be defined as follows:

$$p_g(\mathbf{u}|\mathbf{x}; \theta) = \prod_{i=1}^{|\mathbf{u}|} p(\mathbf{u}_i | \mathbf{x}_{\leq g(i)}, \mathbf{u}_{<i}), \quad (2)$$

where  $\mathbf{u}$  denotes a complete translation hypothesis and  $\mathbf{u}_{<i}$  denotes its partial prefix with  $i$  tokens.

Inspired by Minimum Risk Training (MRT) (Shen et al., 2016; Wieting et al., 2019), we directly optimize the SiMT model towards its bi-objectives (i.e., BLEU and Latency) as follows:

$$\mathcal{L}_g = \sum_{\mathbf{u} \in \mathcal{U}(\mathbf{x})} \operatorname{cost}_g(\mathbf{x}, \mathbf{y}, \mathbf{u}) \frac{p_g(\mathbf{u}|\mathbf{x}; \theta)}{\sum_{\mathbf{u}' \in \mathcal{U}(\mathbf{x})} p_g(\mathbf{u}'|\mathbf{x}; \theta)}, \quad (3)$$

where  $\mathcal{U}(\mathbf{x})$  is a set of candidate hypotheses,  $\mathbf{y}$  is the reference and  $\operatorname{cost}_g(\mathbf{x}, \mathbf{y}, \mathbf{u})$  consists of bi-objectives:

$$\text{cost}_g(\mathbf{x}, \mathbf{y}, \mathbf{u}) = \gamma \cdot \text{AL}_g(\mathbf{x}, \mathbf{u}) + (1 - \gamma) \cdot (1 - \text{BLEU}(\mathbf{y}, \mathbf{u})). \quad (4)$$

The hyperparameter  $\gamma$  is adjustable and allows us to fine-tune for different latency requirements.

**Remark** In Shen et al. (2016); Wieting et al. (2019), the cost is directly defined on a translation candidate  $\mathbf{u}$ , and thus it is trivial to calculate the cost for a given  $\mathbf{u}$ . However, in our scenario,  $\text{AL}_g(\mathbf{x}, \mathbf{u})$  depends not only on  $\mathbf{u}$  but also on  $g(i)$  specified by the read/write policy used in the SiMT system. As a result, during the training process, for each candidate  $\mathbf{u}$  generated via decoding, we access the SiMT model to incrementally compute the  $g(i)$  for all  $i$  and then compute  $\text{AL}_g(\mathbf{x}, \mathbf{u})$  based on all  $g(i)$  for  $\mathbf{u}$ .

### 3.2 Generating $n$ Candidates for Training SiMT

In the conventional training SiMT with cross-entropy loss, the decoding process does not consider multiple candidates. In our scenario, to calculate the objective function defined in (3), we need to generate a set of candidates  $\mathcal{U}$  via decoding which also allows the SiMT model to be exposed to the predictions and thereby mitigates exposure bias during the training stage. To this end, we try two different ways (Beam search and Sampling search) (Holtzman et al., 2020) to generate  $n$ -best candidates in SiMT. Beam search is a maximization-based decoding technique that optimizes output by favoring high-probability tokens. It is widely used in the generation of Full-sentence MT. Sampling search (Holtzman et al., 2020) is a stochastic decoding approach that samples from the top- $p$  portion of the probability distribution. This method excels in enhancing candidate diversity. In our experiments, we generate a set of 5-best candidates and select 0.8 for top- $p$  in the sampling search.

Furthermore, to calculate the  $\text{AL}_g(\mathbf{x}, \mathbf{u})$  of candidates defined in Eq. (1) which is dependent on the  $g(i)$ , we maintain both model score  $p_g$  as well as  $g(i)$  (the number of waited source words before translating  $\mathbf{u}_i$ ) at each timestep  $i$ . Specifically, during the decoding process, the SiMT model uses the value of  $g(i)$  to incrementally specify the source context and produce the next predictive distribution  $p_g$ . From this predictive distribution  $p_g$ , we select the top  $n$ -best (for beam search method) or sample  $n$  (for sampling method) partial candidates along with their respective  $g(i)$  values.

