

Self-Modifying State Modeling for Simultaneous Machine Translation

Anonymous ACL submission

Abstract

Simultaneous Machine Translation (SiMT) generates target outputs while receiving stream source inputs and requires a read/write policy to decide whether to wait for the next source token or generate a new target token, whose decisions form a *decision path*. Existing SiMT methods, which learn the policy by exploring various decision paths in training, face inherent limitations. These methods not only fail to precisely optimize the policy due to the inability to accurately assess the individual impact of each decision on SiMT performance, but also cannot sufficiently explore all potential paths because of their vast number. Besides, building decision paths requires unidirectional encoders to simulate streaming source inputs, which impairs the translation quality of SiMT models. To solve these issues, we propose **Self-Modifying State Modeling (SM²)**, a novel training paradigm for SiMT task. Without building decision paths, SM² individually optimizes decisions at each state during training. To precisely optimize the policy, SM² introduces Self-Modifying process to independently assess and adjust decisions at each state. For sufficient exploration, SM² proposes Prefix Sampling to efficiently traverse all potential states. Moreover, SM² ensures compatibility with bidirectional encoders, thus achieving higher translation quality. Experiments show that SM² outperforms strong baselines. Furthermore, SM² allows offline machine translation models to acquire SiMT ability with fine-tuning.

1 Introduction

Simultaneous Machine Translation (SiMT) (Gu et al., 2017; Ma et al., 2019; Zhang et al., 2020) outputs translation while receiving the streaming source sentence. Different from normal Offline Machine Translation (OMT) (Vaswani et al., 2017), SiMT needs a suitable read/write policy to decide whether to wait for the coming source inputs (READ) or generate target tokens (WRITE).

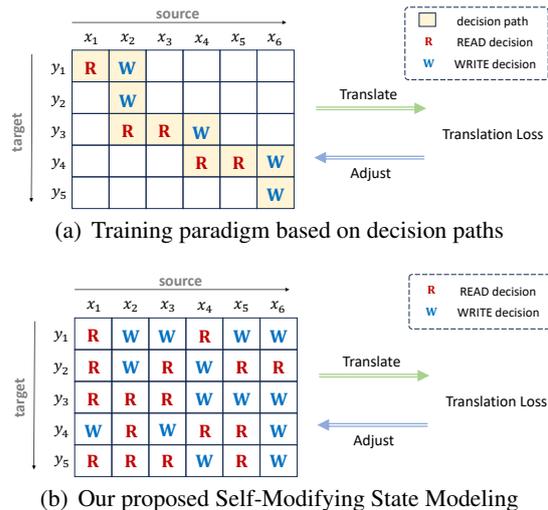


Figure 1: Illustration of different paradigms. (a) Training paradigm based on decision paths. All decisions along a path are optimized in an integrated manner. (b) Self-Modifying State Modeling. The decisions at each state are optimized individually.

As shown in Figure 1(a), to learn a suitable policy, existing SiMT methods usually require building a *decision path* (i.e., a series of READ and WRITE decisions made by the policy) to simulate the complete SiMT process during training (Zhang and Feng, 2022b). Methods of fixed policies (Ma et al., 2019; Zhang and Feng, 2021) build the decision path based on pre-defined rules, and only optimize translation quality along the path. Methods of adaptive policies (Zheng et al., 2019; Miao et al., 2021; Zhang and Feng, 2023) dynamically build the decision path and optimize the policy based on the SiMT performance along this path.

However, the current training paradigm based on decision paths faces inherent limitations. First, it can lead to **unprecise optimization** of the policy during training. For fixed policies, pre-defined rules cannot ensure optimal decisions at each state. For adaptive policies, there exists a credit assign-

ment problem (Minsky, 1961), which means it is difficult to identify the impact of each individual decision on the global SiMT performance along a path, thus hindering the precise optimization of each decision. Second, due to numerous potential decision paths, existing methods (Zheng et al., 2019; Miao et al., 2021; Zhang and Feng, 2023) often prohibit the exploration of some paths during training, but this **insufficient exploration** cannot ensure the optimal policy. Third, for building decision paths in training, existing methods require **unidirectional encoders** to simulate streaming source inputs and avoid the leakage of source future information (Elbayad et al., 2020), which impairs SiMT models’ translation quality (Iranzo-Sánchez et al., 2022; Kim and Cho, 2023).

To address these issues, we propose **Self-Modifying State Modeling (SM²)**, a novel training paradigm for SiMT task. As shown in Figure 1(b), instead of constructing complete decision paths, SM² individually optimizes decisions at all potential states during training. This paradigm necessitates addressing two critical issues: firstly, how to independently optimize each decision based on its own contribution to SiMT performance; and secondly, how to sufficiently explore all potential states during training. To realize the independent optimization, SM² assesses each decision by estimating confidence values which measure the translation credibility. High confidence means the SiMT model can predict a credible target token at current state and WRITE is beneficial for SiMT performance; otherwise, READ is preferred. Since golden confidence values are unavailable, SM² introduces **Self-Modifying** process to learn accurate confidence estimation (DeVries and Taylor, 2018; Lu et al., 2022). Specifically, during training, the SiMT model is allowed to modify its prediction based on the received source prefix with the prediction based on the complete source sentence, and the confidence is estimated to determine whether the modification is necessary to ensure a credible prediction at current state. To sufficiently explore all potential states, SM² conducts **Prefix Sampling** to divide all states into groups according to the number of their received source prefix tokens, and sample one group for optimization in each iteration.

Compared to the training paradigm based on decision paths, SM² presents significant advantages. First, the Self-Modifying process can assess each decision independently, which realizes the precise optimization of policy without the credit assign-

ment problem. Second, Prefix Sampling ensures sufficient exploration of all potential states, promoting the discovery of the optimal policy. These benefits enable SM² to learn a more effective policy. Furthermore, without building decision paths in training, SM² ensures compatibility with bidirectional encoders, thereby improving translation quality. This compatibility also allows OMT models to acquire the SiMT capability via fine-tuning.

Our contributions are outlined in the following:

- We propose **Self-Modifying State Modeling (SM²)**, a novel training paradigm that individually optimizes decisions at all states without building complete decision paths.
- SM² can learn a better policy through precise optimization of each decision and sufficient exploration of all states. With bidirectional encoders, SM² achieves higher translation quality and compatibility with OMT models.
- Experimental results on Zh→En, De→En and En→Ro SiMT tasks show that SM² outperforms strong baselines under all latency levels.

2 Background

Simultaneous machine translation For SiMT task, we respectively denote the source sentence as $\mathbf{x} = (x_1, \dots, x_M)$ and the corresponding target sentence as $\mathbf{y} = (y_1, \dots, y_N)$. Since the source inputs are streaming, we denote the number of source tokens available when generating y_i as g_i , and hence the prediction probability of y_i is $p(y_i | \mathbf{x}_{\leq g_i}, \mathbf{y}_{< i})$ (Ma et al., 2019). Thus, the decoding probability of \mathbf{y} is given by:

$$p(\mathbf{y} | \mathbf{x}) = \prod_{i=1}^N p(y_i | \mathbf{x}_{\leq g_i}, \mathbf{y}_{< i}) \quad (1)$$

Decision state and decision path We define the state s_{ij} as the condition in which the source prefix $\mathbf{x}_{\leq j}$ has been received and the target prefix $\mathbf{y}_{< i}$ has been generated. At s_{ij} , a decision $d_{ij} \in \{\text{WRITE}, \text{READ}\}$ can be made based on the context $(\mathbf{x}_{\leq j}, \mathbf{y}_{< i})$ (Zhao et al., 2023). Specifically, if $\mathbf{x}_{\leq j}$ is sufficient for the SiMT model to predict y_i accurately, d_{ij} should be WRITE; otherwise, d_{ij} should be READ. As shown in Figure 1(a), a series of decisions $[d_{00}, \dots, d_{NM}]$ are made in the SiMT process, which forms a decision path from s_{00} to s_{NM} . Along the decision path, the SiMT model can finish reading the whole \mathbf{x} and outputting the complete \mathbf{y} (Zhang and Feng, 2022b).

