

An Empirical Study of Iterative Refinements for Non-autoregressive Translation

Anonymous ACL submission

Abstract

Iterative non-autoregressive (NAR) models share a spirit of mixed autoregressive (AR) and fully NAR models, seeking a balance between generation quality and inference efficiency. These models have recently demonstrated impressive performance in varied generation tasks, surpassing the autoregressive Transformer. However, they also face several challenges that impede further development. In this work, we target building more efficient and competitive iterative NAR models. Firstly, we produce two simple metrics to identify the potential problems existing in current refinement processes, and look back on the various iterative NAR models to find the key factors for realizing our purpose. Subsequently, based on the analyses of the limitations of previous inference algorithms, we propose a simple yet effective strategy to conduct efficient refinements without performance declines. Experiments on five widely used datasets show that our final models set the new state-of-the-art performance compared to all previous NAR models, even with fewer decoding steps, and outperform AR Transformer by around one BLEU on average.

1 Introduction

Transformer-based models (Vaswani et al., 2017) have achieved promising performance in various tasks, particularly after the emergence and progress of large language models recently (Touvron et al., 2023a; OpenAI, 2023; Touvron et al., 2023b). However, these models adopt an autoregressive (AR) decoding paradigm where tokens are generated one by one in a strict left-to-right order. Consequently, they suffer from low inference efficiency, which even worsens as model parameters increase (Zhao et al., 2023). Non-autoregressive (NAR) models provide an alternative text generation paradigm (Gu et al., 2018). Unlike AR models, NAR models can predict all the target tokens in parallel, significantly accelerating the inference process. However, this

parallel decoding paradigm also leads to performance degradation due to independent predictions lacking target side dependency (Qian et al., 2021; Xiao et al., 2022; Huang et al., 2023).

To balance the generation quality and inference efficiency, researchers have proposed iterative NAR models that utilize multiple decoding steps to generate the final results and retain the non-autoregressive decoding paradigm in each step (Lee et al., 2018; Ghazvininejad et al., 2019; Chan et al., 2020). Results generated from the previous decoding steps can provide partial target side information, and then be refined better in the subsequent steps. Through iterative refinements, the performance of these models can achieve significant improvements, even surpassing their AR counterparts (Huang et al., 2022c; Xiao et al., 2023). However, utilizing multiple decoding steps also increases the decoding time overhead, leading to less efficient inference processes (Kasai et al., 2020b; Helcl et al., 2022), and directly reducing the number of decoding steps for relatively faster decoding always brings inferior performance. Besides, these models have also revealed some flaws in the corresponding research, including the gaps between training and inference (Ghazvininejad et al., 2020; Huang et al., 2022c) and the anisotropic problem (Guo et al., 2023a), which hinders the further developments of iterative NAR models and results in less competitive performance. Therefore, *how to build more efficient and competitive iterative NAR models* deserves further exploration.

In this paper, we closely examine various aspects related to the abilities of iterative NAR models, and conduct systematic studies and analytical experiments to address the above-mentioned question:

- We conduct in-depth explorations of previous iterative NAR models (§3). Specifically, we verify and quantitatively analyze the potential problems existing in current refinement pro-

cesses through two metrics (§3.1). Besides, we conduct analytical experiments based on various iterative NAR models and discover that different enhanced methods play different roles in building efficient and competitive models (§3.2). Then, we attempt to realize our purpose by combining previous superior methods, but notice performance declines with previous efficient strategies (§3.3).

- We trial better strategies for iterative NAR models to become efficient while maintaining competitive performances (§4). We first analyze the limitations of current refinement strategies (§4.1) and then propose a simple yet effective inference algorithm for iterative NAR models (§4.2). Combining it with previous competitive strategies can achieve superior performance with fewer decoding steps.

Experiments on 5 widely used datasets demonstrate the effectiveness of our final models. We yield significant performance improvements (around 0.8 BLEU score on average) over the previous best iterative NAR models and realize completely surpassing AR Transformer (over 1 BLEU score on average). Besides, our models only need 4 decoding steps to set new SOTA performance on all WMT datasets compared with those achieved by previous models with 10 decoding steps, which leads to more efficient inference processes.

2 Preliminaries & Motivation

Non-autoregressive Language Model Up to now, most generative models are autoregressive (AR) models which generate the target sequence one by one from a left-to-right order during inference. They adopt AR factorization during training to maximize the following likelihood: $\mathcal{L}_{AR} = \sum_{t=1}^T \log P(y_t | y_{<t}, X; \theta)$, where $y_{<t}$ denotes the previous generated target tokens, T denotes the target length, X is the source sentence, and θ denotes the model parameters. Unlike these AR models, non-autoregressive (NAR) language models generate the target sequence in parallel during inference, they adopt conditional independent factorization during training to maximize the following likelihood: $\mathcal{L}_{NAR} = \sum_{t=1}^T \log P(y_t | X; \theta)$.

CMLM Conditional Masked Language Model (CMLM) is a typical and widely-used iterative NAR model (Ghazvininejad et al., 2019), which

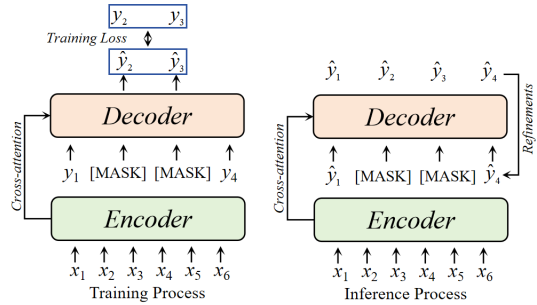


Figure 1: Presentation of the training and inference process of the conditional masked language model.

adopts a Transformer-based encoder-decoder architecture with some specific modifications in the decoder blocks to support NAR generation manner. During each decoding step, CMLM predicts the masked tokens in parallel conditioned on the source sequence and the unmasked part in the target sequence as shown in Figure 1. Then, CMLM will refine the generated output by re-masking and re-predicting several specific tokens according to the Mask-Predict algorithm (Ghazvininejad et al., 2019). During training, CMLM adopts masked language modeling tasks to learn this paradigm. As shown in Figure 1, given a training pair (X, Y) , CMLM selects partial tokens in Y to be masked, denoted as Y_{mask} , while the unmasked tokens as Y_{obs} , then CMLM aims to maximize: $\mathcal{L}_{CMLM} = \sum_{y_t \in Y_{mask}} \log P(y_t | Y_{obs}, X; \theta)$. Besides, CMLM adopts a special token (e.g., [LENGTH]) in its encoder (we omit this in Figure 1) to predict the target length conditional on the source representation following the previous works (Ghazvininejad et al., 2019; Huang et al., 2022c; Xiao et al., 2023).

Follow-up Methods of CMLM Based on CMLM, researchers have proposed many follow-up enhanced methods from different perspectives to improve the training and inference process, e.g., adopting better masking methods (Guo et al., 2020; Xiao et al., 2023) or enhanced modeling mechanism (Kasai et al., 2020a; Cheng and Zhang, 2022; Chen et al., 2024) during training, utilizing extra modules to help training or inference (Hao et al., 2021; Liang et al., 2022; Geng et al., 2021), introducing better inference mechanisms (Ghazvininejad et al., 2020; Huang et al., 2022c), and etc. More details about these variants are included in the Appendix A due to the length limitation. Although we have learned from their own papers that there has been a corresponding improvement in the perfor-

mance of these methods compared to CMLM, we should compare them under more consistent hardware and settings to analyze the effects of different enhanced methods from different perspectives. In this paper, we conduct comprehensive analysis and analytical experiments of CMLM and these follow-up methods, targeting building more efficient and competitive iterative NAR models.