Following Edunov et al. (2018); Wieting et al. (2019), we employ the two-step training paradigm to train SiMT to speed up the training process: we first train the SiMT model with the standard cross-entropy loss, and then, in our context consistency training, we fine-tune the model by optimizing the bi-objectives (translation quality and latency) with the generated  $n$ -best candidates. It is worth noting that we only generate  $n$  candidates in the training stage, but in the inference stage the greedy search is used because of the essence of SiMT.

## 4 Experiments

### 4.1 System Settings

The proposed approach is evaluated on three widely used SiMT benchmarks, including IWSLT14 German  $\rightarrow$  English (De-En), IWSLT15 Vietnamese  $\rightarrow$  English (Vi-En) and German  $\rightarrow$  English (De-En). Experiments are conducted on SiMT systems including two different policies: The fixed read/write system (**wait- $k$  policy**) (Ma et al., 2019); The adaptive read/write system (**wait-info policy**) (Zhang et al., 2022).

**Baseline Training Approaches** The conventional training approach of SiMT systems is the context consistency training based on cross-entropy Ma et al. (2019), denoted **Consistency-CE**. In contrast, context inconsistency training, also based on cross-entropy, involves inconsistent context usage between training and inference stages, called **Inconsistency-CE**. Additionally, we implement a recently widely-used special case of context inconsistency training (Elbayad et al., 2020), termed **Inconsistency-CE-MP**.

**Our Training Approaches** The proposed SiMT systems follow the standard evaluation paradigm (Ma et al., 2019) and report BLEU scores (Papineni et al., 2002) for translation quality and Average Lagging (AL) (Ma et al., 2019) for latency mentioned in §3.1. The proposed context consistency training is based on bi-objectives, called **Consistency-Bi**, and we also implement the context consistency training based on BLEU as the uni-objective, called **Consistency-Uni** for further comparison. For generating  $n$  candidates, we implement Beam search in most cases, except the wait- $k$  policy, for which we utilize the Sampling search strategy. The implementation of all systems is based on Transformer in the Fairseq Library (Ott et al., 2019). Appendix A provides detailed experimental settings.

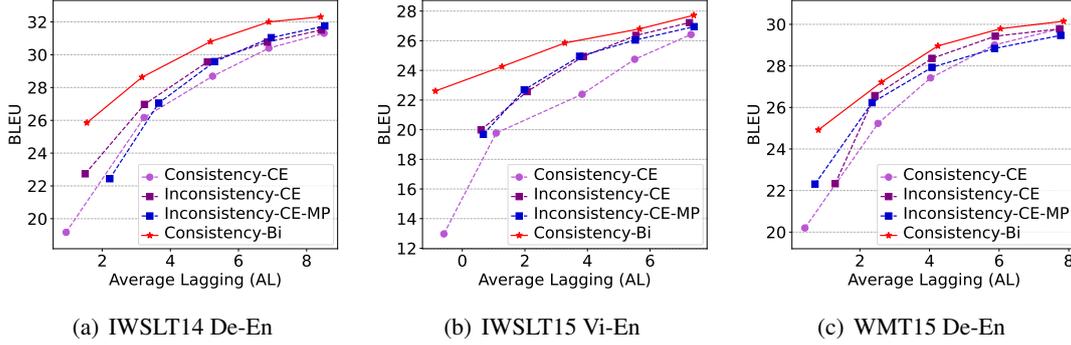


Figure 3: Translation quality (BLEU) v.s. latency (Average Lagging, AL) in Wait- $k$  Policy.

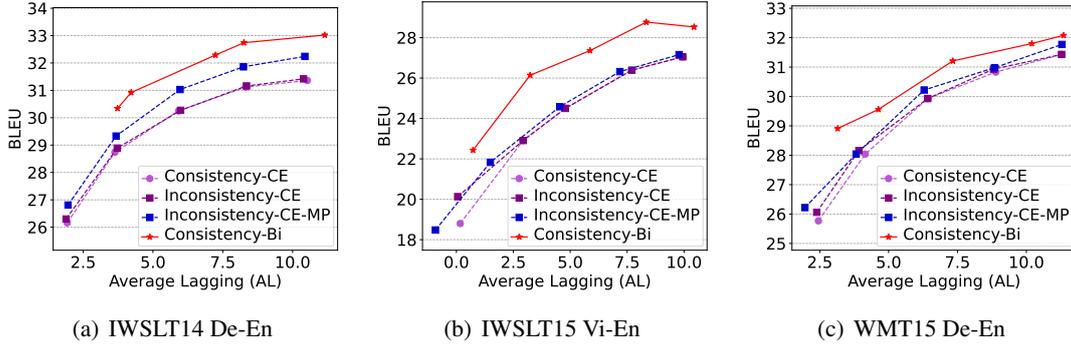


Figure 4: Translation quality (BLEU) v.s. latency (Average Lagging, AL) in Wait-info Policy.