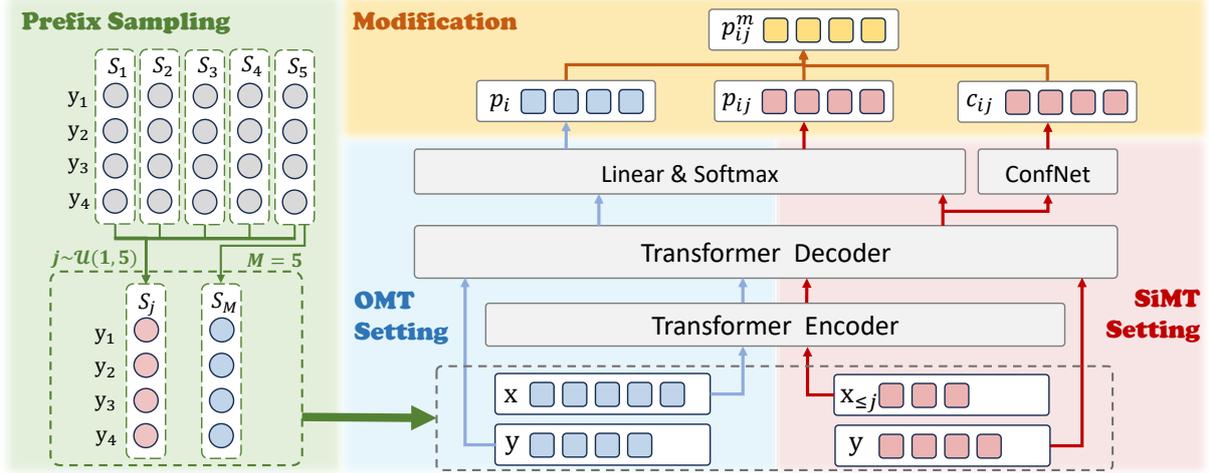


Figure 2: Overview of SM^2 . S_j contains the states where $\mathbf{x}_{\leq j}$ is received. S_M contains the states where complete \mathbf{x} is received. We introduce a confidence net (ConfNet) to estimate the confidence of each state. The model parameters in SiMT setting and OMT setting are shared. In this figure, the sentence lengths of the source and target sides are set to $M = 5$ and $N = 4$ respectively, and $j = 3$ in the Prefix Sampling step.

3 The Proposed Method

We propose **Self-Modifying State Modeling** (SM^2), which individually optimizes decisions at all states. The overview of SM^2 is shown in Figure 2. To independently optimize each decision, SM^2 learns confidence estimation to assess decisions at each state by modeling the Self-Modifying process (Sec.3.1). To ensure sufficient exploration during training, SM^2 conducts Prefix Sampling to traverse all potential states (Sec. 3.2). Then, based on estimated confidence at each state, SM^2 can determine whether the received source tokens are sufficient to generate a credible token and make suitable decisions during inference (Sec.3.3).

3.1 Self-Modifying for Confidence Estimation

Intuitively, when a translation model has access to the complete input \mathbf{x} (i.e., OMT setting), it can produce credible outputs. Therefore, a prediction made by the translation model at s_{ij} (i.e. SiMT setting) is considered credible if it aligns with that in OMT setting. Conversely, if the prediction in SiMT setting is incredible, it will be modified in OMT setting. Based on this insight and *Ask For Hints* (DeVries and Taylor, 2018; Lu et al., 2022), we model the Self-Modifying process to assess the translation credibility of each state. Specifically, we provide the SiMT model an option to modify its prediction in SiMT setting with that in OMT setting, and confidence estimation is defined as a binary classification determining whether the current generation requires the modification to ensure

a credible prediction. Through measuring translation credibility, decisions at each state can be independently assessed. High confidence means the SiMT model can generate a credible token at s_{ij} without modification and the WRITE decision is beneficial for SiMT performance; whereas low confidence indicates the prediction is inaccurate at s_{ij} and the READ decision is preferred.

During training, the Self-Modifying process is conducted in two steps: *prediction in SiMT setting & OMT setting* and *confidence-based modification*.

For *prediction in SiMT setting & OMT setting*, the SiMT model outputs different predictions at s_{ij} in SiMT setting and OMT setting respectively. These predictions are calculated as follows:

$$\begin{aligned} p_{ij} &= p(y_i | \mathbf{x}_{\leq j}, \mathbf{y}_{< i}) \\ p_i &= p(y_i | \mathbf{x}, \mathbf{y}_{< i}) \end{aligned} \quad (2)$$

It is noted that the model parameters in SiMT setting and OMT setting are shared.

For *confidence-based modification*, an additional confidence net is used to predict the confidence c_{ij} at s_{ij} . The confidence net is represented as:

$$c_{ij} = \text{sigmoid}(W^T \cdot h_{ij} + b) \quad (3)$$

where h_{ij} is the hidden representation from the top decoder layer in SiMT setting and $\theta = \{W, b\}$ are trainable parameters. If p_{ij} is credible, c_{ij} should be close to 1; otherwise, c_{ij} should be close to 0. To accurately calibrate c_{ij} in the training process, we integrate the modification into the prediction

Algorithm 1: Confidence-based Policy

Input : Streaming inputs $\mathbf{x}_{\leq j}$, Threshold γ , $i = 1, j = 1, y_0 \leftarrow \langle \text{BOS} \rangle$

Output : Target outputs \mathbf{y}

```
1 while  $y_{i-1} \neq \langle \text{EOS} \rangle$  do
2   calculate confidence  $c_{ij}$  as Eq.(3);
3   if  $c_{ij} \geq \gamma$  then // WRITE
4     generate  $y_i$  with  $\mathbf{x}_{\leq j}, \mathbf{y}_{< i}$ ;
5      $i \leftarrow i + 1$ ;
6   else // READ
7     wait for next source token  $x_{j+1}$ ;
8      $j \leftarrow j + 1$ ;
9   end
10 end
```

probability as follows:

$$p_{ij}^m = c_{ij} \cdot p_{ij} + (1 - c_{ij}) \cdot p_i \quad (4)$$

Subsequently, the translation loss is calculated using the modified probability:

$$\mathcal{L}_{s_{ij}} = -y_i \log(p_{ij}^m) \quad (5)$$

Notably, the SiMT model can enhance the prediction credibility by estimating a lower c_{ij} for more modification. However, this manner may cause an over-reliance on p_i . To avoid that, an additional penalty term for c_{ij} is introduced:

$$\mathcal{L}_{c_{ij}} = -\log(c_{ij}) \quad (6)$$

Through Self-Modifying process, SM^2 independently optimizes each decision based on their individual effect on the SiMT performance, thus realizing the precise optimization of the policy without credit assignment problem. We provide a gradient analysis of the independent optimization in Appendix A for further explanation.

