

Can Multilinguality Benefit Non-autoregressive Machine Translation?

Anonymous ACL submission

Abstract

Non-autoregressive (NAR) machine translation has recently achieved significant improvements, and now outperforms autoregressive (AR) models on some benchmarks, providing an efficient alternative to AR inference. However, while AR translation is often implemented using multilingual models that benefit from transfer between languages and from improved serving efficiency, multilingual NAR models remain relatively unexplored. Taking Connectionist Temporal Classification (CTC) as an example NAR model and Imputer as a semi-NAR model (Saharia et al., 2020), we present a comprehensive empirical study of multilingual NAR. We test its capabilities with respect to positive transfer between related languages and negative transfer under capacity constraints. As NAR models require distilled training sets, we carefully study the impact of bilingual versus multilingual teachers. Finally, we fit a scaling law for multilingual NAR, which quantifies its performance relative to the AR model as model scale increases.

1 Introduction

Non-autoregressive (NAR) models generate output tokens in parallel instead of sequentially, achieving significantly faster inference speed that no longer depends on sequence length. They rely on sequence-level knowledge distillation to reach the quality of autoregressive (AR) models (Gu et al., 2018). As the notion of NAR has expanded to include semi-NAR models that generate their outputs in multiple steps, each time generating several tokens non-autoregressively (Lee et al., 2018; Ghazvininejad et al., 2019), we have begun to see NAR matching the quality of AR. Prior works have benchmarked NAR models for machine translation (MT) on a handful of selected languages like German, Chinese, and Romanian. To efficiently expand this set of languages, it makes sense to explore *multilingual NAR* translation models.

Multilingual MT models (Dong et al., 2015; Firat et al., 2017; Johnson et al., 2017) translate between multiple source and target languages. They offer better parameter efficiency than bilingual models, and they are able to transfer knowledge from high-resource languages to low-resource ones. Therefore they have become an attractive solution for expanding the language coverage of AR MT (Aharoni et al., 2019; Fan et al., 2021). The capability of multilingual modeling is a major feature of the AR regime, and it is one that we should seek to maintain in NAR models.

However, it is unclear to what extent the benefits of multilingual AR models transfer to NAR modeling (Caruana, 1997; Arivazhagan et al., 2019). Do related languages help each other as easily (*positive transfer*)? Do unrelated languages interfere with one another more (*negative transfer*)? Since NAR models tend to trade target-side modeling for improved modeling of the source, the answer to both questions is unclear. Furthermore, NAR modeling raises a new issue of multilingual distillation. To retain the training-time efficiency of multilingual modeling, it is crucial that NAR works well with multilingual teachers; otherwise, the prospect of training many bilingual teachers would greatly increase the effective training cost. It may actually be the case that multilingual teachers are better suited than bilingual ones, as the effective capacity reduction may result in less complex (Zhou et al., 2019) and less multimodal outputs (Gu et al., 2018).

We present an empirical study of multilingual NAR modeling. Taking CTC (Libovický and Helcl, 2018) as our canonical NAR method, and Imputer (Saharia et al., 2020) as our canonical semi-NAR model, we study how they respond to multilinguality in a 6-language scenario designed to emphasize negative transfer, as well as two-language scenarios designed to emphasize positive transfer. In doing so, we make the following contributions:

1. We show that multilingual NAR models suffer

083 more from negative transfer and benefit less
084 from positive transfer than AR models.

- 085 2. We fit a scaling law for our 6-language NAR
086 scenario, demonstrating that this trend contin-
087 ues as model size increases.
- 088 3. We find that multilingual NAR performs
089 equally well with multilingual and bilingual
090 teachers, even in scenarios where the multilin-
091 gual teacher has lower BLEU.

092 Our results indicate that scaling up is not going
093 solve the challenges of multilingual NAR, at least
094 for the models studied here, but our analysis points
095 to promising directions for future work.

096 2 Non-Autoregressive Multilingual NMT

097 Let, $D^l = (x, y) \in X \times Y$ denote the bilingual cor-
098 pus of a language pair, l . Given an input sequence x
099 of length T' , an AR model (Bahdanau et al., 2015;
100 Vaswani et al., 2017) predicts the target y with
101 length T sequentially based on the conditional dis-
102 tribution $p(y_t | y_{<t}, x_{1:T'}; \theta)$. NAR models assume
103 conditional independence in the output token space;
104 that is, they model $p(y_t | x_{1:T'}; \phi)$. Due to this con-
105 ditional independence assumption, training NAR
106 models directly on the true target distribution leads
107 to degraded performance (Gu et al., 2018). Hence,
108 NAR models are typically trained with sequence-
109 level knowledge distillation (Kim and Rush, 2016)
110 to reduce the modeling difficulty.