3 In-depth Explorations of Previous Iterative NAR Models

In this section, we conduct in-depth explorations of current iterative NAR models. Specifically, we introduce several sub-problems and make detailed analyses. We aim to find the key factors for building efficient and competitive iterative NAR models.

Problems and Explorations. Firstly, previous works always adopt the final generated output after pre-defined decoding steps to evaluate iterative NAR models, but overlook the fine-grained analysis of intermediate states throughout the refinement process. Consequently, some potential problems (e.g., useless and negative decoding steps) during the refinement process can not be reflected based on the current evaluation process. Therefore, we introduce two metrics (DRR and ROR) to quantitatively analyze the potential problems mentioned above, we aim to answer: *how to evaluate the intermediate states of the whole refinement process of different iterative NAR models* (§3.1). Based on our proposed two metrics, we compare different iterative NAR models under a consistent re-implementation. We aim to find *what are the key components for iterative NAR models to perform better* (§3.2). Finally, we conduct extended experiments to answer *can combining superior methods brings benefits* (§3.3), and make a summary (§3.4).

Experimental Settings. We adopt the vanilla CMLM and several typical variants which contain different improving strategies from different respects as mentioned in Section 2 for exploration. We summarize them as different categories: adopting enhanced training skills (JM-NAT, AMOM, Multitask-NAT), using adaptive inference algorithms (Disco, Rewrite-NAT), and introducing self-correction mechanisms (SMART, CORR, CMLMC). To make more consistent comparisons, we re-implement all these models based on the same hardware and training hyper-parameters. For the evaluation dataset, we select the IWSLT’14

DE→EN dataset containing about 170k training sentence pairs, 7k valid pairs, and 7k test pairs. We train each model on the training set and then evaluate them on the test set. Following the previous work (Kasai et al., 2020a), we apply sequence-level knowledge distillation (Kim and Rush, 2016) for all backbone models. All experiments use the Fairseq library (Ott et al., 2019) with GTX 3090 GPU cards. We adopt the same training hyper-parameters following CMLM realization in Fairseq. During inference, we average the 5 best checkpoints chosen by validation BLEU as our final model. Finally, we evaluate the generation quality with BLEU score (Papineni et al., 2002). Besides, to eliminate the effects of randomness, we follow the previous works to use statistical significance tests (Koehn, 2004) to detect if the difference in BLEU score between the traditional CMLM and other enhanced iterative NAR models is significant.

3.1 How to Evaluate the Intermediate States of the whole Refinement Process of Different Iterative NAR Models?

The current common evaluation process of iterative NAR models, where we compare the generated output of the last decoding step with ground truth to compute a score (e.g., BLEU), can not directly reflect the potential problems as mentioned above. Therefore, we introduce Decline Risks of Refinements (DRR) and Ratio of Over-Refinements (ROR) to respectively measure the ratio of these potential problems, and evaluate the stability and reliability of the refinement process.

Decline Risks of Refinements. Decline Risks of Refinements (DRR) evaluates the stability of the refinement process of iterative NAR models. It measures the performance decline rate after one specific decoding step, i.e., the extent of the negative decoding step. Specifically, given a test set with N examples, a fixed decoding step T , we compute the ratio of each example during the whole refinement process where the performance declines compared with the previous iteration, formatted as:

$$\text{DRR} = \frac{1}{T-1} \sum_{t=1}^{T-1} \frac{|\text{Score}_i^t > \text{Score}_i^{t+1}|}{N}, \quad (1)$$

where Score_i^t denotes the performance of sample i in the t th step.

Ratio of Over-Refinements. Ratio of Over-Refinements (ROR) evaluates the reliability of the

final generated output in iteration T . It measures the failure rate of the output from the last decoding step to be the best, i.e., the extent of the useless decoding steps. Specifically, given a test set with N examples, a fixed decoding step T , we compute the ratio of each example whose best performance is achieved in the intermediate steps of the refinement process, formatted as:

$$\text{ROR} = \frac{1}{T-1} \sum_{t=1}^{T-1} \frac{|\text{Score}_i^t > \text{Score}_i^T|}{N}, \quad (2)$$

where Score_i^t denotes the performance of sample i in the t th step, Score_i^T denotes the performance of sample i in the final iteration T .

3.2 What are the Key Components for Iterative NAR Models to Perform Better?

Exploration Process. We look for the key components for two aspects, i.e., *competitive* and *efficient*. The former can be directly reflected in the final performances, and the latter is reflected in the efficiency of achieving relatively promising performances. We re-implement and evaluate the related enhanced CMLM methods. For these with adaptive inference algorithms (i.e., Disco and RewriteNAT), we set T as the number of adaptive decoding steps of each sentence pair in Equation 1 and Equation 2 during inference, and 10 for other methods following the previous works. Besides, for the models that support two inference algorithms (e.g., CMLMC can omit the self-correction process and transform to the Mask-Predict algorithm), we both report the results with the Mask-Predict algorithm and the corresponding enhanced inference strategy.

Main Findings. The results are presented in Table 1, we find that: (1) *DRR and ROR are relatively lower while decoding with adaptive inference algorithms*. These models aim to find more suitable methods to decide how many and which tokens to mask, and when to stop refinements during inference. They can achieve comparable performance with fewer decoding steps, indicating that adaptive inference algorithms bring benefits to building more efficient iterative NAR models. (2) *Enhanced training skills bring benefits on performance, but there is no evident improvement on DRR and ROR*. These models trained with enhanced training skills can improve performance compared with the vanilla CMLM, but DRR and ROR are still high, indicating that enhanced training skills are

Methods	Iteration	BLEU	DRR (%)	ROR (%)
<i>Enhanced Training Skills</i>				
CMLM	10	33.55	13.4	19.1
JM-NAT	10	32.60	14.4	17.5
Multitask-NAT	10	33.60	16.5	18.4
Disco	10	33.22	14.6	13.1
RewriteNAT †	10	33.88	12.1	14.4
CORR †	10	33.65	13.3	14.1
CMLMC †	10	34.02	13.1	13.8
AMOM †	10	34.68	16.3	17.9
<i>Adaptive Inference Algorithms</i>				
Disco	Adv.	33.32	11.8	6.9
RewriteNAT †	Adv.	33.91	7.9	1.1
<i>Self-correction Mechanism</i>				
SMART	10	33.17	14.5	16.6
CORR †	10	33.76	15.0	15.3
CMLMC †	10	34.40	15.2	14.9
<i>Combining Superior Methods</i>				
AMOMC †	10	35.08	16.8	16.7
w/ Locator †	Adv.	34.68	5.9	6.0

Table 1: DRR and ROR of different models. Adv. denotes adaptive decoding steps, which are always less than 10. † denotes that the BLEU improvements over CMLM are statistically significant with $p < 0.05$.

useful for building more competitive iterative NAR models, the performance improvements of these models come from the better ability to model token dependency during training rather than benefiting the refinement process. (3) *Introducing the self-correction mechanism can improve performance, but DRR gets higher*. These models with the self-correction mechanism achieve performance improvement. However, DRR increases, indicating that this mechanism may bring more unstable factors during the refinement process.

3.3 Can Combining Superior Methods Brings Benefits?

Exploration Process. We can learn from the explorations in Section 3.2 that different enhanced methods are independently beneficial to making the models more efficient and competitive. Naturally, we wonder: can combining superior methods bring benefits? We further explore the following questions: (1) Since the adaptive inference algorithms can bring promising performance with fewer decoding steps, can they further improve the performance with more steps? (2) Since adopting enhanced training skills and the self-correction mechanism can boost performance but not stabilize the refinement process, can we incorporate the adaptive inference

337 algorithms into these models to make them more
338 efficient? Specifically, for question 1, we force
339 these models (Disco and RewriteNAT) to continue
340 the refinement process until reaching the maximal
341 T decoding step. For question 2, we first com-
342 bine the previous superior methods of enhanced
343 training skills and adaptive inference algorithms
344 (AMOM and CMLMC, denoted as AMOMC), and
345 then we further apply the Locator module proposed
346 in RewriteNAT into AMOMC.