3.2 Prefix Sampling

To sufficiently explore all potential states during training, Prefix Sampling is conducted in SM^2 . As shown in Figure 2, states are categorized into groups, and one group is randomly sampled for optimization in each iteration. Specifically, all possible states of (\mathbf{x}, \mathbf{y}) are divided into M groups according to the number of their received source prefix tokens, and each group comprises N states, which can be formulated as follows:

$$S_j = \{s_{ij} \mid 1 \leq i \leq N\}, j \in [1, M] \quad (7)$$

In each iteration, we sample $j \sim \mathcal{U}(1, M)$. Then, SM^2 respectively predicts target translation in SiMT setting based on S_j and those in OMT setting based on S_M , where the complete source sentence is received. Thus, the modified translation loss and the penalty item of each iteration are computed as follows:

$$\begin{aligned} \mathcal{L}_{S_j} &= \sum_{i=1}^N \mathcal{L}_{s_{ij}} \\ \mathcal{L}_{C_j} &= \sum_{i=1}^N \mathcal{L}_{c_{ij}} \end{aligned} \quad (8)$$

Besides, to ensure the p_i in OMT setting can provide effective modification, the translation loss in OMT setting is required, which is formulated as:

$$\mathcal{L}_{omt} = -\sum_{i=1}^N \log(p_i) \quad (9)$$

The total training loss is the following:

$$\mathcal{L} = \mathcal{L}_{omt} + \mathcal{L}_{S_j} + \lambda \mathcal{L}_{C_j} \quad (10)$$

where λ is the super parameter. We discuss the effect of λ in Appendix B.

Through Prefix Sampling, SM^2 explores all potential states without building any decision paths. Therefore, SM^2 can employ bidirectional encoders without the leakage of source future information in the training process.

3.3 Confidence-based Policy in Inference

During inference, SM^2 utilizes c_{ij} to assess the credibility of current prediction, thus making suitable decisions between READ and WRITE at s_{ij} . Specifically, a confidence threshold γ is introduced to serve as a criterion for making decisions. As shown in Algorithm 1, if $c_{ij} > \gamma$, SM^2 selects WRITE; otherwise, SM^2 selects READ. This decision process is constantly repeated until the complete translation is finished. It is noted that we only utilize SiMT setting in the inference process.

By adjusting γ , SM^2 can perform the SiMT task under different latency levels. A higher γ encourages the SiMT model to predict more credible target tokens and the latency will be longer. Conversely, a lower γ reduces the latency but may lead to a decrease in translation quality. The values of γ employed in our subsequent experiments are detailed in Appendix C.

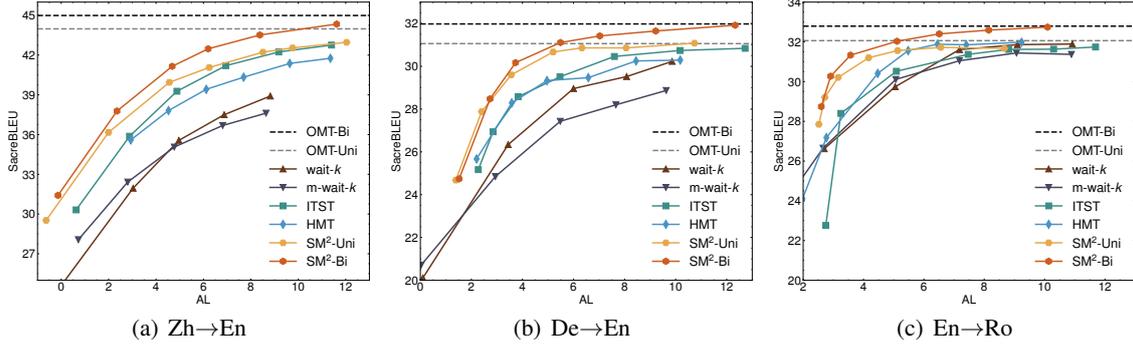


Figure 3: SacreBLEU against Average Lagging (AL) on Zh→En, De→En and En→Ro.

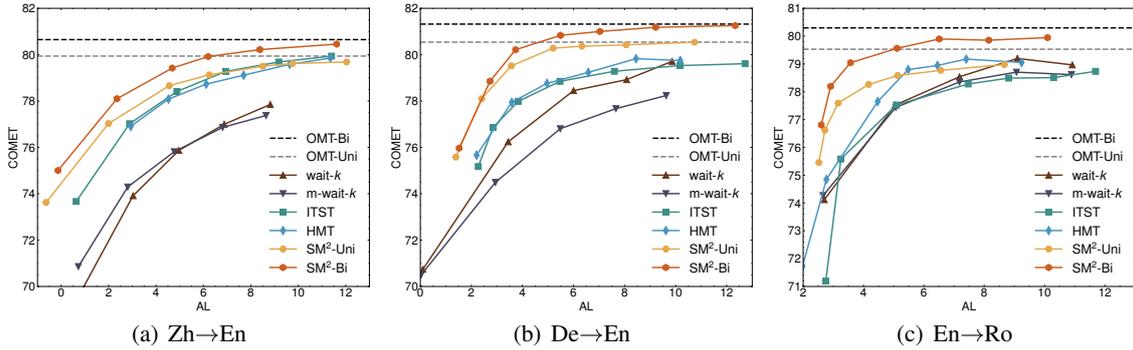


Figure 4: COMET against Average Lagging (AL) on Zh→En, De→En and En→Ro.

4 Experiments

4.1 Datasets

We conduct experiments on three datasets:

Zh→En We use LDC corpus which contains 2.1M sentence pairs as the training set, NIST 2008 for the validation set and NIST 2003, 2004, 2005, and 2006 for the test sets.

De→En We choose WMT15 for training, which contains 4.5M sentence pairs. Newstest 2013 are used as the validation set and newstest 2015 are used as the test set.

En→Ro WMT16 (0.6M) is used as the training set. We choose newsdev 2016 as the validation set and newstest 2016 as the test set.

We apply BPE (Sennrich et al., 2016) for all language pairs. In Zh→En, the vocabulary size is 30k for Chinese and 20k for English. In both De→En and En→Ro, a shared vocabulary is learned with 32k merge operations.

4.2 System Settings

The models used in our experiments are introduced as follows. All baselines are built based on Transformer (Vaswani et al., 2017) with the unidirec-

tional encoder unless otherwise stated. More details are presented in Appendix C.

OMT-Uni/OMT-Bi (Vaswani et al., 2017): OMT model with an unidirectional/bidirectional encoder.

wait-*k* (Ma et al., 2019): a fixed policy, which first reads *k* tokens, then writes one token and reads one token in turns.

m-wait-*k* (Elbayad et al., 2020): a fix policy, which improves wait-*k* by randomly sampling different *k* during training.

ITST (Zhang and Feng, 2022a): an adaptive policy, which models the SiMT task as a transport problem of information from source to target.

HMT (Zhang and Feng, 2023): an adaptive policy, which models the SiMT task as a hidden Markov model, by treating the states as hidden events and the predicted tokens as observed events.

SM²-Uni/SM²-Bi: Our proposed method with an unidirectional/bidirectional encoder.