111 2.1 Non-Autoregressive NMT with CTC

112 In this work, we focus on NAR modelling via
113 CTC (Graves et al., 2006) due to its superior per-
114 formance on NAR generation and the flexibility
115 of variable length prediction (Libovický and Helcl,
116 2018; Saharia et al., 2020; Gu and Kong, 2021).

117 CTC models an alignment a that provides a map-
118 ping between a sequence of predicted and target
119 tokens. Alignments can be constructed by inserting
120 special *blank tokens* (" ") and token repetitions into
121 the target sequence. The alignment is monotonic
122 with respect to the target sequence and is always
123 the same length as the source sequence x . How-
124 ever, in MT, the target sequence y can be longer
125 than the source sequence x . This is handled via
126 upsampling the source sequence x , to s times its
127 original length. An alignment is valid only if when
128 collapsed, i.e., merging repeated tokens and remov-
129 ing blank tokens, it results in the original target

sequence. The CTC loss marginalizes over all pos-
130 sible valid alignments $\Gamma(y)$ compatible with the
131 target y and is defined as:
132

$$p(y | x) = \sum_{a \in \Gamma(y)} \prod_{1 \leq t' \leq T'} p(a_{t'} | x_{1:T'}; \phi).$$

Note that each alignment token $a_{t'}$ is modeled inde-
134 pendently. This conditional independence allows
135 CTC to predict the single most likely alignment
136 non-autoregressively at inference time, which can
137 then be efficiently collapsed to an output sequence.
138 This same independence assumption enables effi-
139 cient minimization of the CTC loss via dynamic
140 programming (Graves et al., 2006). While CTC
141 enforces monotonicity between the alignment and
142 the target, it does not require any cross- or self-
143 attention layers inside the model to be monotonic.
144 Hence, CTC should still be able to model language
145 pairs with different word orders between the source
146 and the target sequence. Following Saharia et al.
147 (2020), we train encoder-only CTC models, using
148 a stack of self-attention layers to map the source
149 sequence directly to the alignments.
150

151 2.2 Iterative Decoding with Imputer

152 IMPUTER (Saharia et al., 2020) extends NAR
153 CTC modeling by iterative refinement (Lee et al.,
154 2018). At each inference step, it conditions on
155 a previous partially generated alignment to emit
156 a new alignment. While IMPUTER, like CTC,
157 generates all tokens at each inference step, only
158 a subset of these tokens are selected to generate
159 a partial alignment, similar to iterative masking
160 approaches (Ghazvininejad et al., 2019). This is
161 achieved by training with marginalization over par-
162 tial alignments:

$$p(y | x) = \sum_{a \in \Gamma(a)} p(a | a_{\text{Mask}}, x; \phi),$$

163 where a_{Mask} is a partially masked input-alignment.
164 At training time, the a_{Mask} alignment is generated
165 using a CTC model trained on the same dataset,
166 and its masked positions are selected randomly.
167 This training procedure enables IMPUTER to it-
168 eratively refine a partial alignment over multiple
169 decoding steps at inference time — consuming its
170 own alignments as input to the next iteration. With
171 $k > 1$ decoding steps, the IMPUTER becomes *semi*-
172 autoregressive, requiring k times more inference
173 passes than pure CTC models.
174

175 IMPUTER differs from Conditional Masked Lan-
176 guage Modeling (CMLM) (Ghazvininejad et al.,

	TGT WORD ORDER	SIZE	SCRIPT DIFFERENCE	WHITE SPACE	SRC LENGTH	TGT LENGTH
EN-KK	SOV	150K	✓	✓	26.7	20.0
EN-DE	SVO/SOV	4.6M	✗	✓	25.7	24.3
EN-PL	SVO	5M	✗	✓	16.2	14.6
EN-HI	SOV	8.6M	✓	✓	18.3	19.8
EN-JA	SOV	17.9M	✓	✗	21.4	25.9
EN-RU	Free	33.5M	✓	✓	23.2	21.5
EN-FR	SVO	38.1M	✗	✓	29.2	32.8

Table 1: Details on training data used. Target word orders are the ones that are dominating within the language according to (Dryer and Haspelmath, 2013), but there may be sentence-specific variations. English follows predominantly SVO (Subject-Verb-Object) order. Size is measured as the number of parallel sentences in the training data. Source (Src) and Target (Tgt) length are averaged across sentences after word-based tokenization.

2019) in that it utilizes the CTC loss instead of the standard cross-entropy loss, removing the need for explicit output length prediction. Also, IMPUTER is an encoder-only model that makes one prediction per source token, just like CTC. The cross-attention component from encoder-decoder is replaced by a simple sum between the embeddings of the source sequence and the input alignment (a_{Mask}) before the first self-attention layer.¹

2.3 Multilingual Modeling

Multilingual AR and NAR models are trained on datasets from multiple language pairs, $\{D^l\}_{l=1}^L$. We prepend each source sequence with the desired target language tag (<2tgt>) and generate a shared vocabulary across all languages (Johnson et al., 2017). The models encode this tag as any other token, and uses it to guide the generation of the output sequence in the desired target language.