347 **Main Findings.** The results are shown in Table 1,
348 we can find that: (1) Concerning the models with
349 adaptive inference algorithms, the performance
350 even declines once we adopt more decoding steps
351 for them, e.g., the performance declines from 33.32
352 to 33.22 for Disco, from 33.91 to 33.88 for Rewrite-
353 NAT. Besides, DRR and ROR get much higher with
354 more decoding steps, indicating that models with
355 adaptive inference algorithms do not need many de-
356 coding steps to achieve the best performance during
357 inference. (2) Further utilizing the Locator mod-
358 ule for AMOMC can make the refinement process
359 more efficient since it can achieve comparable per-
360 formance with fewer decoding steps and get lower
361 DRR and ROR, but it also leads to performance
362 declines compared with the original AMOMC.

363 3.4 Summary

364 Now, we summarize our above explorations. We
365 first propose two simple metrics to analyze the
366 potential problems existing in current refinement
367 methods. We encourage the researchers to pay
368 more attention to the intermediate stages of the re-
369 finement process. Next, we conduct comparative
370 experiments to look for the key components for
371 building more efficient and competitive iterative
372 NAR models, and then further combine superior
373 methods to realize our purpose. However, we find
374 that the current efficient strategy leads to perfor-
375 mance declines. This motivates us to explore better
376 strategies for building efficient iterative NAR mod-
377 els while maintaining competitive performance.

378 4 Trials for Better Efficient Strategies

379 In this section, we explore better strategies for
380 building more efficient iterative NAR models while
381 keeping them maintain competitive performance.
382 We conduct a detailed analysis of original refine-
383 ment methods and then propose a simple yet effec-
384 tive strategy to realize our purpose.

385 **Problems and Explorations.** Firstly, the Mask-
386 Predict algorithm exhibits higher DRR and ROR
387 than adaptive inference algorithms as shown in Ta-
388 ble 1. Therefore, we aim to analyze: *what makes*
389 *the Mask-Predict algorithm fail to do efficient re-*
390 *finements* (§4.1). Besides, noticing that although
391 current adaptive inference algorithms are advanta-
392 geous for reducing the decoding steps, they also
393 lead to performance declines. Therefore, we ana-
394 lyze the corresponding reasons and further investi-
395 gate: *are there better efficient strategies for iter-*
396 *ative NAR models* (§4.2). Finally, we summarize
397 the aforementioned questions and point out future
398 directions for iterative NAR models (§4.3).

399 **Experimental Settings.** During the analysis of
400 the failure of the Mask-Predict algorithm, we
401 adopt the CMLM checkpoint achieved from the
402 above exploration process. To explore more effec-
403 tive inference algorithms, we adopt more datasets,
404 i.e., WMT’16 English↔Roman (En↔Ro) and
405 WMT’14 English↔German (En↔De) language
406 pairs which are widely used in previous NAR
407 works, to evaluate our proposed methods. The
408 training data sizes are respectively about 0.6M and
409 4.5M, and the test data are from the correspond-
410 ing newest data containing around 2,000 and 3,000
411 samples. Besides, the training and evaluation set-
412 tings are the same as those mentioned in Section 3.

413 4.1 What Makes the Mask-Predict Algorithm 414 Fail to Do Efficient Refinements?

415 We attribute the success of the adaptive inference
416 algorithm to the reasonable strategy to determine
417 "which token should be masked in the next decoding
418 step?" Comparatively, the Mask-Predict algorithm
419 relies on predicted confidence to select masked
420 tokens in the subsequent decoding step. How-
421 ever, we have identified two shortcomings with
422 this confidence-based refinement process:

423 1) *The independent confidence updating strat-*
424 *egy for each token is sub-optimal.* In the Mask-
425 Predict algorithm, the prediction confidence is up-
426 dated only for masked tokens during each decod-
427 ing step, i.e., those for unmasked tokens remain
428 unchanged after the last decoding step when they
429 were predicted. This denotes that the prediction
430 confidences of masked and unmasked tokens are
431 derived from different decoding steps and under
432 different masking conditions. Consequently, this in-
433 consistency poses challenges in determining which
434 tokens to be masked in the subsequent decoding

step. This shortcoming is also supported by the comparison presented in Table 1. The models that update the confidence scores of all tokens in the same decoding step can alleviate this problem, e.g., Disco, RewriteNAT, and CMLM all achieve lower DRR and ROR even without adopting adaptive inference algorithms during inference.

2) *The prediction confidence of CMLM is not strongly related to the generation quality.* As discussed in Section 2, CMLM selects the prediction probability as the confidence to choose newly masked tokens. This approach assumes that tokens with higher prediction probabilities are more reliable. However, previous works have highlighted several issues. Ding et al. observe that some specific tokens, such as high-frequency words and conjunctions, consistently exhibit high confidence, leading to repetitive output and neglect of low-frequency but important words. Additionally, Liang et al. notes that the function words dominate the high probability region of the output distribution, making it challenging to generate informative tokens using the Mask-Predict algorithm with CMLM. We also perform a simple experiment to exhibit the irrelevance between the prediction confidence and final generation output.

Set	Win (%)	Lose (%)
Valid	54.61	45.39
Test	54.10	45.90

Table 2: **Win** denotes the model predicts the ground truth token as the final results, **Lose** denotes the vice.

Exploration Process. Since the traditional Mask-Predict algorithm always selects tokens with the highest prediction probability as output during each decoding step, we verify whether the probability of ground truth tokens ranks first. Specifically, we first randomly mask several tokens in the target sequence and send them into CMLM to obtain the prediction probability, then we find the probability place of ground truth tokens, e.g., given the test sentence "Thank you." We first replace the token "you" with the [MASK] token, then we send the sequence "Thank [MASK]. " into CMLM, and verify whether the prediction probability of token "you" ranks first. If not, the highest prediction confidence does not equal the correct token.

Main Findings. We conduct experiments on the validation and test set of IWSLT'14 DE→ dataset

and show the results in Table 2. We find that only around 54 percent of tokens meet our expectations, i.e., these ground truth tokens have the highest prediction probabilities. This shows that the model's own prediction probabilities are not strongly related to the correct tokens. We attribute this failure to the conditional independent factorization during training, which causes CMLM to fail to capture the target-side dependency well (Gu and Kong, 2021).

4.2 Are There Better Efficient Strategies for Iterative NAR Models?

The explorations in Section 4.1 explain why adaptive inference algorithms are more effective than the traditional Mask-Predict algorithm. However, noticing that adopting the Locator leads to performance declines as shown in Table 1 (34.68 v.s. 35.08). We further analyze the corresponding reason. Since the Locator module assigns zero-one discrete scores for predicted tokens, i.e., these tokens scored as zero will be masked again, and one will not be masked in the next decoding step. We point out that this scoring mechanism is too absolute, i.e., there is no difference for unreliable tokens that are scored as zero, and once all the tokens are scored as one, there are no subsequent actions to further improve the generation quality. To explore the potential of a more effective scoring module for iterative NAR models, we intended to replace the zero-one discrete score with a zero-one continuous distribution, in which we can design the refinement process more flexibly and constantly.