4.3 Evaluation Metric

For SiMT, both translation quality and latency require evaluation. Since existing datasets mainly focus on the OMT task, the metric based on n-gram may cause inaccurate evaluation (Rei et al., 2020). Therefore, we measure the translation quality with

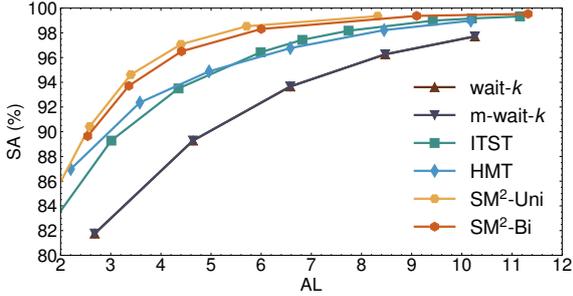


Figure 5: Evaluation of different SiMT policies. We calculate SA (\uparrow) under different latency levels.

Correlation Coefficient	Pearson	Spearman	Kendall's τ
Value	0.82	0.84	0.65

Table 1: Correlation between c_{ij} and p_{ij}^c .

both SacreBLEU (Post, 2018) and COMET¹ scores. For latency evaluation, we choose Average Lagging (AL) (Ma et al., 2019) as the metric.

Furthermore, to assess the quality of read/write policy in different SiMT models, we follow Zhang and Feng (2022a) and Kim and Cho (2023) to use Satisfied Alignments (SA), the proportion of the ground-truth aligned source tokens received before translating. Specifically, when generating y_i , the number of received source tokens g_i should be no less than the golden-truth aligned source position a_i , so that the alignment between y_i and x_{a_i} can be satisfied in the SiMT process. Thus, SA(\uparrow) can be calculated as:

$$SA = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(a_i \leq g_i) \quad (11)$$

5 Results and Analysis

5.1 Simultaneous Translation Quality

We present the translation quality under various latency levels of different SiMT models in Figure 3 and Figure 4. These results indicate that SM² outperforms previous methods across three language pairs in terms of both SacreBLEU and COMET scores. With the unidirectional encoder, SM-Uni achieves higher translation quality compared to current state-of-the-art SiMT models (ITST, HMT) at low and medium latency levels ($AL \in [0, 6]$), and maintains comparable performance at high latency level ($AL \in [6, 12]$). We attribute this improvement to the effectiveness of learning a better policy

¹Unbabel/wmt22-cometkiwi-da

during training. Furthermore, with the superior capabilities of the bidirectional encoder, SM²-Bi outperforms previous SiMT models more significantly across all latency levels. All SiMT models with unidirectional encoders can approach the translation quality of OMT-Uni at high latency levels, but only SM²-Bi achieves similar performance to OMT-Bi as the latency increases. These experimental results prove that SM² achieves better performance than other SiMT methods for learning better policy and improving translation quality. Detailed numerical results are provided in Appendix D, supplemented with additional evidence demonstrating the robustness of SM² to sentence length variations.

5.2 Superiority of SM² in Learning Policy

To verify whether SM² can learn a more effective policy, we compare SA(\uparrow) under various latency levels of different SiMT models. Following Zhang and Feng (2022a) and Kim and Cho (2023), we conduct the analysis on RWTH², a De \rightarrow En alignment dataset. The results are presented in Figure 5. Compared with existing methods, both SM²-Uni and SM²-Bi receive more aligned source tokens before generating target tokens under the same latency. Especially at medium latency level ($AL \in [4, 6]$), SM² can receive about 8% more source tokens than fixed policies (wait- k , m-wait- k) and 3.6% more than adaptive policies (ITST, HMT). We attribute these improvements to the advantages of SM² in learning policy. Through precise optimization, SM² can make more suitable decisions at each state, which generates faithful translations once receiving sufficient source tokens and waits for more source inputs when the predicted tokens are incredible. With sufficient exploration, SM² can investigate all possible situations and reduce unnecessary latency in the SiMT process.

5.3 Precise Optimization for Each Decision

To validate whether the confidence-based policy is precisely optimized at each state, we examine the relationship between estimated confidence c_{ij} and the probability of the correct token y_i in the prediction, denoted as p_{ij}^c . Specifically, we employ SM² to decode the validation set in a teacher-forcing manner, calculating the c_{ij} and p_{ij}^c for all possible states. Subsequently, a correlation analysis is performed between c_{ij} and p_{ij}^c . The results in Table 1 demonstrate a strong correlation, evidenced by

²<https://www-i6.informatik.rwth-aachen.de/goldAlignment/>

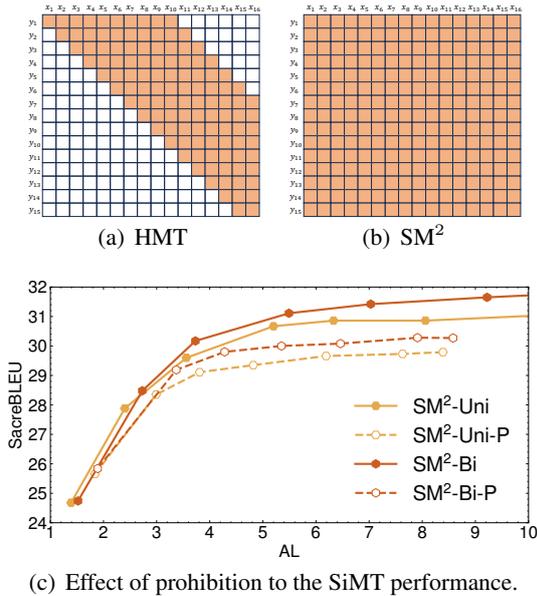
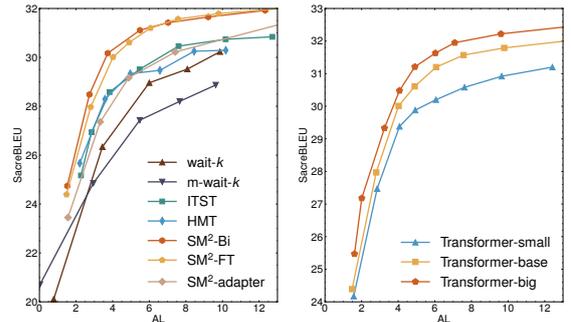


Figure 6: The visualization and effect of prohibition. In (a)(b), the shaded areas represent the states allowed for exploration in training. We apply the same prohibition in HMT (Zhang and Feng, 2023) to train SM²-Uni-P and SM²-Bi-P.

high values in Pearson (0.82) and Spearman (0.84) coefficients, with a slightly moderate but significant Kendall’s τ coefficient (0.65). These results suggest a robust linear and monotonic relationship between c_{ij} and p_{ij}^c , indicating the capacity of c_{ij} to accurately assess the credibility of the current predicted token. Consequently, this confirms the effectiveness of the confidence-based policy in making precise decisions at each state.

5.4 Advantage of Sufficient Exploration

Existing methods often prohibit the exploration of some paths due to the possible decision paths being numerous (Zheng et al., 2019; Miao et al., 2021; Zhang and Feng, 2023). To investigate the impact of the prohibition on SiMT models and the superiority of SM² in sufficiently exploring all states, we attempt to train these methods without prohibition, but they fail to converge. Therefore, we analyze the impact by employing the same prohibition in HMT (Zhang and Feng, 2023) and RIL (Zheng et al., 2019) to train SM², which restricts SM² to explore states only between wait- k_1 and wait- k_2 paths in training. As shown in Figure 6(a), we set $k_1 = 1$ and $k_2 = 10$ in our experiments. The performances of SM² with prohibition (SM²-Uni-P and SM²-Bi-P) are shown in Figure 6(c), indicating a decline in performance. These results suggest



(a) Comparison with SiMT (b) Comparison between different OMT models.