## 4.2 Main Results

The results are illustrated in Figure 3 and Figure 4. Within our proposed context consistency training approach (Consistency-Bi), all implemented SiMT systems (wait- $k$  and wait-info) exhibit significant improvements in both translation quality and latency, as evidenced by an increase in BLEU score and a decrease in AL across all the benchmarks. This reveals that our proposed methods not only yield substantial performance improvements but also demonstrate strong generalization capabilities for SiMT systems.

In contrast to the original consistency training (Consistency-CE) of the wait- $k$  policy, our proposed Consistency-Bi achieves over 5 BLEU improvement at low latency ( $k=1$ ) across all datasets. Specifically, our method improves 2.68 BLEU on the IWSLT14 De-En task, 4.39 BLEU on the IWSLT15 Vi-En task, and 1.91 on the WMT15 De-En task, respectively (average on all latency). Furthermore, compared with inconsistency training (Inconsistency-CE and Inconsistency-CE-MP), the proposed method also demonstrates significant improvements, especially at low latency ( $k=1$ ), achieving over 3 BLEU score increases. This

suggests that incorporating our proposed context consistency training enables a wait- $k$  model trained consistently under the same wait- $k$  inference setting to outperform an inconsistently trained model.

To evaluate whether our method could achieve improvements with advanced adaptive SiMT systems, we apply our proposed training method to wait-info policy (Zhang et al., 2022). The results are depicted in Figure 4. Similarly, in comparison to the three baseline training methods, we observe a significant enhancement in translation quality across all latencies. However, in IWSLT15 Vi-En and WMT15 De-En tasks, Inconsistency-CE and Inconsistency-CE-MP are not significantly better than Consistency-CE. This can be attributed to the advanced policy, which makes more informed read/write decisions based on information.

## 4.3 Ablation Study

**Ablation Studies on Consistency-Bi and Consistency-Uni** To validate the effectiveness of Consistency-Bi, we perform the ablation studies on Consistency-Bi (Both BLEU and AL) and Consistency-Uni (BLEU only) in Figure

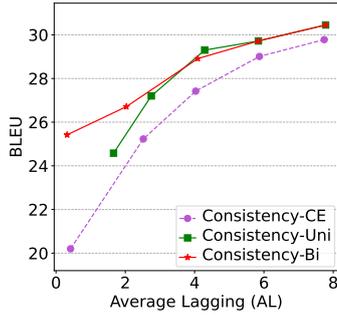


Figure 5: Ablation studies between Consistency-Bi and Consistency-Uni on WMT15 De-En test set of wait- $k$ .

5. The experiments reveal that compared with Consistency-Uni, Consistency-Bi not only results in lower latency but also yields superior translation quality, especially in low latency scenarios ( $k=1$ ), except for  $k=3$ , where Consistency-Uni is slightly better than Consistency-Bi. It is largely attributed to the latency as part of the training objectives.

**Ablation studies on  $n$ -best candidates generations** We conduct the ablation studies on two types of  $n$ -best generation methods (Beam search and Sampling search) under both wait- $k$  and wait-info policies, as depicted in Figure 6. The results reveal that under the wait- $k$  policy, the performance of Consistency-Bi using sampling search is slightly superior to that using beam search. Conversely, under the wait-info policy, employing beam search yields slightly better results compared to sampling search. These findings suggest the choice of generation method is not notably sensitive.