Exploration Process. We aim to find a simple yet effective mechanism to score each token within a sentence, and then we can depend on these scores to determine which tokens should be masked in the subsequent decoding step. Motivated by the previous practice that a pre-trained AR model can successfully serve as an effective scorer on the sentence-level to evaluate the fluency of sentences, we can extend it as a token-level scorer, named ARSCORER in the remaining space of this paper. Specifically, we utilize the generated tokens from each decoding step as inputs for a pre-trained AR model. The AR model conducts its prediction on this input sequence in an autoregressive manner. Subsequently, we obtain the corresponding prediction distribution and use the probability associated with the input token index as the final score. The scores range from zero to one after undergoing the normalized softmax operation. Comparatively,

Model		Iter.	WMT'14		WMT'16		Speedup
			EN→DE	DE→EN	EN→RO	RO→EN	
AR	Transformer (Vaswani et al., 2017)	<i>N</i>	27.30	31.29	-	-	-
	Transformer*	<i>N</i>	28.41	32.28	<u>34.23</u>	34.28	1.0x
	Transformer (SCORER 12-1)*	<i>N</i>	<u>28.91</u>	<u>32.65</u>	33.75	<u>34.34</u>	2.5x
Fully NAR	GLAT (Qian et al., 2021)	1	25.21	29.84	31.19	32.04	15.3x
	Fully NAT (Gu and Kong, 2021)	1	27.49	31.39	33.79	34.16	16.5x
	latent-GLAT (Bao et al., 2022)	1	26.64	29.93	-	-	11.3x
	DA-Transformer (Huang et al., 2022b)	1	27.49	31.37	-	-	13.9x
	RenewNAT (Guo et al., 2023b)	1	26.65	30.65	33.02	33.74	11.2x
	FA-DAT (Ma et al., 2023)	1	27.49	31.37	-	-	14.0x
	PCFG-NAT (Gui et al., 2024)	1	27.02	31.29	32.72	33.07	12.6x
Iterative NAR	Levenshtein (Gu et al., 2019)	Adv.	27.73	-	33.02	-	4.0x
	CMLM (Ghazvininejad et al., 2019)	10	27.03	30.53	33.08	33.31	1.7x
	DisCo (Kasai et al., 2020a)	Adv.	27.34	-	33.25	33.22	3.5x
	SMART (Ghazvininejad et al., 2020)	10	27.65	31.27	33.85	33.53	1.7x
	JM-NAT (Guo et al., 2020)	10	27.69	32.24	33.52	33.72	5.7x
	RewriteNAR (Geng et al., 2021)	Adv.	27.83	31.52	33.63	34.09	3.9x
	MvCR-NAT (Xie et al., 2021)	10	27.39	31.18	33.38	33.56	3.8x
	CORR (Huang et al., 2022c)	10	28.19	31.31	34.31	34.08	-
	CMLMC (Huang et al., 2022c)	10	<u>28.37</u>	31.41	34.57	34.13	-
	AMOM (Xiao et al., 2023)	10	27.57	<u>31.67</u>	<u>34.62</u>	<u>34.82</u>	1.7x
	EECR (Chen et al., 2024)	10	28.04	31.65	34.33	34.32	3.8x
	DCMCL (Liao et al., 2024)	10	28.14	31.47	33.76	33.64	-
Ours*	AMOMC	4	28.39	32.83	34.73	35.11	4.1x
		10	28.81	33.22	34.99	35.20	1.9x
	AMOMC w/ ARSCORER	4	29.08	33.22	35.06	35.78	3.1x
		10	29.18	33.41	35.30	36.02	1.4x

Table 3: Results on 4 WMT machine translation tasks. * denotes the results of our implementations. † denotes that the BLEU improvements over AMOMC are statistically significant with $p < 0.05$. **Bold** denotes the best performance, underline denotes the previous best AR or NAR performance.

adopting ARSCORER offers several advantages over the Mask-Predict algorithm, which have also been mentioned in the previous section: (1) The AR model can assess the validity of each token in the whole sentence and update the corresponding prediction probability of each token after each decoding step of NAR model. (2) Previous studies have shown that models trained with autoregressive factorization excel in capturing target side dependencies compared to NAR models (Huang et al., 2022a). Besides, these AR models do not suffer from the multi-modality problem. Thus, adopting extra ARSCORER to provide the prediction score is more robust and effective. We should recognize that adopting the AR model to achieve the prediction score leads to extra inference time. Therefore, we adopt the structure of deep encoder and shallow decoder as mentioned in Kasai et al. (2020b) for ARSCORER to reduce inference efficiency.

Main Findings. The results on various WMT datasets are shown in Table 3, we can find that: (1) Combining superior methods (AMOMC) achieves significant performance improvements, outper-

forming all baseline models around 0.8 BLEU score. (2) Further adopting ARSCORER can quickly achieve competitive performance, i.e., it can achieve new state-of-the-art performance with only 4 decoding steps while getting 3.1 times speedup compared to the vanilla AR Transformer. (3) Adopting ARSCORER increases the inference time compared to AMOMC, but it makes up for this deficiency by achieving superior performance with relatively fewer decoding steps, which still indicates that ARSCORER can bring benefits for building efficient iterative NAR models.

Further Analysis. We further conduct detailed analytical experiments of our proposed methods. Firstly, we compare the backbone models with and without ARSCORER based on our proposed two metrics, DRR and ROR, as mentioned in Section 3.1. Results on IWSLT'14 DE→EN and WMT'16 RO→EN datasets are presented in Table 4. We can find that: (1) The models with ARSCORER can achieve lower DRR and ROR compared with the corresponding baselines. (2) DRR and ROR are higher on the WMT'16 RO→EN

Methods	Iter.	BLEU	DRR (%)	ROR (%)
IWSLT'14 DE→EN				
CMLM	10	33.55	13.4	19.1
w/ ARSCORER	10	34.05	10.0	13.4
AMOMC	10	35.08	16.8	16.7
w/ ARSCORER	10	35.61	9.8	13.6
WMT'16 RO→EN				
CMLM	10	33.87	21.0	36.1
w/ ARSCORER	10	34.51	13.5	19.4
AMOMC	10	35.20	19.4	28.7
w/ ARSCORER	10	36.02	14.0	19.1

Table 4: Results of DRR and ROR with ARSCORER.

dataset across all models, indicating that this dataset is relatively difficult to learn. Besides, we compare the results using different AR models to serve as ARSCORER. As we adopt the model with the structure of deep encoder and shallow (denoted as DESD ARSCORER) to reduce inference efficiency in our main result, we also adopt the model with common layers (i.e., 6 encoders and 6 decoders, denoted as vanilla ARSCORER) here. Results are presented in Table 5. We can find that (1) Adopting vanilla ARSCORER leads to lower decoding speed due to the high cost of passing the AR model to achieve the output probabilities. (2) The DESD ARSCORER can reduce the decoding latency due to the shallow decoder structure but still brings extra expense compared to the vanilla NAR model without ARSCORER with the same decoding steps. (3) The DESD ARSCORER can effectively reduce the decoding steps while achieving significant performance improvements, which can make up for the flaw of the extra expense of the AR model. Besides, considering that BLEU is sensitive to the predicted length, we also adopt Comet (Rei et al., 2020) to evaluate our methods, results are also shown in Table 5, we can find a similar phenomenon, i.e., adopting vanilla and DESD ARSCORER can both achieve better Comet score.