Figure 7: The SiMT performance of different OMT models after fine-tuning according to SM².

OMT model	Parameters	SacreBLEU	
		before FT	after FT
Transformer-small	47.9M	30.86	31.33
Transformer-base	60.5M	31.93	31.87
Transformer-big	209.1M	32.99	32.75

Table 2: The OMT performance of different OMT models before/after fine-tuning according to SM².

that the prohibition causes insufficient exploration, leading to diminished performance. In contrast, SM² ensures comprehensive exploration, which is shown in Figure 6(b), thereby achieving higher performance. Further analysis of the policy quality is provided in Appendix E.

5.5 Compatibility with OMT Models

SM² allows for the parallel training of the bidirectional encoder. Due to this compatibility, SM²-Bi achieves superior translation quality than existing SiMT methods with unidirectional encoders (Figure 3,4). To further present the superiority of this compatibility, we propose fine-tuning OMT models according to SM², so that the translation ability in OMT models can be easily utilized to gain SiMT models. Specifically, two distinct methods are used: fine-tuning all model parameters (SM²-FT) and fine-tuning with adapters (SM²-adapter)³. As shown in Figure 7(a), SM²-adapter can achieve comparable performance with current state-of-the-art SiMT models, and SM²-FT closely matches the performance of SM²-Bi.

Additionally, we further explore the effect of the OMT models’ translation abilities on the cor-

³We add adapters after the feed-forward networks of each encoder and decoder layer. For each adapter, the input dimension and output dimension are 512, and the hidden layer dimension is 128.

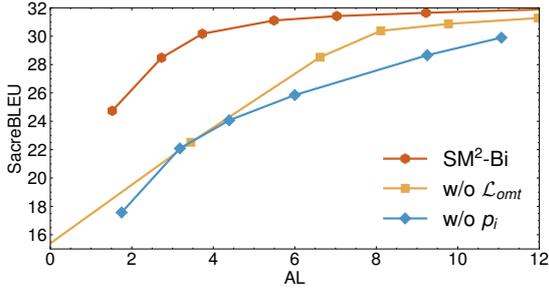


Figure 8: Effect of \mathcal{L}_{omt} and modification from OMT setting on the SM². "w/o \mathcal{L}_{omt} " is SM² trained without \mathcal{L}_{omt} , and "w/o p_i " means SM² trained using one-hot rather than OMT setting for modification.

Model	SM ² -Bi	w/o \mathcal{L}_{omt}	w/o p_i	OMT
SacreBLEU	31.87	30.33	31.95	31.93

Table 3: Effect of \mathcal{L}_{omt} and p_i on the OMT ability.

responding SiMT abilities after fine-tuning. We conduct the full-parameter fine-tuning on OMT models with Transformer-small, Transformer-base, and Transformer-big respectively. The OMT and SiMT capabilities of these models are illustrated in Table 2 and Figure 7(b), which reveal that models with stronger OMT abilities achieve better SiMT performance after fine-tuning. Besides, the results in Table 2 show that these models' original OMT abilities are not hurt, indicating that SM² enables models to support both OMT and SiMT abilities.

5.6 Ablation Study

We conduct ablation studies on SM² to analyze the effect of \mathcal{L}_{omt} and modification from OMT setting.

Effect of \mathcal{L}_{omt} As shown in Figure 8, the SiMT model without \mathcal{L}_{omt} drops quickly. We argue this is because training without \mathcal{L}_{omt} may cause a worse modification. The results in Table 3 show that the OMT performance of SM² trained without \mathcal{L}_{omt} is significantly affected, even worse than its SiMT performance in the high latency levels. This poor OMT ability cannot provide accurate modification, thus disrupting the policy learning process.

Effect of OMT modification Following *Ask For Hints* (DeVries and Taylor, 2018; Lu et al., 2022), we use the one-hot label as the "hints" to modify the prediction in SiMT setting. Specifically, we denote t_i as the ground-truth label of the i -th target token, and hence the modification in SM² is adjusted as:

$$p_{ij}^m = c_{ij} \cdot p_{ij} + (1 - c_{ij}) \cdot t_i \quad (12)$$

As shown in Figure 8, the performance of SM² trained with modification in Eq.(12) also drops.

We argue that the modification from t_i cannot reflect the real available gain from the modification after receiving the complete source sentence, thus learning a worse policy.

6 Related Work

Simultaneous Machine Translation Existing SiMT methods are divided into fixed policy and adaptive policy. For fixed policy, Ma et al. (2019) proposed wait- k , which starts translation after receiving k tokens. Elbayad et al. (2020) proposed multipath wait- k , which randomly samples k during training. For adaptive policy, heuristic rules Cho and Esipova (2016) and reinforcement learning Gu et al. (2017) are used to realize the SiMT task. Ma et al. (2020b) integrated monotonic attention to model the decision process. Miao et al. (2021) proposed a generative framework to learn a read/write policy. Zhang and Feng (2022a) measured the information SiMT had received and proposed an information-based policy. Zhang and Feng (2023) used the Hidden Markov model in SiMT task to learn an adaptive policy.

Previous methods based on decision paths are limited in policy learning and model structure. Our proposed SM² individually explores all states during training, overcoming these limitations.

Confidence Estimation for OMT Confidence estimation is used to measure the models' credibility. Wang et al. (2019) use Monte Carlo dropout to propose an uncertainty-based confidence estimation. Wan et al. (2020) utilize the confidence score to guide self-paced learning. DeVries and Taylor (2018) evaluates the confidence by measuring the level it asks for hints from the ground-truth label, and Lu et al. (2022) transfers it to OMT to improve the out-of-distribution detection.

7 Conclusion

In this paper, we propose Self-Modifying State Modeling (SM²), a novel training paradigm for SiMT task. Instead of constructing complete decision paths, SM² individually optimizes decisions at all potential states during training. Experiments on three language pairs show the superiority of SM² in terms of read/write policy, translation quality, and compatibility with OMT models.

8 Limitations

In this paper, we propose SM², a novel paradigm that individually optimizes decisions at each state.