2.4 Efficiency

Inference We refrain from wallclock inference time measurements since these are dependent on implementation, low-level optimization and machines (Dehghani et al., 2021). We instead compare generation speed in terms of the number of tokens that get generated per iteration N_{gen} (Kreutzer et al., 2020), which is < 1 for AR models,² T for fully non-autoregressive models like CTC and $\frac{T}{k}$ for iterative semi-autoregressive models like IMPUTER. We acknowledge that other factors like model-depth play a role for inference time, but we assume that both NAR and AR models can be optimized for this aspect (Kasai et al., 2020).

¹We experimented with an encoder-decoder variant of IMPUTER but it did not change the overall output quality in multilingual scenarios or otherwise.

²1 for greedy search, < 1 to account for scoring and expansion of multiple hypotheses in beam search.

Training At training time, NAR models are less efficient than AR models because their quality depends on distillation (Gu and Kong, 2021). Extra cost is incurred to train a teacher model (usually AR) and to use it to decode the training set.

Multilinguality Multilingual models multi-task over language pairs, so that a single multilingual model can replace several bilingual models. Thanks to transfer across languages, model size needs to be increased less than m -fold for modeling m language pairs instead of a single one.

Considering all of the above factors, an ideal model needs only a few iterations (decoder passes or steps), requires no teacher, and covers several languages, while incurring the smallest drop in quality compared to less efficient models. CTC is desirable as it uses only one pass, while IMPUTER gives up some efficiency to improve quality. Both require a teacher, but we can try to reduce the cost by training fewer teachers.

3 Experimental Setup

Data We perform our main experiments on six language pairs, translating from English into WMT-14 German (DE) (Bojar et al., 2014), WMT-15 French (FR) (Bojar et al., 2015), WMT-19 Russian (RU) (Barrault et al., 2019), WMT-20 Japanese (JA), WMT-20 Polish (PL) (Barrault et al., 2020) and Samanantar Hindi (HI) (Ramesh et al., 2021). The lower-resourced WMT-19 English-Kazakh (KK) (Barrault et al., 2019) is used for an additional transfer experiment in Section 5. The properties of the datasets are listed in Table 1. Target word order and writing script notably differ across these languages, so we focus on translating *into* these languages as this is the more challenging direction. A shared sub-word vocabulary of 32k is trained with SentencePiece (Kudo and Richardson, 2018).

MODEL	TEACHER	N_{gen}	EN-FR	EN-DE	EN-PL	EN-RU	EN-HI	EN-JA	AVG.
AR-big		< 1	38.8	29.0	21.4	27.2	34.6	35.4	31.1
multi-AR-big			38.5	27.0	21.6	25.3	32.6	33.6	29.3
Bilingual Models									
AR-base		< 1	38.2	27.6	21.2	26.2	33.8	34.8	30.3
CTC	AR-big	T	35.7	25.2	18.0	21.4	31.6	31.6	27.3
	multi-AR-big		35.1	24.0	17.7	20.8	30.8	28.9	26.2
IMPUTER	AR-big	$\frac{T}{8}$	38.5	27.2	21.2	25.6	32.0	32.0	29.4
Multilingual Models									
multi-AR-base		< 1	35.2	24.8	19.7	23.2	30.8	31.2	27.5
CTC	AR-big	T	31.6	20.5	13.0	17.7	28.2	28.1	23.2
	multi-AR-big		31.2	20.5	13.7	18.0	27.8	27.5	23.1
IMPUTER	AR-big	$\frac{T}{8}$	34.4	22.8	14.9	21.3	29.9	29.6	25.5
	multi-AR-big		34.1	21.2	16.4	21.7	29.9	27.9	25.2

Table 2: Test BLEU scores for multilingual and bilingual AR and NAR models.

The proportion of sub-words allocated for each language is proportional to its data size.

Evaluation Metrics Translation quality is evaluated with BLEU (Papineni et al., 2002) as calculated by Sacrebleu (Post, 2018) with default tokenization except for EN-JA, where we use character-level tokenization.

Architecture We train the IMPUTER model using the same setup as described in Saharia et al. (2020): We follow their base model with $d_{model} = 512$, $d_{hidden} = 2048$, $n_{heads} = 8$, $n_{layers} = 12$, and $p_{dropout} = 0.1$. AR models follow Transformer-base (Vaswani et al., 2017) and have similar parameter counts. We train both models using Adam with learning rate of 0.0001. We train CTC models with a batch size of 2048 and 8192 sentences for 300K steps for the bilingual and multilingual models respectively. We train the IMPUTER using CTC loss using a Bernoulli masking policy for next 300K steps with a batch size of 1024 and 2048 sentences for the bilingual and multilingual models respectively. We upsample the source sequence by a factor of 2 for all our experiments.³ We pick the best checkpoint based on validation BLEU for bilingual models, and the last checkpoint for multilingual models.