4.3 Summary

In this section, we aim to explore the potential for better efficient strategies. We begin by examining the limitations of the Mask-Predict algorithm in facilitating consistent and efficient refinements. Through thorough analysis and corresponding experimentation, we attribute these limitations to the independent confidence updating strategies and the unrelated prediction confidence to generation output. Consequently, we endeavor to identify a superior strategy to address these issues. Fortunately,

Methods	Iter.	BLEU	Comet	Latency
WMT'14 EN→DE				
w/o ARSCORER	4	28.39	0.584	130.8
	10	28.81	0.585	282.1
w/ vanilla ARSCORER	4	28.82	0.587	222.0
	10	29.17	0.587	522.7
w/ DESD ARSCORER	4	29.08	0.587	181.6
	10	29.18	0.588	402.1
WMT'16 RO→EN				
w/o ARSCORER	4	35.11	0.765	87.9
	10	35.20	0.768	211.7
w/ vanilla ARSCORER	4	35.26	0.768	147.6
	10	35.65	0.775	341.7
w/ DESD ARSCORER	4	35.78	0.772	116.0
	10	36.02	0.772	253.9

Table 5: Results of different scorer models. Iter. denotes the decoding step. Latency denotes the total seconds.

by adopting the pre-trained AR models to serve as a scorer, iterative NAR models can conduct steady and effective refinements, thereby achieving superior performance with even fewer decoding steps, and getting closer to the efficient iterative NAR models. It is worth noting that there are other viable options for scoring, such as adopting a pre-trained language model or even current well-known large language models, we leave this as future work.

5 Conclusion and Future Outlook

In this paper, we conduct extensive experiments and detailed analysis based on various iterative NAR models to explore *how to build more efficient and competitive iterative NAR models*. By combining competitive strategies and the newly proposed ARSCORER, our final models set the new state-of-the-art results on several widely used datasets even with fewer decoding steps, which leads to completely outperforming their AR counterparts.

In the future, we will extend our explorations to more scenarios since CMLM-based iterative NAR models have been successfully applied in speech and video-related fields (Higuchi et al., 2021). Besides, there is also a need to explore methods for conducting efficient denoising steps for diffusion models (Sohl-Dickstein et al., 2015) since they suffer greatly from low efficiency with numerous denoising steps (Tang et al., 2023; Gong et al., 2023). Lastly, recent advancements in LLMs (Touvron et al., 2023b) hold promise in serving as better scorers for iterative NAR models.

643 Limitations

644 Firstly, since CMLM-based iterative NAR models
645 have been applied to various language generation
646 tasks, we only conduct our explorations on ma-
647 chine translation task. Besides, although CMLM-
648 based methods are one of the most widely-used and
649 well-known iterative NAR models, there exist other
650 categories of iterative NAR models, such as editing-
651 based models (Stern et al., 2019; Gu et al., 2019),
652 denoising based models (Lee et al., 2018; Savinov
653 et al., 2021), we only consider CMLM-based meth-
654 ods in this paper. Besides, our proposed efficient
655 strategy, ARSCORER, relies on a pre-trained AR
656 model to serve as a scorer for each token, it brings
657 some extra costs to achieve this AR model and the
658 corresponding prediction confidence.

659 References

660 Yu Bao, Hao Zhou, Shujian Huang, Dongqi Wang, Li-
661 hua Qian, Xinyu Dai, Jiajun Chen, and Lei Li. 2022.
662 *latent-glat: Glancing at latent variables for parallel*
663 *text generation*. In *Proceedings of the 60th Annual*
664 *Meeting of the Association for Computational Lin-*
665 *guistics*, volume 1, pages 8398–8409.

666 William Chan, Chitwan Saharia, Geoffrey Hinton, Mo-
667 hammad Norouzi, and Navdeep Jaitly. 2020. Imputer:
668 Sequence modelling via imputation and dynamic pro-
669 gramming. In *ICML*, pages 1403–1413. PMLR.

670 Xinran Chen, Sufeng Duan, and Gongshen Liu. 2024.
671 Improving non-autoregressive machine translation
672 with error exposure and consistency regularization.
673 *arXiv preprint arXiv:2402.09725*.

674 Hao Cheng and Zhihua Zhang. 2022. Con-nat: Con-
675 trastive non-autoregressive neural machine transla-
676 tion. In *Findings of the Association for Computa-*
677 *tional Linguistics: EMNLP 2022*, pages 6219–6231.

678 Liang Ding, Longyue Wang, Xuebo Liu, Derek F Wong,
679 Dacheng Tao, and zhaopeng Tu. 2021. Rejuvenat-
680 ing low-frequency words: Making the most of par-
681 allel data in non-autoregressive translation. In *ACL-*
682 *IJCNLP*, pages 3431–3441.

683 Xinwei Geng, Xiaocheng Feng, and Bing Qin. 2021.
684 Learning to rewrite for non-autoregressive neural ma-
685 chine translation. In *Proceedings of the 2021 Con-*
686 *ference on Empirical Methods in Natural Language*
687 *Processing*, pages 3297–3308.

688 Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and
689 Luke Zettlemoyer. 2019. Mask-predict: Parallel de-
690 coding of conditional masked language models. In
691 *Proceedings of the 2019 Conference on Empirical*
692 *Methods in Natural Language Processing and the 9th*
693 *International Joint Conference on Natural Language*
694 *Processing*, pages 6112–6121.

695 Marjan Ghazvininejad, Omer Levy, and Luke Zettle-
696 moyer. 2020. Semi-autoregressive training im-
697 proves mask-predict decoding. *arXiv preprint*
698 *arXiv:2001.08785*.

699 Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu,
700 and Lingpeng Kong. 2023. Diffuseq-v2: Bridg-
701 ing discrete and continuous text spaces for accel-
702 erated seq2seq diffusion models. *arXiv preprint*
703 *arXiv:2310.05793*.

704 Jiatao Gu, James Bradbury, Caiming Xiong, Victor OK
705 Li, and Richard Socher. 2018. Non-autoregressive
706 neural machine translation. In *International Confer-*
707 *ence on Learning Representations*.

708 Jiatao Gu and Xiang Kong. 2021. Fully non-
709 autoregressive neural machine translation: Tricks of
710 the trade. In *Findings of the Association for Com-*
711 *putational Linguistics: ACL-IJCNLP 2021*, pages
712 120–133.

713 Jiatao Gu, Changan Wang, and Junbo Zhao. 2019. Lev-
714 enshtein transformer. In *Advances in Neural Infor-*
715 *mation Processing Systems*, volume 32, pages 11181–
716 11191.

717 Shangdong Gui, Chenze Shao, Zhengrui Ma, Yunji Chen,
718 Yang Feng, et al. 2024. Non-autoregressive machine
719 translation with probabilistic context-free grammar.
720 *Advances in Neural Information Processing Systems*,
721 36.

722 Junliang Guo, Linli Xu, and Enhong Chen. 2020.
723 Jointly masked sequence-to-sequence model for non-
724 autoregressive neural machine translation. In *Pro-*
725 *ceedings of the 58th Annual Meeting of the Associa-*
726 *tion for Computational Linguistics*, pages 376–385.

727 Pei Guo, Yisheng Xiao, Juntao Li, Yixin Ji, and
728 Min Zhang. 2023a. Isotropy-enhanced conditional
729 masked language models. In *Findings of the Associa-*
730 *tion for Computational Linguistics: EMNLP 2023*,
731 pages 8278–8289.

732 Pei Guo, Yisheng Xiao, Juntao Li, and Min Zhang.
733 2023b. Renewnat: renewing potential translation
734 for non-autoregressive transformer. In *Proceedings*
735 *of the AAAI Conference on Artificial Intelligence*,
736 volume 37, pages 12854–12862.

737 Yongchang Hao, Shilin He, Wenxiang Jiao, Zhaopeng
738 Tu, Michael Lyu, and Xing Wang. 2021. Multi-task
739 learning with shared encoder for non-autoregressive
740 machine translation. In *Proceedings of the 2021*
741 *Conference of the North American Chapter of the*
742 *Association for Computational Linguistics: Human*
743 *Language Technologies*, pages 3989–3996.