545	Although our experiments show the superiority of	framework. In <i>Proceedings of the 57th Annual Meet-</i>	598
546	not building decision paths during training, there	<i>ing of the Association for Computational Linguistics</i> ,	599
547	are still some parts to be further improved. For ex-	pages 3025–3036.	600
548	ample, using a more effective way to independently	Shuming Ma, Dongdong Zhang, and Ming Zhou. 2020a.	601
549	assess the individual effect of each decision on the	A simple and effective unified encoder for document-	602
550	SiMT performance. Besides, how to leverage other	level machine translation. In <i>Proceedings of the 58th</i>	603
551	pre-trained encoder-decoder models like BART and	<i>annual meeting of the association for computational</i>	604
552	T5, to gain SiMT models, is still a promising di-	<i>linguistics</i> , pages 3505–3511.	605
553	rection to explore. These will be considered as	Xutai Ma, Juan Miguel Pino, James Cross, Liezl Pu-	606
554	objectives for our future work.	zon, and Jiatao Gu. 2020b. Monotonic multithread	607
555	Acknowledgements	attention. In <i>International Conference on Learning</i>	608
556		<i>Representations</i> .	609
557	References	Yishu Miao, Phil Blunsom, and Lucia Specia. 2021.	610
558	Kyunghyun Cho and Masha Esipova. 2016. Can neu-	A generative framework for simultaneous machine	611
559	ral machine translation do simultaneous translation?	translation. In <i>Proceedings of the 2021 Conference</i>	612
560	<i>arXiv e-prints</i> , pages arXiv–1606.	<i>on Empirical Methods in Natural Language Process-</i>	613
561	Terrance DeVries and Graham W Taylor. 2018. Learn-	<i>ing</i> , pages 6697–6706.	614
562	ing confidence for out-of-distribution detection in	Marvin Minsky. 1961. Steps toward artificial intelli-	615
563	neural networks. <i>arXiv preprint arXiv:1802.04865</i> .	gence. <i>Proceedings of the IRE</i> , 49(1):8–30.	616
564	Maha Elbayad, Laurent Besacier, and Jakob Verbeek.	Masato Neishi and Naoki Yoshinaga. 2019. On the	617
565	2020. Efficient wait-k models for simultaneous ma-	relation between position information and sentence	618
566	chine translation.	length in neural machine translation. In <i>Proceedings</i>	619
567	Jiatao Gu, Graham Neubig, Kyunghyun Cho, and Vic-	<i>of the 23rd Conference on Computational Natural</i>	620
568	tor OK Li. 2017. Learning to translate in real-time	<i>Language Learning (CoNLL)</i> , pages 328–338.	621
569	with neural machine translation. In <i>Proceedings of</i>	Matt Post. 2018. A call for clarity in reporting bleu	622
570	<i>the 15th Conference of the European Chapter of the</i>	scores. In <i>Proceedings of the Third Conference on</i>	623
571	<i>Association for Computational Linguistics: Volume</i>	<i>Machine Translation: Research Papers</i> , page 186.	624
572	<i>1, Long Papers</i> , pages 1053–1062.	Association for Computational Linguistics.	625
573	Javier Iranzo-Sánchez, Jorge Civera, and Alfons Juan.	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	626
574	2022. From simultaneous to streaming machine	Lavie. 2020. Comet: A neural framework for mt eval-	627
575	translation by leveraging streaming history. In <i>Pro-</i>	uation. In <i>Proceedings of the 2020 Conference on</i>	628
576	<i>ceedings of the 60th Annual Meeting of the Associa-</i>	<i>Empirical Methods in Natural Language Processing</i>	629
577	<i>tion for Computational Linguistics (Volume 1: Long</i>	<i>(EMNLP)</i> , pages 2685–2702.	630
578	<i>Papers)</i> , pages 6972–6985.	Rico Sennrich, Barry Haddow, and Alexandra Birch.	631
579	Xiaomian Kang, Yang Zhao, Jiajun Zhang, and	2016. Neural machine translation of rare words with	632
580	Chengqing Zong. 2020. Dynamic context selection	subword units. In <i>Proceedings of the 54th Annual</i>	633
581	for document-level neural machine translation via	<i>Meeting of the Association for Computational Lin-</i>	634
582	reinforcement learning. In <i>Proceedings of the 2020</i>	<i>guistics (Volume 1: Long Papers)</i> , pages 1715–1725.	635
583	<i>Conference on Empirical Methods in Natural Lan-</i>	Dušan Variš and Ondřej Bojar. 2021. Sequence length	636
584	<i>guage Processing (EMNLP)</i> , pages 2242–2254.	is a domain: Length-based overfitting in transformer	637
585	Kang Kim and Hankyu Cho. 2023. Enhanced simulta-	models. <i>arXiv preprint arXiv:2109.07276</i> .	638
586	neous machine translation with word-level policies.	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	639
587	<i>arXiv preprint arXiv:2310.16417</i> .	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	640
588	Yu Lu, Jiali Zeng, Jiajun Zhang, Shuangzhi Wu, and	Kaiser, and Illia Polosukhin. 2017. Attention is all	641
589	Mu Li. 2022. Learning confidence for transformer-	you need. <i>Advances in neural information processing</i>	642
590	-based neural machine translation. In <i>Proceedings</i>	<i>systems</i> , 30.	643
591	<i>of the 60th Annual Meeting of the Association for</i>	Yu Wan, Baosong Yang, Derek F Wong, Yikai Zhou,	644
592	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	Lidia S Chao, Haibo Zhang, and Boxing Chen. 2020.	645
593	pages 2353–2364.	Self-paced learning for neural machine translation.	646
594	Mingbo Ma, Liang Huang, Hao Xiong, Renjie Zheng,	<i>arXiv preprint arXiv:2010.04505</i> .	647
595	Kaibo Liu, Baigong Zheng, Chuanqiang Zhang,	Shuo Wang, Yang Liu, Chao Wang, Huanbo Luan, and	648
596	Zhongjun He, Hairong Liu, Xing Li, et al. 2019.	Maosong Sun. 2019. Improving back-translation	649
597	Stacl: Simultaneous translation with implicit antici-	with uncertainty-based confidence estimation. <i>arXiv</i>	650
	pation and controllable latency using prefix-to-prefix	<i>preprint arXiv:1909.00157</i> .	651

Ruiqing Zhang, Chuanqiang Zhang, Zhongjun He, Hua Wu, and Haifeng Wang. 2020. Learning adaptive segmentation policy for simultaneous translation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2280–2289.

Shaolei Zhang and Yang Feng. 2021. Universal simultaneous machine translation with mixture-of-experts wait-k policy. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7306–7317.

Shaolei Zhang and Yang Feng. 2022a. **Information-transport-based policy for simultaneous translation**. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 992–1013, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Shaolei Zhang and Yang Feng. 2022b. Modeling dual read/write paths for simultaneous machine translation. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2461–2477.

Shaolei Zhang and Yang Feng. 2023. Hidden markov transformer for simultaneous machine translation. *arXiv preprint arXiv:2303.00257*.

Libo Zhao, Kai Fan, Wei Luo, Wu Jing, Shushu Wang, Ziqian Zeng, and Zhongqiang Huang. 2023. Adaptive policy with wait-k model for simultaneous translation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4816–4832.

Baigong Zheng, Renjie Zheng, Mingbo Ma, and Liang Huang. 2019. Simultaneous translation with flexible policy via restricted imitation learning. *arXiv preprint arXiv:1906.01135*.

A Gradient Analysis

In this section, we provide a gradient analysis of the independent optimization in SM². The training loss function \mathcal{L} of SM² is formulated in Eq. (10). During training, this loss function adjusts each decision d_{ij} at state s_{ij} by changing the value of corresponding confidence c_{ij} . Specifically, the gradient of \mathcal{L} with respect to c_{ij} is calculated as:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial c_{ij}} &= \frac{\partial \mathcal{L}_{s_{ij}}}{\partial c_{ij}} + \lambda \frac{\partial \mathcal{L}_{c_{ij}}}{\partial c_{ij}} \\ &= -\frac{y_i}{p_{ij}^m} \cdot \frac{\partial p_{ij}^m}{\partial c_{ij}} - \frac{\lambda}{c_{ij}} \\ &= -\frac{y_i(p_{ij} - p_i)}{c_{ij} \cdot p_{ij} + (1 - c_{ij}) \cdot p_i} - \frac{\lambda}{c_{ij}} \end{aligned} \quad (13)$$

It is evident that this gradient does not contain any $c_{i'j'}$ ($i' \neq i$ or $j' \neq j$). Therefore, in the training

process, the estimated value of c_{ij} is adjusted only based on its current value and the prediction probability of the current state, without being affected by the decisions at other states, thus allowing for the independent optimization of c_{ij} .