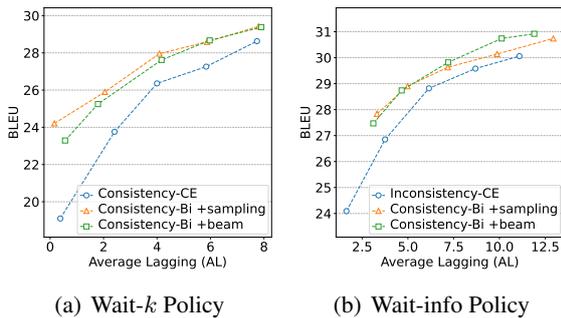


Figure 6: Ablation studies on  $n$ -best candidates generations (Beam search and Sampling search) on the valid set of WMT15 De-En.

**Variation in hyperparameter  $\gamma$**  Fine-tuning hyperparameter  $\gamma$  defined in (4) aims to achieve a better trade-off between BLEU and latency in our proposed Consistency-Bi. As illustrated in Table 5, as  $\gamma$  increases, AL decreases while the

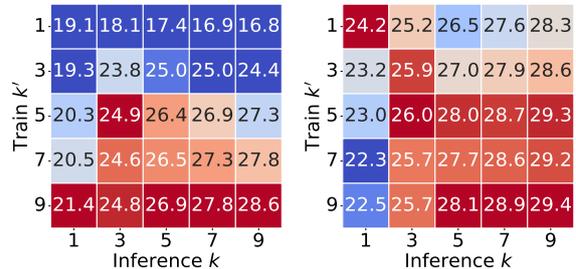
BLEU score improves, reaching its peak at  $\gamma = 0.4$ . This indicates that our proposed method can simultaneously optimize two objectives and achieve a value that is relatively optimally balanced between BLEU and AL, which can effectively enhance both translation quality and latency.

$\gamma$	0.0	0.1	0.2	0.3	0.4	0.5	0.6
BLEU	23.5	23.37	23.08	23.56	<b>24.21</b>	21.09	17.74
AL	1.68	1.62	1.53	1.14	0.16	-1.48	<b>-2.93</b>

Table 5: Ablation studies on various  $\gamma$  in wait-1 training with wait-1 inference of Consistency-Bi.

#### 4.4 Analysis

**Counterintuitive Phenomenon Mitigation** To explore whether the counterintuitive phenomenon described in §2.1 is alleviated, we conduct experiments using models trained with wait- $k'$  but tested with wait- $k$ , as illustrated in Figure 7. Figure 7(a) presents the results of the original training method. Optimal results for inference with  $k$  are generally achieved when  $k'=9$ , except for  $k=3$ , where  $k'=5$  yields the best. In contrast, our proposed training method demonstrates that the best results tested with wait- $k$  closely match with the diagonal line as depicted in Figure 7(b). Specifically, when inference with  $k=1$  and 9, the best results match the models trained with the same value of  $k'$ . For  $k=3, 5$ , and 7, although the best results come from different models, the differences are not significant. These findings suggest that our method exhibits improved consistency between training and inference compared with the original one.



(a) origin(train w/ CE-obj) (b) proposed(train w/ Bi-obj)

Figure 7: BLEU score comparison between the original and proposed training methods using wait- $k'$  during training and wait- $k$  during inference on the WMT15 De-En valid set. The diagonal line indicates consistency between training  $k'$  and inference  $k$ .

**Correlation between training loss and translation quality** We analyze the correlation between BLEU score and training loss, similar to the analysis described in §2.2. The results shown in Figure 8 demonstrate that, compared with Consistency-CE, proposed Consistency-Bi exhibits a strong correlation between training loss and translation quality, even when using a small  $k$ .

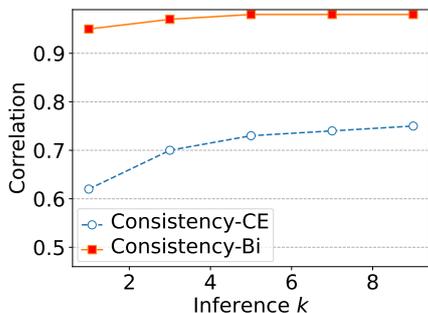


Figure 8: Comparison of correlation between BLEU score and training loss (cross-entropy loss for Consistency-CE and bi-objectives loss for Consistency-Bi) on training subset of WMT15 De-En task.