Distillation We apply sequence-level knowledge distillation (Kim and Rush, 2016) from AR teacher

³We do not vary the upsampling ratio due to small difference in the performance of the resulting NAR models (see Table 6, Gu and Kong (2021)).

models as widely used in NAR generation (Gu et al., 2018). Specifically, when training the NAR models, we replace the reference sequences during training with translation outputs from Transformer-Big AR teacher model with a beam width of four. We also report the quality of the AR teacher models, both bilingual and multilingual.

4 Negative Transfer Scenario

Our main experiment compares English-to-X models for the six high-resource languages in Table 1. These languages are typologically diverse, and each have enough data so that we do not expect them to benefit substantially from joint modeling. We use this challenging scenario to test the impact of multilingual teachers, and to measure each paradigm’s ability to model several unrelated languages. Results are shown in Table 2.

4.1 Multilingual Teacher Comparison

The top two rows of Table 2 show that in this negative transfer scenario, multilingual teachers have substantially reduced BLEU compared to bilingual teachers. However, as we look at the impact on bilingual students, we see that CTC models trained from the multilingual teacher, multi-AR-big, do not reflect the entirety of this drop in teacher quality when compared to training with the bilingual AR-big. An average teacher gap of -1.8 BLEU is mapped to -1.1 in the corresponding students. The comparison becomes more interesting as we shift to multilingual students: mul-

304 tilingual CTC does not suffer at all from having a
 305 multilingual teacher (average BLEU gap of -0.1),
 306 and multilingual Imputer likewise suffers very lit-
 307 tle (-0.3). These three results taken together sug-
 308 gest that *datasets distilled from multilingual mod-*
 309 *els are likely simpler and easier to model non-*
 310 *autoregressively*, which makes up for their lower
 311 BLEU. Our analysis in Section 4.3 supports this
 312 hypothesis. We hope that highly multilingual mod-
 313 els, trained with similar target language pairs to en-
 314 hance positive transfer (Tan et al., 2019), are even
 315 better suited to serve as teachers for multilingual
 316 NAR models, which we leave to future work.

317 4.2 Multilingual Model Comparison

318 Returning to the “Bilingual Models” section of Ta-
 319 ble 2 with AR-big teachers, we can see that we
 320 have reproduced the expected results of Saharia
 321 et al. (2020). Bilingual CTC performs well for a
 322 fully NAR method, but does not reach AR quality.
 323 IMPUTER ably closes the gap with AR, surpassing
 324 or coming within 0.2 BLEU of the AR-base mod-
 325 els on 3/6 language pairs, with the largest gap in
 326 performance for the distant EN-JA. Does this story
 327 hold as we move to multilingual NAR students?

328 To understand each model’s multilingual capa-
 329 bilities, we can compare its bilingual performance
 330 to its multilingual performance. Comparing AR-
 331 base to multilingual AR-base gives us a baseline
 332 average drop of -2.8 BLEU, confirming that this
 333 is indeed a difficult multilingual scenario that leads
 334 to negative transfer. Comparing bilingual CTC to
 335 multilingual CTC, both with AR-big teachers, we
 336 see an average drop of -4.1 . This larger drop indi-
 337 cates that CTC *suffers more from negative interfer-*
 338 *ence than its AR counterpart*. We hypothesize that
 339 CTC models need more capacity than AR models
 340 to achieve similar multilingual performance, mo-
 341 tivating our scaling law experiments in Section 6.
 342 Performing the same bilingual-to-multilingual com-
 343 parison for IMPUTER shows a similar -3.9 average
 344 drop due to negative transfer. So although IM-
 345 PUTER is indeed substantially better than CTC, it
 346 does not seem to be necessarily better suited for
 347 multilingual modeling in this difficult scenario.

348 4.3 How do the distilled datasets differ?

349 Table 3 summarizes different statistics for the origi-
 350 nal (R) and distilled datasets from both multilin-
 351 gual (M) and bilingual (B) AR teacher models.
 352 We report the number of types and average se-
 353 quence length (in tokens) for the target side of