In contrast, existing SiMT methods usually conduct training on decision paths and can not ensure independent optimization. Taking ITST (Zhang and Feng, 2022a) as an example, whose loss function \mathcal{L}' is formulated as:

$$\begin{aligned} \mathcal{L}' &= \mathcal{L}_{ce} + \mathcal{L}_{latency} + \mathcal{L}_{norm} \\ \mathcal{L}_{latency} &= \sum_{i=1}^I \sum_{j=1}^J T_{ij} \times C_{ij} \\ \mathcal{L}_{norm} &= \sum_{i=1}^I \left\| \sum_{j=1}^J T_{ij} - 1 \right\|_2 \end{aligned} \quad (14)$$

where \mathcal{L}_{ce} is the cross-entropy for learning translation ability, and C_{ij} is the latency cost for each state. During training, the decision is dominated by T_{ij} . The gradient of \mathcal{L}' with respect to T_{ij} is calculated as:

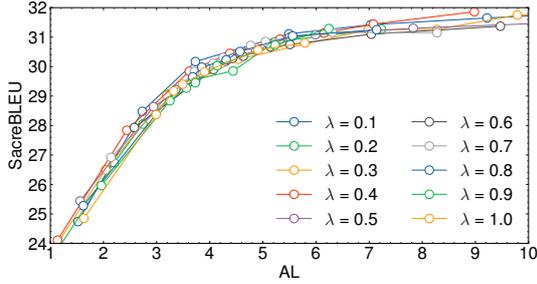
$$\frac{\partial \mathcal{L}'}{\partial T_{ij}} = \frac{\partial \mathcal{L}_{ce}}{\partial T_{ij}} + C_{ij} + 2 \left(\sum_{j=1}^J T_{ij} - 1 \right) \quad (15)$$

It is noted that the gradient of T_{ij} is also affected by the current values of $T_{ij'}$ ($j' = 1, 2, \dots, J$). These decisions are coupled in the optimization, thus not enabling the independent optimization of each decision. This can trigger mutual interference during training (Zhang and Feng, 2023) and lead to a credit assignment problem.

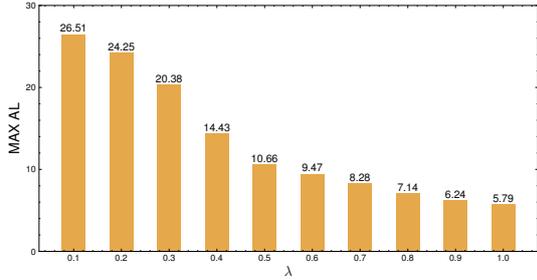
B Effect of λ

We analyze the effect of λ , which is the weight of the penalty during training. We train SM² with different λ ranging from 0.1 to 1, in increments of 0.1. As shown in Figure 9(a), the SM² models trained with different λ show comparable performance across all latency. This indicates that the performance of SM² is robust to variations in hyperparameters λ .

When λ becomes larger, the corresponding γ at the same latency will also increase. Therefore, we further analyze the effect of λ on the applicable latency range of SM². We denote the "MAX AL" as the latency of SM² when γ is set as 0.99 during inference. The results are shown in Figure 9(b). When λ becomes larger, "MAX AL" also decreases, which means a smaller applicable latency



(a) Performance of SM² with different λ .



(b) Max Latency of SM² with different λ .

Figure 9: Effect of λ on SM².

range. For example, when $\lambda = 1.0$ in training, it is hard for SM² to perform SiMT task under the latency levels where AL is larger than 5.79 since the threshold γ has been close to 1.

C Hyper-parameters

The system settings in our experiments are shown in Table 4. We set $\lambda = 0.1$ during training. Besides, we follow (Ma et al., 2020b) to use greedy search during inference for all baselines. The values of γ we used are 0.3,0.4,0.5,0.55,0.6,0.65 for Zh→En, 0.3,0.4,0.5,0.55,0.6,0.65,0.7 for De→En, and 0.3,0.4,0.5,0.6,0.65,0.7,0.75 for En→Ro.

D Main Results Supplement

D.1 Numerical Results

Table 5, 6, 7 respectively report the numerical results on LDC Zh→En, WMT15 De→En, MWT16 En→Ro measured by AL, SacreBLEU and COMET.

D.2 Robustness of SM² to Sentence Length

To validate the robustness of SM² to Sentence Length, we conduct additional experiments on De→En SiMT tasks. Specifically, we divide the test set into two groups based on sentence length: LONG group and SHORT group. The average lengths and the number of sentences in each group are shown in Table 8. Then, we test SM²-Bi and

Hyper-parameter	
encoder layers	6
encoder attention heads	8
encoder embed dim	512
encoder ffn embed dim	1024
decoder layers	6
decoder attention heads	8
decoder embed dim	512
decoder ffn embed dim	1024
dropout	0.1
optimizer	adam
adam- β	(0.9, 0.98)
clip-norm	1e-7
lr	5e-4
lr scheduler	inverse sqrt
warmup-updates	4000
warmup-init-lr	1e-7
weight decay	0.0001
label-smoothing	0.1
max tokens	8192

Table 4: Hyper-parameters of our experiments.

SM²-Uni separately on these two groups. The translation quality under different latency levels for SM²-Bi and SM²-Uni are presented in Figure 10. For clearer comparison, we also provide the performances of OMT models (OMT-Bi, OMT-Uni) on LONG and SHORT groups.

The results in Figure 10 indicate that when applied to longer sentences, the performance changes of SM² are similar to OMT models in both unidirectional and bidirectional encoder settings. Since the performance of OMT models unavoidably drops as the sentences become longer (Neishi and Yoshinaga, 2019; Kang et al., 2020; Ma et al., 2020a; Variš and Bojar, 2021), it is not SM² that triggers the decrease of translation quality. Therefore, SM² is still effective on long sentences.

E Effect of Prohibition on Policy

To further validate that the prohibition of exploration negatively affects the policy. We compare the SA of SM² with and without the prohibition on RWTH dataset. The results in Figure 11 indicate that the prohibition makes SM² learn a worse policy. Therefore, we can conclude that the prohibition will hurt the quality of policy. This further presents the advantage of SM² in sufficiently exploring all states through Prefix Sampling.