**Exposure Bias** To assess whether our method successfully mitigates exposure bias discussed in §2.2, we conduct wait-1 decoding experiments using both Consistency-CE and Consistency-Bi under the prefix-constrained decoding setting (Wuebker et al., 2016). The detailed experimental settings are as described in §2.2. Figure 9 reveals that as the number of gold prefixes decreases, the performance of Consistency-Bi improves, while the performance of Consistency-CE deteriorates. This suggests that the proposed method effectively mitigates exposure bias, enhancing the model’s performance when relying on prediction rather than on gold prefixes.

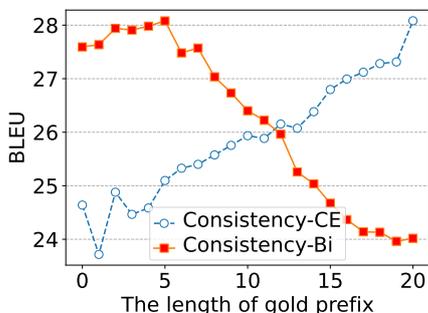


Figure 9: BLEU score comparison between original Consistency-CE model and ours proposed Consistency-Bi model for wait-1 decoding under the prefix-constrained decoding setting.

## 5 Related Work

Existing SiMT studies can be mainly categorized into two types (i.e., fixed or adaptive policy) according to the READ/WRITE policy.

As the fixed policy, Dalvi et al. (2018) introduced STATIC-RW, and Ma et al. (2019) proposed the wait- $k$  policy, which consistently generates target tokens lagging behind the source by  $k$  positions. Building upon this, Elbayad et al. (2020) enhanced the wait- $k$  policy by introducing the practice of sampling different values of  $k$  during training. Additionally, Han et al. (2020) incorporated meta-learning into the wait- $k$  policy, and Zhang et al. (2021) proposed future-guided training for the wait- $k$  policy.

Alternatively, many notable works develop an adaptive policy for SiMT (Zheng et al., 2019; Zhang et al., 2020; Wilken et al., 2020; Miao et al., 2021; Zhang and Feng, 2022; Zhang et al., 2022). For instance, Zheng et al. (2020) propose the adaptive policy through a heuristic ensemble of multiple wait- $k$  models. Other studies (Zheng et al., 2019; Arivazhagan et al., 2019; Ma et al., 2020; Zhang and Zhang, 2020; Zhang et al., 2020) resort to an adaptive policy controller to determine the READ/WRITE action and then integrate the controller into the SiMT model.

The above studies overlook the counterintuitive phenomenon about the context usage between training and inference, and our work thereby provides comprehensive analysis on this phenomenon and propose an effective approach to address this phenomenon, which is general enough to be applied into both policies.

## 6 Conclusion

This paper pays attention to a counterintuitive phenomenon in the context of usage between training and inference in SiMT. Subsequently, we conduct a comprehensive analysis and make the noteworthy discovery that this phenomenon primarily stems from the weak correlation between translation quality and training loss as well as exposure bias between training and inference. Based on our findings, we propose a context consistency training method that incorporates both translation quality and latency as bi-objectives and alleviates the exposure bias issue during the training. Experiments verify the effectiveness of the proposed approach, making the context-consistent SiMT successful for the first time.

## 560 Limitations

561 Our context consistency training approach neces-  
562 sitates a search for an appropriate hyperparameter,  
563 denoted as  $\gamma$ , to strike a balance between  
564 translation quality and latency. Further research  
565 is required to establish an efficient method for this  
566 purpose.

## 567 References

568 Naveen Arivazhagan, Colin Cherry, Wolfgang  
569 Macherey, Chung-Cheng Chiu, Semih Yavuz,  
570 Ruoming Pang, Wei Li, and Colin Raffel.  
571 2019. [Monotonic infinite lookback attention for  
572 simultaneous machine translation](#). In *Proceedings  
573 of the 57th Annual Meeting of the Association  
574 for Computational Linguistics*, pages 1313–1323,  
575 Florence, Italy. Association for Computational  
576 Linguistics.

577 Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam  
578 Shazeer. 2015. [Scheduled sampling for sequence  
579 prediction with recurrent neural networks](#). In  
580 *Advances in Neural Information Processing Systems  
581 28: Annual Conference on Neural Information  
582 Processing Systems 2015, December 7-12, 2015,  
583 Montreal, Quebec, Canada*, pages 1171–1179.