744 Jindřich Helcl, Barry Haddow, and Alexandra Birch.
745 2022. Non-autoregressive machine translation:
746 It’s not as fast as it seems. *arXiv preprint*
747 *arXiv:2205.01966*.

748	Yosuke Higuchi, Hirofumi Inaguma, Shinji Watanabe,	Zhengrui Ma, Chenze Shao, Shangtong Gui, Min Zhang,	801
749	Tetsuji Ogawa, and Tetsunori Kobayashi. 2021. Im-	and Yang Feng. 2023. Fuzzy alignments in directed	802
750	proved mask-ctc for non-autoregressive end-to-end	acyclic graph for non-autoregressive machine trans-	803
751	asr. In <i>ICASSP 2021</i> , pages 8363–8367. IEEE.	lation. <i>arXiv preprint arXiv:2303.06662</i> .	804
752	Fei Huang, Pei Ke, and Minlie Huang. 2023. [tacl]	OpenAI. 2023. Gpt-4 technical report.	805
753	directed acyclic transformer pre-training for high-	Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan,	806
754	quality non-autoregressive text generation. In <i>The</i>	Sam Gross, Nathan Ng, David Grangier, and Michael	807
755	<i>61st Annual Meeting Of The Association For Computa-</i>	Auli. 2019. fairseq: A fast, extensible toolkit for	808
756	<i>tional Linguistics</i> .	sequence modeling. In <i>Proceedings of NAACL-HLT</i>	809
757	Fei Huang, Tianhua Tao, Hao Zhou, Lei Li, and Minlie	<i>2019: Demonstrations</i> .	810
758	Huang. 2022a. On the learning of non-autoregressive	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	811
759	transformers. In <i>International Conference on Ma-</i>	Jing Zhu. 2002. Bleu: a method for automatic evalu-	812
760	<i>chine Learning</i> , pages 9356–9376. PMLR.	ation of machine translation. In <i>Proceedings of the</i>	813
761	Fei Huang, Hao Zhou, Yang Liu, Hang Li, and Minlie	<i>40th annual meeting of the Association for Computa-</i>	814
762	Huang. 2022b. Directed acyclic transformer for non-	<i>tional Linguistics</i> , pages 311–318.	815
763	autoregressive machine translation. In <i>International</i>	Lihua Qian, Hao Zhou, Yu Bao, Mingxuan Wang, Lin	816
764	<i>Conference on Machine Learning</i> , pages 9410–9428.	Qiu, Weinan Zhang, Yong Yu, and Lei Li. 2021.	817
765	PMLR.	Glancing transformer for non-autoregressive neural	818
766	Xiao Shi Huang, Felipe Perez, and Maksims Volkovs.	machine translation. In <i>Proceedings of the 59th An-</i>	819
767	2022c. Improving non-autoregressive translation	<i>nual Meeting of the Association for Computational</i>	820
768	models without distillation. In <i>International Con-</i>	<i>Linguistics and the 11th International Joint Confer-</i>	821
769	<i>ference on Learning Representations</i> .	<i>ence on Natural Language Processing</i> , pages 1993–	822
770	Jungo Kasai, James Cross, Marjan Ghazvininejad, and	2003.	823
771	Jiatao Gu. 2020a. Parallel machine translation with	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	824
772	disentangled context transformer. <i>arXiv preprint</i>	Lavie. 2020. Comet: A neural framework for mt	825
773	<i>arXiv:2001.05136</i> .	evaluation. <i>arXiv preprint arXiv:2009.09025</i> .	826
774	Jungo Kasai, Nikolaos Pappas, Hao Peng, James Cross,	Nikolay Savinov, Junyoung Chung, Mikolaj Binkowski,	827
775	and Noah Smith. 2020b. Deep encoder, shallow	Erich Elsen, and Aaron van den Oord. 2021. Step-	828
776	decoder: Reevaluating non-autoregressive machine	unrolled denoising autoencoders for text generation.	829
777	translation. In <i>ICLR</i> .	<i>arXiv preprint arXiv:2112.06749</i> .	830
778	Yoon Kim and Alexander M Rush. 2016. Sequence-	Jascha Sohl-Dickstein, Eric Weiss, Niru Mah-	831
779	level knowledge distillation. In <i>EMNLP</i> , pages 1317–	eswaranathan, and Surya Ganguli. 2015. Deep un-	832
780	1327.	supervised learning using nonequilibrium thermody-	833
781	Philipp Koehn. 2004. Statistical significance tests for	namics. In <i>ICML</i> , pages 2256–2265. PMLR.	834
782	machine translation evaluation. In <i>Proceedings of</i>	Mitchell Stern, William Chan, Jamie Kiros, and Jakob	835
783	<i>the 2004 conference on empirical methods in natural</i>	Uszkoreit. 2019. Insertion transformer: Flexible se-	836
784	<i>language processing</i> , pages 388–395.	quence generation via insertion operations. In <i>ICML</i> ,	837
785	Jason Lee, Elman Mansimov, and Kyunghyun Cho.	pages 5976–5985. PMLR.	838
786	2018. Deterministic non-autoregressive neural se-	Zecheng Tang, Pinzheng Wang, Keyan Zhou, Juntao	839
787	quence modeling by iterative refinement. In <i>Proceed-</i>	Li, Ziqiang Cao, and Min Zhang. 2023. Can diffu-	840
788	<i>ings of the 2018 Conference on Empirical Methods</i>	sion model achieve better performance in text gener-	841
789	<i>in Natural Language Processing</i> , pages 1173–1182.	ation? bridging the gap between training and infer-	842
790	Xiaobo Liang, Zecheng Tang, Juntao Li, and Min Zhang.	ence! <i>arXiv preprint arXiv:2305.04465</i> .	843
791	2023. Open-ended long text generation via masked	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	844
792	language modeling. In <i>ACL</i> .	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	845
793	Xiaobo Liang, Lijun Wu, Juntao Li, and Min Zhang.	Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal	846
794	2022. Janus: Joint autoregressive and non-	Azhar, Aurelien Rodriguez, Armand Joulin, Edouard	847
795	autoregressive training with auxiliary loss for se-	Grave, and Guillaume Lample. 2023a. Llama: Open	848
796	quence generation. In <i>EMNLP</i> , pages 1067–1073.	and efficient foundation language models. <i>arXiv</i>	849
797	Yusheng Liao, Yanfeng Wang, and Yu Wang. 2024.	<i>preprint arXiv:2302.13971</i> .	850
798	Leveraging diverse modeling contexts with collabo-	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	851
799	rating learning for neural machine translation. <i>arXiv</i>	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	852
800	<i>preprint arXiv:2402.18428</i> .	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	853
		Bhosale, et al. 2023b. Llama 2: Open founda-	854
		tion and fine-tuned chat models. <i>arXiv preprint</i>	855
		<i>arXiv:2307.09288</i> .	856