Chinese→English			
wait- k			
k	AL	SacreBLEU	COMET
1	-0.60	23.14	67.06
3	3.03	31.94	73.91
5	4.96	35.56	75.87
7	6.87	37.50	76.99
9	8.82	38.90	77.85
m-wait- k			
k	AL	SacreBLEU	COMET
1	0.72	28.06	70.85
3	2.80	32.41	74.29
5	4.76	35.05	75.81
7	6.81	36.68	76.86
9	8.64	37.61	77.37
HMT			
(L, K)	AL	SacreBLEU	COMET
(2,4)	2.93	35.59	76.90
(3,6)	4.52	37.81	78.08
(5,6)	6.11	39.41	78.73
(7,6)	7.69	40.33	79.11
(9,8)	9.64	41.37	79.58
(11,8)	11.35	41.75	79.85
ITST			
δ	AL	SacreBLEU	COMET
0.2	0.62	30.31	73.66
0.3	2.88	35.87	77.02
0.4	4.88	39.27	78.41
0.5	6.94	41.20	79.27
0.6	9.17	42.23	79.68
0.7	11.40	42.75	79.93
SM ² -Uni			
γ	AL	SacreBLEU	COMET
0.3	-0.63	29.52	73.62
0.4	1.99	36.16	77.02
0.5	4.56	39.94	78.66
0.55	6.24	41.06	79.13
0.6	8.51	42.21	79.50
0.65	9.75	42.54	79.61
SM ² -Bi			
γ	AL	SacreBLEU	COMET
0.3	-0.14	31.41	75.00
0.4	2.35	37.77	78.09
0.5	4.68	41.15	79.42
0.55	6.19	42.47	79.91
0.6	8.37	43.51	80.21
0.65	11.61	44.34	80.45

Table 5: Numerical results on LDC Zh→En.

German→English			
wait- k			
k	AL	SacreBLEU	COMET
1	0.10	20.11	70.74
3	3.44	26.34	76.24
5	6.00	28.96	78.44
7	8.08	29.52	78.92
9	9.86	30.23	79.71
m-wait- k			
k	AL	SacreBLEU	COMET
1	0.03	20.71	70.49
3	2.94	24.85	74.49
5	5.48	27.43	76.80
7	7.66	28.2	77.67
9	9.63	28.87	78.23
HMT			
(L, K)	AL	SacreBLEU	COMET
(2,4)	2.20	25.67	75.66
(3,6)	3.58	28.29	77.94
(5,6)	4.96	29.33	78.76
(7,6)	6.58	29.47	79.23
(9,8)	8.45	30.25	79.82
(11,8)	10.18	30.29	79.74
ITST			
δ	AL	SacreBLEU	COMET
0.2	2.27	25.17	75.17
0.3	2.85	26.94	76.86
0.4	3.83	28.58	77.98
0.5	5.47	29.51	78.85
0.6	7.60	30.46	79.28
0.7	10.17	30.74	79.53
0.8	12.72	30.84	79.61
SM ² -Uni			
γ	AL	SacreBLEU	COMET
0.3	1.39	24.68	75.58
0.4	2.4	27.88	78.09
0.5	3.56	29.6	79.51
0.55	5.2	30.67	80.28
0.6	6.33	30.86	80.36
0.65	8.06	30.89	80.42
0.7	10.74	31.08	80.53
SM ² -Bi			
γ	AL	SacreBLEU	COMET
0.3	1.52	24.74	75.96
0.4	2.73	28.48	78.85
0.5	3.73	30.17	80.21
0.55	5.49	31.11	80.83
0.6	7.03	31.42	81.00
0.65	9.22	31.65	81.18
0.7	12.33	31.92	81.25

Table 6: Numerical results on WMT15 De→En.

English→Romanian			
wait- k			
k	AL	SacreBLEU	COMET
1	2.70	26.62	74.12
3	5.05	29.74	77.52
5	7.18	31.61	78.54
7	9.10	31.86	79.20
9	10.92	31.89	78.97
m-wait- k			
k	AL	SacreBLEU	COMET
1	2.66	26.65	74.27
3	5.07	30.11	77.44
5	7.18	31.05	78.35
7	9.07	31.44	78.71
9	10.89	31.37	78.62
HMT			
(L, K)	AL	SacreBLEU	COMET
(1,2)	1.98	24.11	71.73
(2,2)	2.77	27.18	74.85
(4,2)	4.47	30.41	77.65
(5,4)	5.48	31.56	78.80
(6,4)	6.45	31.88	78.94
(7,6)	7.41	31.85	79.17
(9,6)	9.24	31.98	79.05
ITST			
δ	AL	SacreBLEU	COMET
0.1	2.75	22.76	71.19
0.2	3.25	28.40	75.58
0.3	5.09	30.52	77.53
0.4	7.47	31.37	78.28
0.45	8.81	31.62	78.49
0.5	10.30	31.63	78.51
0.55	11.69	31.74	78.73
SM ² -Uni			
γ	AL	SacreBLEU	COMET
0.3	2.52	27.85	75.45
0.4	2.72	29.21	76.62
0.5	3.16	30.21	77.59
0.6	4.17	31.20	78.26
0.65	5.13	31.56	78.58
0.7	6.56	31.72	78.77
0.75	8.67	31.67	78.98
SM ² -Bi			
γ	AL	SacreBLEU	COMET
0.3	2.60	28.74	76.81
0.4	2.91	30.27	78.20
0.5	3.57	31.33	79.04
0.6	5.11	32.03	79.56
0.65	6.51	32.40	79.90
0.7	8.15	32.59	79.85
0.75	10.10	32.74	79.95

Table 7: Numerical results on WMT16 En→Ro.

	LONG	SHORT
Average Sentence Length	36.95	14.07
Number of Sentences	1085	1084

Table 8: Statistics on the average sentence length and number of sentences for LONG and SHORT groups.

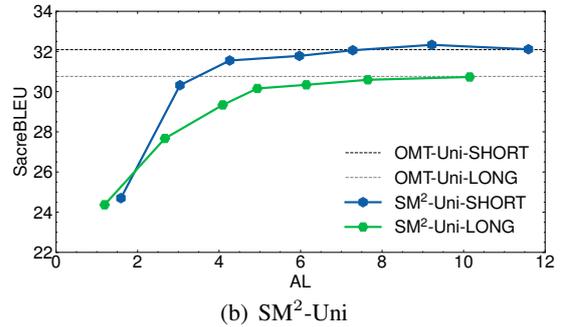
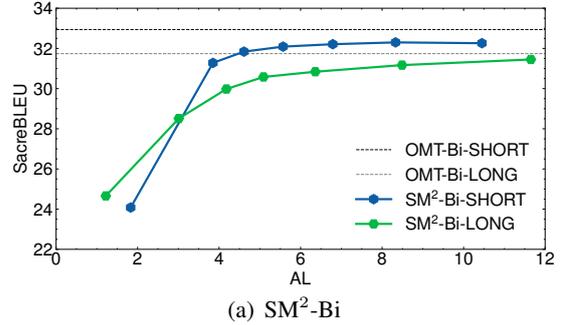


Figure 10: Translation quality against latency of SM² on LONG and SHORT groups. We provide the performance of OMT models for a clearer comparison.

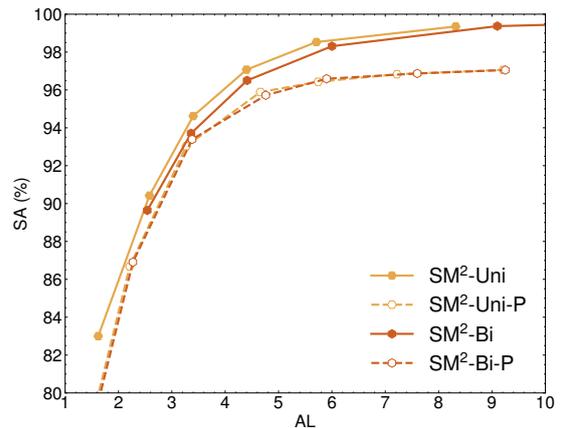


Figure 11: Evaluation of policies in SM² with and without prohibition. We calculate SA (\uparrow) under different latency levels.