584 Chris Callison-Burch, Philipp Koehn, Christof Monz,  
585 and Josh Schroeder. 2009. [Findings of the 2009  
586 Workshop on Statistical Machine Translation](#). In  
587 *Proceedings of the Fourth Workshop on Statistical  
588 Machine Translation*, pages 1–28, Athens, Greece.  
589 Association for Computational Linguistics.

590 Mauro Cettolo, Jan Niehues, Sebastian Stüker, Luisa  
591 Bentivogli, and Marcello Federico. 2013. [Report  
592 on the 10th IWSLT evaluation campaign](#). In  
593 *Proceedings of the 10th International Workshop  
594 on Spoken Language Translation: Evaluation  
595 Campaign*, Heidelberg, Germany.

596 Kyunghyun Cho and Masha Esipova. 2016. [Can neural  
597 machine translation do simultaneous translation?](#)  
598 *ArXiv preprint*, abs/1606.02012.

599 Fahim Dalvi, Nadir Durrani, Hassan Sajjad, and  
600 Stephan Vogel. 2018. Incremental decoding and  
601 training methods for simultaneous translation in  
602 neural machine translation. In *Proceedings of the  
603 2018 Conference of the North American Chapter  
604 of the Association for Computational Linguistics:  
605 Human Language Technologies, Volume 2 (Short  
606 Papers)*, pages 493–499.

607 Sergey Edunov, Myle Ott, Michael Auli, David  
608 Grangier, and Marc’Aurelio Ranzato. 2018. [Clas-  
609 sical structured prediction losses for sequence to  
610 sequence learning](#). In *Proceedings of the 2018  
611 Conference of the North American Chapter of the  
612 Association for Computational Linguistics: Human  
613 Language Technologies, Volume 1 (Long Papers)*,

pages 355–364, New Orleans, Louisiana. Association  
for Computational Linguistics.

Maha Elbayad, Laurent Besacier, and Jakob Verbeek.  
2020. [Efficient wait-k models for simultaneous  
machine translation](#). In *Interspeech 2020, 21st  
Annual Conference of the International Speech  
Communication Association, Virtual Event, Shanghai,  
China, 25-29 October 2020*, pages 1461–1465.  
ISCA.

Jiatao Gu, Graham Neubig, Kyunghyun Cho, and  
Victor O.K. Li. 2017. [Learning to translate in real-  
time with neural machine translation](#). In *Proceedings  
of the 15th Conference of the European Chapter  
of the Association for Computational Linguistics:  
Volume 1, Long Papers*, pages 1053–1062, Valencia,  
Spain. Association for Computational Linguistics.

Shoutao Guo, Shaolei Zhang, and Yang Feng. 2022.  
[Turning fixed to adaptive: Integrating post-evaluation  
into simultaneous machine translation](#). In *Findings  
of the Association for Computational Linguistics:  
EMNLP 2022*, pages 2264–2278, Abu Dhabi, United  
Arab Emirates. Association for Computational  
Linguistics.

Shoutao Guo, Shaolei Zhang, and Yang Feng. 2023.  
[Glancing future for simultaneous machine translation](#).  
*ArXiv preprint*, abs/2309.06179.

Hou Jeung Han, Mohd Abbas Zaidi, Sathish Reddy  
Indurthi, Nikhil Kumar Lakumarapu, Beomseok Lee,  
and Sangha Kim. 2020. End-to-end simultaneous  
translation system for iwslt2020 using modality  
agnostic meta-learning. In *Proceedings of the  
17th International Conference on Spoken Language  
Translation*, pages 62–68.

Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and  
Yejin Choi. 2020. [The curious case of neural text  
degeneration](#). In *8th International Conference on  
Learning Representations, ICLR 2020, Addis Ababa,  
Ethiopia, April 26-30, 2020*. OpenReview.net.

Taku Kudo and John Richardson. 2018. [SentencePiece:  
A simple and language independent subword  
tokenizer and detokenizer for neural text processing](#).  
In *Proceedings of the 2018 Conference on Empirical  
Methods in Natural Language Processing: System  
Demonstrations*, pages 66–71, Brussels, Belgium.  
Association for Computational Linguistics.