PROPERTY	R	B	M
EN-FR			
# TYPES	522K	430K	396K
AVG. LENGTH	32.8	31.2	29.2
COMPLEXITY	1.529	1.167	0.944
FRS	0.463	0.541	0.536
BLEU (Train)	-	40.8	37.8
EN-DE			
# TYPES	812K	616K	573K
AVG. LENGTH	24.3	23.4	22.2
COMPLEXITY	1.243	0.819	0.709
FRS	0.490	0.606	0.605
BLEU (Train)	-	35.0	26.4
EN-PL			
# TYPES	636K	516K	503K
AVG. LENGTH	14.6	13.4	12.7
COMPLEXITY	1.435	0.942	0.591
FRS	0.590	0.678	0.695
BLEU (Train)	-	26.3	22.0
EN-RU			
# TYPES	636K	516K	503K
AVG. LENGTH	21.5	20.5	19.5
COMPLEXITY	1.083	0.882	0.819
FRS	0.640	0.719	0.716
BLEU (Train)	-	43.2	40.0
EN-HI			
# TYPES	346K	200K	185K
AVG. LENGTH	19.8	18.8	17.8
COMPLEXITY	1.438	1.256	1.138
FRS	0.347	0.363	0.366
BLEU (Train)	-	34.6	28.0
EN-JA			
# TYPES	547K	440K	402K
AVG. LENGTH	25.9	23.5	22.2
COMPLEXITY	1.541	1.369	1.338
FRS	0.344	0.337	0.340
BLEU (Train)	-	35.9	30.6

Table 3: Comparison of datasets distilled from bilin-
 gual (B) or multilingual (M) AR models on a subset of
 1M samples: Multilingual distilled datasets have fewer
 types, are less complex and more monotonic than bilin-
 gual distilled datasets.

the dataset. We compute the complexity of the
 dataset based on probabilities from a statistical
 word aligner (Zhou et al., 2019). The FRS (Talbot
 et al., 2011) score represents the average fuzzy re-
 ordering score over all the sentence pairs for the
 respective language pair as measured in Xu et al.
 (2021), with higher values suggesting that the target

354
 355
 356
 357
 358
 359
 360

is more monotonic with the source sequence. We also report BLEU for the distilled datasets relative to the original training references.

The datasets distilled from the bilingual AR models (B) are shorter, less complex, have reduced lexical diversity (in number of types) and are more monotonic compared to the original corpora (R), which corroborates findings from prior work (Zhou et al., 2019; Xu et al., 2021). One exception is EN-JA, where the distilled translations are slightly less monotonic than the original references. Moving to multilingual teachers (M), the resulting datasets have further reduced types, are shorter and less complex than those distilled from bilingual teachers. In particular, their monotonicity increased (FRS) for the more distant language pairs, EN-JA and EN-HI. As shown in Xu et al. (2021), reduced lexical diversity and reordering complexity can help NAR models to learn better alignments between source and target, improving the translation quality of the outputs.

4.4 Which translation errors are made?

In this section, we analyze quantitatively how the output quality of NAR models differs across language pairs when trained in isolation (bilingual) or with other language pairs (multilingual).

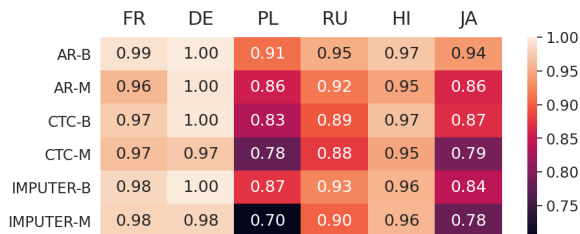


Figure 1: Brevity penalty scores for bilingual (-B) and multilingual (-M) models, the closer to 1 the better.

Effect of Length Figure 1 shows the brevity penalty (BP) scores (Papineni et al., 2002) for all languages. EN-PL and EN-JA have lowest BP scores across the board, meaning that their translations are shorter than the references. Manual inspection reveals that this could be attributed to the subject pronouns being dropped in both of these target languages. Multilingual modeling results in shorter outputs relative to bilingual models for both AR and NAR models and most language pairs. While IMPUTER models tend to have fewer issues with output length compared to CTC models, they still lag behind AR models, suggesting that the length might need to be controlled explicitly for

these language pairs (Gu and Kong, 2021).

Invalid Words CTC frequently generates *invalid* words, i.e. tokens that are not present in the target side of the bitext but are being composed from multiple sub-words. These sub-words represent alternative translations that the model fails to distinguish. In the Hindi example below, the invalid (or made-up) word in the sentence is marked in red. The correct word should be जहरीले as the dependent vowel “ी” can only be used once.

Hindi: इससे ग्रामीण महिलाओं को जहरीले धुएं से मुक्ति मिली है।

English: This has relieved the rural women from the poisonous smoke.

Figure 2 reports the percentage of sequences that include at least one invalid word in the test set. CTC generates many invalid words compared to both AR and IMPUTER, with multilingual modeling leading to an average increase in invalid words by 37%. The shared vocabulary of the multilingual model results in shorter sub-words, hence longer sequences, and the conditionally independent generation leads to more clashing adjacent sub-words.⁴ IMPUTER’s iterative decoding alleviates this for some languages. Increasing the number of iterations could help, but would also erode the efficiency arguments that make NAR models attractive.