857	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob	tokens whose prediction confidence is higher than	909
858	Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz	Y_i in the previous iteration). Disco stops decoding	910
859	Kaiser, and Illia Polosukhin. 2017. Attention is all	when no new tokens are generated in one specific	911
860	you need. In <i>Advances in Neural Information Pro-</i>	decoding step. This easy-first policy can largely	912
861	<i>cessing Systems</i> , pages 5998–6008.	improve the inference latency.	913
862	Yisheng Xiao, Lijun Wu, Junliang Guo, Juntao Li,	Multitask-NAT Hao et al. introduces Multitask-	914
863	Min Zhang, Tao Qin, and Tie-yan Liu. 2022. A	NAT which utilizes a shared encoder and separated	915
864	survey on non-autoregressive generation for neural	decoders for both AR and NAR modeling during	916
865	machine translation and beyond. <i>arXiv preprint</i>	training. They assume that AR training can bring	917
866	<i>arXiv:2204.09269</i> .	benefits for NAR training and aim to adopt multi-	918
867	Yisheng Xiao, Ruiyang Xu, Lijun Wu, Juntao Li, Tao	task learning to transfer the AR knowledge to NAR	919
868	Qin, Tie-Yan Liu, and Min Zhang. 2023. Amom:	models through encoder sharing.	920
869	Adaptive masking over masking for conditional	RewriteNAT Geng et al. propose RewriteNAT, a	921
870	masked language model. <i>Proceedings of the AAAI</i>	new framework that contains a Locator and Revisor	922
871	<i>Conference on Artificial Intelligence</i> , 37(11):13789–	module that locate the incorrect words within pre-	923
872	13797.	viously generated translations and then revise them,	924
873	Pan Xie, Zexian Li, and Xiaohui Hu. 2021. Mvsr-	respectively. Specifically, the Locator module can	925
874	nat: Multi-view subset regularization for non-	transform the problem of determining which tokens	926
875	autoregressive machine translation. <i>arXiv preprint</i>	to be masked in the next decoding step into a binary	927
876	<i>arXiv:2108.08447</i> .	classification problem instead of depending on the	928
877	Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,	self-predicted confidence, i.e., the Locator will pre-	929
878	Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen	dict a special symbol ([MASK] or [KEEP]) for each	930
879	Zhang, Junjie Zhang, Zican Dong, et al. 2023. A	token. Once the token is predicted as [MASK], it	931
880	survey of large language models. <i>arXiv preprint</i>	will be masked again, and vice versa. RewriteNAT	932
881	<i>arXiv:2303.18223</i> .	can finish the generation process once the Locator	933
882	A Details for Follow-up Methods	module predicts all the target tokens as [KEEP].	934
883	We supplement the details for follow-up methods of	SMART Ghazvininejad et al. introduce Semi-	935
884	CMLM we adopted for explorations as mentioned	Autoregressive Training (SMART) to help the train-	936
885	in Section 2.	ing process better match the Mask-Predict algo-	937
886	JM-NAT Guo et al. introduce a jointly masked	rithm with multiple decoding steps. Specifically,	938
887	sequence-to-sequence model. Unlike the tradi-	since the model can not see the ground truth tokens	939
888	tional CMLM which only masks the target se-	during inference, it only takes the model prediction	940
889	quence during training, JM-NAT also masks the	in the previous decoding steps as partially observed	941
890	source sequence to help train the encoder more	tokens to make predictions. This leads to inconsis-	942
891	rigorously. Besides, in order to alleviate the prob-	tency compared with training methods. Thus	943
892	lem of translating duplicate words, they propose to	SMART first constructs a mixed training example	944
893	train the decoder based on the consecutive masking	and then encourages the model to recover from the	945
894	of the decoder input with an ngram loss function	model prediction errors during training,	946
895	rather than the original uniform masking.	CMLMC Huang et al. propose Condi-	947
896	Disco Kasai et al. propose an attention-masking	tional Masked Language Model with Correction	948
897	based model, Disentangled Context (DisCo) trans-	(CMLMC) which incorporates a self-correction	949
898	former. During training, Disco is learned to pre-	mechanism into traditional CMLM and several	950
899	dict each target token given an arbitrary subset of	modifications on the decoder structure such as ex-	951
900	the other reference tokens, which is more efficient	posing the positional encodings and incorporat-	952
901	than just predicting masked tokens in the origi-	ing causal attention layers to differentiate adja-	953
902	nal CMLM. During inference, unlike the previous	cent tokens. CORR is the corresponding variant	954
903	Mask-Predict algorithm which just updates masked	which only adopts the self-correction mechanism	955
904	tokens in each decoding step (i.e., predicting Y_{mask}	without the structure modifications in CMLMC.	956
905	based on Y_{obs}), Disco introduces an easy-first pol-	Specifically, except for adopting masking meth-	957
906	icy where each token will be predicted in each	ods in target sequence during training, CMLMC	958
907	step dependent on relatively easier tokens (i.e., pre-		
908	dicting each Y_i based on $Y_{<i}$, where $Y_{<i}$ denotes		

Models	Parameters	IWSLT'14 DE→EN	WMT'14 EN↔DE	WMT'16 EN↔RO
CMLM	learning rate	5e-4	7e-4	5e-4
	warmup_step	4k	10k	10k
	dropout	0.3	0.2	0.3
	update_step	300k	300k	300k
	GPU	1xGTX 3090	4xGTX 3090	4xGTX 3090
AMOMC	learning rate	5e-4	7e-4	5e-4
	warmup_step	30k	40k	15k
	dropout	0.3	0.2	0.3
	update_step	175k	150k	120k
	GPU	1xGTX 3090	4xGTX 3090	4xGTX 3090

Table 6: Training hyper-parameters for CMLM and AMOMC.

also replaces the partially unmasked tokens with model predictions based on a fully masked target sequence. Then CMLMC learns to predict the masked tokens and correct the replaced tokens simultaneously during training. During inference, this self-correction mechanism helps the model to correct the unreliable tokens in the unmasked subset.

AMOM Xiao et al. propose an Adaptive Masking Over Masking (AMOM) strategy based on CMLM which contains two different adaptive masking mechanisms which work on the inputs of encoder and decoder respectively. Specifically, based on the ratio of the target sequence, AMOM also masks the specific number of tokens in the source sequence to make the encoder optimization easier. Besides, AMOM conducts an extra masking step where the masking ratio of the target sequence in this step is adaptive to the correction ratio of the model prediction. This two-step masking strategy can help the model capture the masking ratio changes in various decoding steps during inference.

B Training Hyper-parameters

During our experiments, we set training hyper-parameters for CMLM in the same way as CMLM realization in the Fariseq library, and for AMOMC, we follow those adopted in CMLMC (Huang et al., 2022c). Now, we present these training hyper-parameters in Table 6.

C More Explorations between Iterative NAR and AR Models

Researchers commonly assess the effectiveness of iterative NAR models by comparing their performance (generation quality and inference efficiency) with their AR counterparts. This section delves into further exploration of experimental settings for a comprehensive comparison.

Problems and Explorations. Firstly, the structure of deep encoder and shallow decoder (DESD) has been proven effective for AR models to increase inference speed and improve generation quality but does not work well for iterative NAR models (Kasai et al., 2020b). We assume that their experimental setting, which only adopts the decoder with one layer, is too strict for iterative NAR models, resulting in their relatively biased conclusion that iterative NAR models fail to benefit from DESD. Therefore, we experiment with more layer allocation schemes based on the structure of DESD to explore: *do iterative NAR models really fail with deep encoder and shallow decoder?* (§C.1) Besides, most previous works on iterative NAR models adopt the same decoding steps (4 and 10) for all the testing instances with different lengths (Ghazvininejad et al., 2019; Huang et al., 2022c). Comparatively, as AR models only predict the next token in each decoding step, they need decoding steps proportional to the length of the target sentence to achieve the final results. Thus, they adopt more decoding steps for longer sentences. Since the general knowledge for iterative NAR models is that longer sentences need more decoding steps, adopting the fixed decoding steps for all the target sentences naturally causes inconsistent comparisons with AR models. Consequently, we explore: *can iterative NAR models achieve better performance beyond fixed decoding steps?* (§C.2) Finally, we summarize the above two questions and conclude briefly (§C.3).

Experimental Settings. We adopt two popular iterative NAR models (CMLM and CORR) as mentioned in Section 2) and their AR counterpart (the vanilla Transformer) to conduct detailed analytical experiments. The training and inference settings are the same as mentioned in Section 3. We report the decoding time during inference, i.e., we select S_1 measured by L_1^{GPU} to compare the inference efficiency following the previous work (Kasai et al., 2020b; Helcl et al., 2022)), S_1 denotes the translation latency by running the model with one sentence at a time on a single GPU.