Minh-Thang Luong and Christopher Manning. 2015.  
[Stanford neural machine translation systems for  
spoken language domains](#). In *Proceedings of the  
12th International Workshop on Spoken Language  
Translation: Evaluation Campaign*, pages 76–79, Da  
Nang, Vietnam.

Mingbo Ma, Liang Huang, Hao Xiong, Renjie Zheng,  
Kaibo Liu, Baigong Zheng, Chuanqiang Zhang,  
Zhongjun He, Hairong Liu, Xing Li, Hua Wu,  
and Haifeng Wang. 2019. [STACL: Simultaneous  
translation with implicit anticipation and controllable](#)



784 *the Association for Computational Linguistics*,  
785 pages 4334–4343, Florence, Italy. Association for  
786 Computational Linguistics.

787 Baigong Zheng, Kaibo Liu, Renjie Zheng, Mingbo Ma,  
788 Hairong Liu, and Liang Huang. 2020. *Simultaneous*  
789 *translation policies: From fixed to adaptive*. In  
790 *Proceedings of the 58th Annual Meeting of the*  
791 *Association for Computational Linguistics*, pages  
792 2847–2853, Online. Association for Computational  
793 Linguistics.

794 Baigong Zheng, Renjie Zheng, Mingbo Ma, and  
795 Liang Huang. 2019. *Simpler and faster learning of*  
796 *adaptive policies for simultaneous translation*. In  
797 *Proceedings of the 2019 Conference on Empirical*  
798 *Methods in Natural Language Processing and the 9th*  
799 *International Joint Conference on Natural Language*  
800 *Processing (EMNLP-IJCNLP)*, pages 1349–1354,  
801 Hong Kong, China. Association for Computational  
802 Linguistics.

## 803 A Detailed Experimental Settings

804 We conduct experiments on the following datasets,  
805 which are the widely-used SiMT benchmarks.

806 **IWSLT14 German → English (De→En)** (Cet-  
807 tolo et al., 2013) we train on 160K pairs,  
808 develop on 7K held-out pairs, and test on TED  
809 dev2010+tst2010-2013 (6,750 pairs). Following  
810 the previous setting (Elbayad et al., 2020), all  
811 data is tokenized and lower-cased and we segment  
812 sequences using byte pair encoding (Sennrich et al.,  
813 2016) with 10K merge operations. The resulting  
814 vocabularies are of 8.8K and 6.6K types in German  
815 and English respectively.

816 **IWSLT15<sup>3</sup> Vietnamese → English**  
817 **(Vi→En)** (Luong and Manning, 2015) we  
818 train on 133K pairs, develop on TED tst2012  
819 (1,553 pairs), and test on TED tst2013 (1,268  
820 pairs). The corpus is simply tokenized by  
821 SentencePiece (Kudo and Richardson, 2018),  
822 resulting in 16K and 8K word vocabularies in  
823 English and Vietnamese respectively.

824 **WMT15<sup>4</sup> German → English**  
825 **(De→En)** (Callison-Burch et al., 2009) is a  
826 parallel corpus with 4.5M training pairs. We use  
827 newstest2013 (3003 pairs) as the dev set and  
828 newstest2015 (2169 pairs) as the test set. The  
829 corpus is simply tokenized by SentencePiece  
830 resulting in 32k shared word vocabularies.

831 The implementation of all systems is based on  
832 Transformer (Vaswani et al., 2017) and adapted  
833 from Fairseq Library (Ott et al., 2019). Follow-  
834 ing Ma et al. (2019); Elbayad et al. (2020), we

835 apply Transformer-Small (4 heads) for IWSLT15  
836 Vi-En and IWSLT14 De-En, Transformer-Base  
837 (8 heads) for WMT15 De-En. To avoid the  
838 recalculation of the encoder hidden states when  
839 a new source token is read, unidirectional  
840 encoder (Elbayad et al., 2020) is proposed to make  
841 each source token only attend to its previous words.

<sup>3</sup>[nlp.stanford.edu/projects/nmt/](http://nlp.stanford.edu/projects/nmt/)

<sup>4</sup>[www.statmt.org/wmt15/translation-task](http://www.statmt.org/wmt15/translation-task)