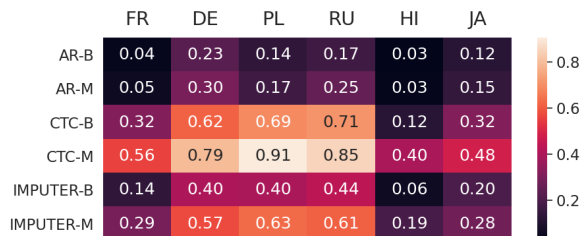


Figure 2: % of outputs with invalid words for bilingual (-B) and multilingual (-M) models, the lower the better.

5 Positive Transfer Scenario

In this section we present two experimental setups designed to emphasize positive transfer, where languages are related and training data is limited.

English→{German, French} To isolate the effect of transfer via multilingual modelling, we relax the capacity bottleneck and competition for parameters: We combine the two most related languages

⁴One might hope to alleviate this by increasing vocabulary size, but preliminary experiments showed that an increased vocabulary was less efficient in improving quality than increasing overall model size, which is explored in Section 6.

(DE, FR) and give them smaller, balanced training sets (1M sentences). We compare bilingual and multilingual AR and NAR models trained on this reduced data.

Table 4 shows that NAR models benefit from training with multiple language pairs in this relaxed scenario — all models exhibit positive transfer (in green). IMPUTER achieves higher positive transfer than CTC for both languages, but lags behind the AR multilingual model in EN-FR. However, for EN-FR the bilingual IMPUTER is already ahead of the bilingual AR model by 0.4 BLEU.

MODEL	EN-DE	EN-FR
Bilingual Models		
AR	22.8	27.7
CTC	21.5	26.5
IMPUTER	22.8	28.1
Multilingual Models		
AR	24.3 +1.5	29.0 +1.3
CTC	22.1 +0.6	26.9 +0.4
IMPUTER	23.7 +1.3	28.5 +0.4

Table 4: Results on subsampled (1M) training data.

English→{Russian, Kazakh} Does this positive transfer survive data imbalance? We test the performance of the multilingual NAR model on the low-resource task of translating English into Kazakh, for which the size of clean training data is insufficient to train a bilingual AR model from scratch. We instead distill translations from the publicly available multilingual AR model PRISM (Thompson and Post, 2020). We then pair it with the higher-resource but related language Russian to encourage positive transfer to Kazakh. Given the huge difference in data sizes for Russian and Kazakh (see Table 1), we sample training data from the two languages based on the data size scaled by a temperature value τ , $p_l^{1/\tau}$ (Arivazhagan et al., 2019), where, $p_l = \frac{D_l}{\sum_k D_k}$. We experiment with multiple temperature values (1, 3, 5, 10, 20) and pick the best value ($\tau = 5$; $p_{RU}^{1/\tau} = 0.75$, $p_{KK}^{1/\tau} = 0.25$) based on the performance on the validation set.

As can be seen in Table 5, both AR and CTC show positive transfer when translating into Kazakh when trained in combination with Russian. The multilingual CTC model is able to improve over the bilingual CTC model, but the overall quality of the outputs is very low compared to

MODEL	TEACHER	EN-KK	EN-RU
PRISM	-	8.9	27.0
Bilingual Models			
AR	PRISM	4.4	-
CTC	PRISM	1.2	-
Multilingual Models			
AR	PRISM	7.1 +2.7	26.0
CTC	PRISM	2.8 +1.6	20.4

Table 5: Results on *English* → *Kazakh, Russian*.

the teacher model (BLEU: -5.3). This experiment showcases that *current NAR models do not perform well on very low-resource language pairs* and might need further data augmentation (Anonymous, 2022) or transfer from other similar languages.⁵

6 Impact of Model Scale

We hypothesized in Section 4 that CTC might require more capacity than AR models. If we increase the parameters for NAR models sufficiently, could we reach AR quality? Scaling laws can characterize the relationship between MT output quality, the cross-entropy loss and the number of parameters used for training the model (Ghorbani et al., 2021; Gordon et al., 2021).

We derive the relationship between BLEU and the number of parameters for our AR and CTC models directly from the scaling laws proposed by Gordon et al. (2021) and Ghorbani et al. (2021) as follows:

$$L(N) \approx L_0 + \alpha_n (1/N)^{\alpha_k} \text{ (Ghorbani et al., 2021)}$$

$$\text{BLEU}(L) \approx C e^{-kL} \text{ (Gordon et al., 2021)}$$

$$\text{BLEU}(N) \approx a e^{-b(1/N)^c} \text{ (this work)}$$

where L is the test loss, $\{\alpha_n, \alpha_k, L_0, C, k\}$ are fitted parameters from previous power laws, and $\{a, b, c\}$ are the collapsed fitted parameters of our power law. Ghorbani et al. (2021)’s L_0 corresponds to the irreducible loss of the data, which becomes a in our formulation.