C.1 Do Iterative NAR Models Really Fail with Deep Encoder and Shallow Decoder?

Exploration Process. We adopt the structure of DESD on two backbone iterative NAR models and the vanilla AR Transformer to make comparisons.

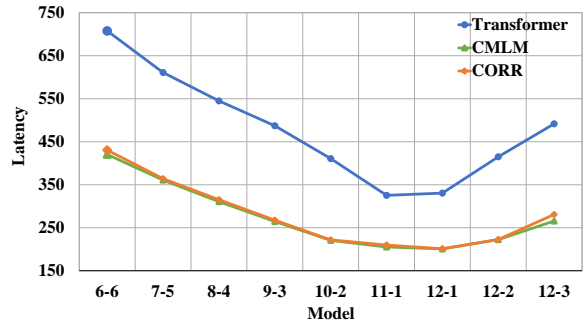
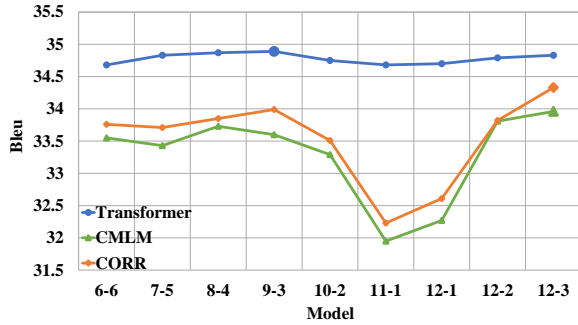


Figure 2: Results with different models and layer allocations.

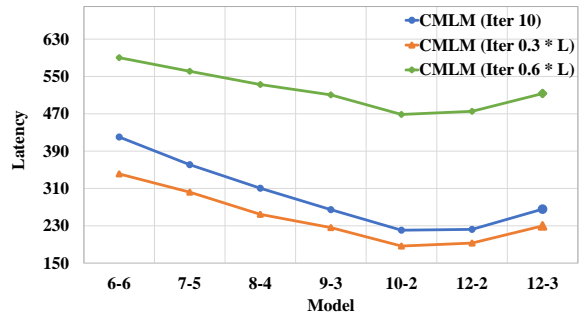
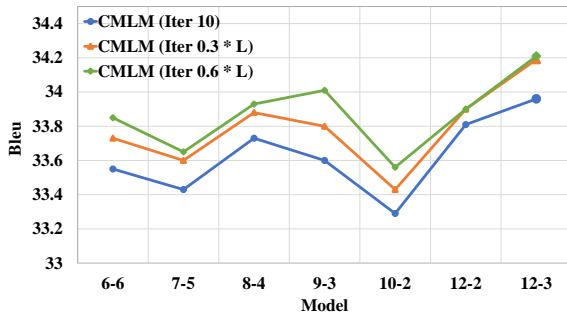


Figure 3: Results with different decoding steps (Iter) during inference. L denotes the length of the source sentence.

1045 We conduct experiments with different layer allocation schemes, each marked as $x-y$. x and y denote the number of encoder and decoder layers, respectively. Specifically, we design the schemes based on two rules: (1) We keep the total layers of encoder and decoder as 12 following the common setting (6-6), then assign different layer allocations for them, e.g., 7-5, 8-4, 9-3, etc. (2) We keep the layers of encoder as 12 (12-1) following the setting in Kasai et al. (2020b), and adopt more decoder layers, e.g., 12-2 and 12-3.

1056 **Main Findings.** Figure 2 presents the corresponding results, we can find that: (1) The number of decoder layers of iterative NAR models significantly affects the inference speed. Comparatively, the number of encoder layers has little effect. (2) Iterative NAR models can also perform well with the structure of DESD (8-4, 9-3, 12-2, 12-3), and accelerate the inference process (more than 1.5x Speedup). (3) Decoder layers are essential for iterative NAR models; at least two decoder layers are needed for iterative NAR models, and more decoder layers achieve better performance. Only one decoder layer harms performance seriously, which is consistent with the findings in Kasai et al. (2020b). Generally, the model with 12-3 layer allocation achieves the best trade-off between genera-

tion quality and inference efficiency.

1072 C.2 Can Iterative NAR Models Achieve Better Performance Beyond Fixed Decoding Steps? 1073 1074 1075

1076 **Exploration Process** We aim to design an adaptive algorithm to decide the decoding steps for each testing instance. To keep consistent with AR models, we can design a mapping function according to the sequence length. We select the relatively simple and understandable linear mapping function, i.e., given the sequence length L , the decoding step is $\alpha * L$, α is 1 in AR models. We can further set it smaller (e.g., 0.3 and 0.6) to realize relatively fast decoding for iterative NAR models. 1077 1078 1079 1080 1081 1082 1083 1084 1085

1086 **Main Findings.** Figure 3 presents the corresponding results, we can find that: (1) While adopting 0.3 * L decoding steps can improve the performance on BLEU score slightly (about 0.2), the inference latency of different models consistently reduces compared with adopting fixed 10 decoding steps. (2) Adopting 0.6 * L decoding steps can further improve the performance of all models but hurt the speed, indicating that blindly increasing the decoding steps is not advisable. (3) Adopting feasible adaptive decoding steps for different test instances can benefit both generation 1087 1088 1089 1090 1091 1092 1093 1094 1095 1096 1097

1098 quality and inference efficiency. We recognize that
 1099 the simple linear mapping function with $\alpha = 0.3$
 1100 may not be optimal, but it has been verified that
 1101 improvements can be achieved beyond fixed de-
 1102 coding steps. Furthermore, we also compare the
 1103 improvements in generation quality based on dif-
 1104 ferent source lengths. Specifically, we divide the
 1105 source sentences into five intervals by the corre-
 1106 sponding length and then compare the performance
 1107 improvements through increasing iterations. We
 1108 plot the performance improvements from fixed 10
 1109 to $0.3 * L$ decoding steps in Figure 4, demonstrating
 1110 that long sequences need more decoding steps.

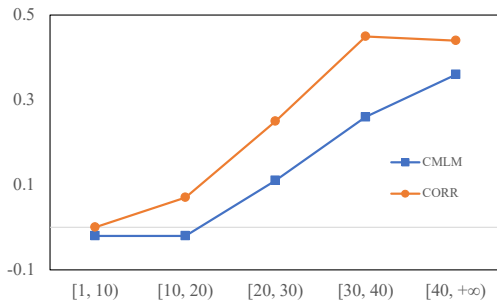


Figure 4: The performance improvements with more iterations based on different source lengths.

1111 C.3 Summary

1112 In this section, we explore more settings during
 1113 comparison between iterative NAR and AR models.
 1114 Different from the findings in Kasai et al. (2020b),
 1115 our experiments demonstrate that the structure of
 1116 DESD can also work well for iterative NAR mod-
 1117 els by simply adopting one more decoder layer and
 1118 more performance improvement can be achieved
 1119 with more layer allocation schemes. Besides, adopt-
 1120 ing the fixed decoding steps leads to an inconsis-
 1121 tent setting since AR models adopt decoding steps
 1122 adaptive to the sequence length. As a result, we de-
 1123 sign a simple adaptive mapping function to decide
 1124 the decoding steps for different testing instances,
 1125 and achieve both performance and inference speed
 1126 improvements. In general, we point out that com-
 1127 monly used and recognized settings still need ex-
 1128 ploration for iterative NAR models, and we should
 1129 consider more consistent settings when comparing
 1130 iterative NAR and AR models.