Setup We train seven different models with varying capacity for AR and CTC models. The number of layers and model size are varied as: (6, 128), (6, 256), (12, 256), (12, 512),⁶ (24, 512), (12, 1024), (24, 1024). The feed-forward size is $4\times$ the model size. AR models have equal numbers

⁵We do not train IMPUTER for KK as the quality of the distilled dataset and alignments from CTC is very low.

⁶Size for experiments in Section 4.

of encoder and decoder layers. The number of attention heads is given by $(8/(512/\text{Model Size}))$. For a fair comparison, we use the same number of layers and dimensions and train both AR and CTC models on distilled outputs from a bilingual teacher (AR-big).

Results Figure 3 shows the fitted parameters using the scaling law, which can almost perfectly describe the relationship between the number of parameters and the development BLEU (R^2 : 0.99) averaged across all language pairs from Section 4. When the number of parameters is below 10M, AR and CTC model yield similar translation quality. However, the gap in BLEU increases with the number of parameters. We can also see that CTC needs many more parameters to achieve comparable BLEU to AR models and plateaus early at a BLEU of 26.7, while AR models plateau at 30.8. By projecting the curves out to 1 billion parameters, we show that increasing the capacity of NAR is insufficient to reach the quality of AR models.

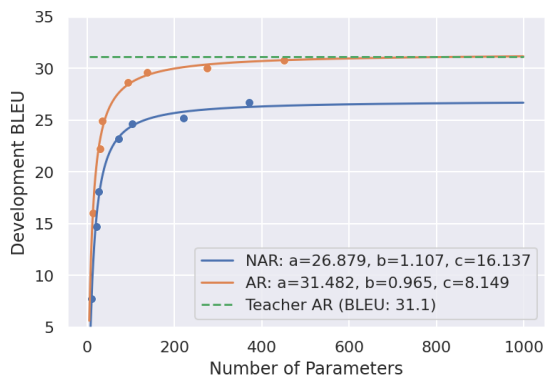


Figure 3: BLEU versus number of parameters and fitted power-law curves (R^2 AR: 0.99, R^2 CTC: 0.99).

7 Related Work

Multiple approaches with varying architectures (Gu et al., 2018, 2019; Chan et al., 2020; Xu and Carpuat, 2021), custom loss functions (Ghazvininejad et al., 2020; Du et al., 2021) and training strategies (Ghazvininejad et al., 2019; Qian et al., 2021) have been used to enable parallel generation of output tokens for MT with sequence-level knowledge distillation as one of the key ingredient in the training of NAR models. While most prior work focuses on bilingual NAR modeling, we investigate multilingual NAR MT models. One limitation of our study is that we choose one representative system for NAR and semi-NAR modeling, rather than exploring the full breadth of NAR options.

Both supervised and unsupervised (Sun et al., 2020) MT have benefitted from training with multiple languages, especially those that have very little (Siddhant et al., 2020) to no training data (Zhang et al., 2020). However, multilingual modelling has not yet received any attention in the NAR literature. Concurrent to our work, Anonymous (2022) investigate a non-autoregressive multilingual MT model with a code-switch decoder. They show that adding code-switched back-translation data to the training of multilingual models improves performance. Our work instead focuses on understanding multilinguality for both the student and the teacher model in the context of NAR training, without using any additional data augmentation strategies.

Recent works have derived empirical scaling laws that govern the relationship between the performance of language models or translation models to the scale of the model or dataset (Kaplan et al., 2020; Hernandez et al., 2021; Bahri et al., 2021; Gordon et al., 2021; Ghorbani et al., 2021). We extend this work to examine the impact of model scale on multilingual NAR MT.

8 Conclusion

The capability for multilingual MT is a valuable feature of AR models, therefore, we have tested NAR models for that same capability. We focus on challenging scenarios to discover potential weaknesses and to identify areas for future work. In a relaxed setting with little interference between languages and balanced data, multilingual NAR models nicely exhibit positive transfer, practically closing the gap to AR models with a few decoding iterations. However, we do not see the same positive transfer in a true low-resource scenario. Experiments in a six-language scenario reveal that multilingual NAR models suffer proportionally more from negative interference than AR models. Our derived scaling laws show that scaling up CTC model parameters is not a sufficient remedy. Our analysis identified two issues that hurt translation quality and worsen with multilinguality, namely output length control and the generation of invalid words. We have also shown beneficial properties of using multilingual teachers for distillation. We hope that this work will serve as a call for increased focus on multilingual modeling in NAR research.

584
585
586
587

588
589
590
591

592
593
594
595
596
597

598
599
600
601
602

603
604
605

606
607
608
609
610
611
612
613
614
615
616
617

618
619
620
621
622
623
624
625

626
627
628
629
630
631
632

633
634
635
636
637
638
639

References

Roei Aharoni, Melvin Johnson, and Orhan Firat. 2019. Massively multilingual neural machine translation. In *Proceedings of NAACL-HLT*, pages 3874–3884.

Anonymous. 2022. [Non-autoregressive models are better multilingual translators](#). In *Submitted to The Tenth International Conference on Learning Representations*. Under review.

Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun, Mia Xu Chen, Yuan Cao, George Foster, Colin Cherry, et al. 2019. Massively multilingual neural machine translation in the wild: Findings and challenges. *arXiv preprint arXiv:1907.05019*.

Dzmitry Bahdanau, Kyung Hyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *3rd International Conference on Learning Representations, ICLR 2015*.

Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. 2021. Explaining neural scaling laws. *arXiv preprint arXiv:2102.06701*.

Loïc Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. 2020. [Findings of the 2020 conference on machine translation \(WMT20\)](#). In *Proceedings of the Fifth Conference on Machine Translation*, pages 1–55, Online. Association for Computational Linguistics.

Loïc Barrault, Ondřej Bojar, Marta R Costa-Jussa, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, et al. 2019. Findings of the 2019 conference on machine translation (wmt19). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61.

Ondřej Bojar, Christian Buck, Christian Federmann, Barry Haddow, Philipp Koehn, Johannes Leveling, Christof Monz, Pavel Pecina, Matt Post, Herve Saint-Amand, et al. 2014. Findings of the 2014 workshop on statistical machine translation. In *Proceedings of the ninth workshop on statistical machine translation*, pages 12–58.

Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Barry Haddow, Matthias Huck, Chris Hokamp, Philipp Koehn, Varvara Logacheva, Christof Monz, Matteo Negri, Matt Post, Carolina Scarton, Lucia Specia, and Marco Turchi. 2015. [Findings of the 2015 workshop on statistical machine translation](#). In *Proceedings of the Tenth Workshop on Statistical*

Machine Translation, pages 1–46, Lisbon, Portugal. Association for Computational Linguistics. 640
641

Rich Caruana. 1997. Multitask learning. *Machine learning*, 28(1):41–75. 642
643

William Chan, Chitwan Saharia, Geoffrey Hinton, Mohammad Norouzi, and Navdeep Jaitly. 2020. Imputer: Sequence modelling via imputation and dynamic programming. In *International Conference on Machine Learning*, pages 1403–1413. PMLR. 644
645
646
647
648

Mostafa Dehghani, Anurag Arnab, Lucas Beyer, Ashish Vaswani, and Yi Tay. 2021. The efficiency misnomer. *arXiv preprint arXiv:2110.12894*. 649
650
651

Daxiang Dong, Hua Wu, Wei He, Dianhai Yu, and Haifeng Wang. 2015. [Multi-task learning for multiple language translation](#). In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1723–1732, Beijing, China. Association for Computational Linguistics. 652
653
654
655
656
657
658
659

Matthew S. Dryer and Martin Haspelmath, editors. 2013. *WALS Online*. Max Planck Institute for Evolutionary Anthropology, Leipzig. 660
661
662

Cunxiao Du, Zhaopeng Tu, and Jing Jiang. 2021. Order-agnostic cross entropy for non-autoregressive machine translation. *arXiv preprint arXiv:2106.05093*. 663
664
665
666

Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, et al. 2021. Beyond english-centric multilingual machine translation. *Journal of Machine Learning Research*, 22(107):1–48. 667
668
669
670
671
672

Orhan Firat, Kyunghyun Cho, Baskaran Sankaran, Fatos T. Yarman Vural, and Yoshua Bengio. 2017. [Multi-way, multilingual neural machine translation](#). *Computer Speech & Language*, 45:236–252. 673
674
675
676

Marjan Ghazvininejad, Vladimir Karpukhin, Luke Zettlemoyer, and Omer Levy. 2020. [Aligned cross entropy for non-autoregressive machine translation](#). *CoRR*, abs/2004.01655. 677
678
679
680

Marjan Ghazvininejad, Omer Levy, Yinhan Liu, and Luke Zettlemoyer. 2019. Mask-predict: Parallel decoding of conditional masked language models. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6112–6121. 681
682
683
684
685
686
687
688

Behrooz Ghorbani, Orhan Firat, Markus Freitag, Ankur Bapna, Maxim Krikun, Xavier Garcia, Ciprian Chelba, and Colin Cherry. 2021. Scaling laws for neural machine translation. *arXiv preprint arXiv:2109.07740*. 689
690
691
692
693

694	Mitchell A Gordon, Kevin Duh, and Jared Kaplan.	<i>Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5774–5782, Online. Association for Computational Linguistics.	749
695	2021. Data and parameter scaling laws for neural machine translation.		750
696	In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 5915–5922.		751
698			
699	Alex Graves, Santiago Fernández, Faustino Gomez, and Jürgen Schmidhuber. 2006. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In <i>Proceedings of the 23rd international conference on Machine learning</i> , pages 369–376.		752
700			753
701			754
702			755
703			
704			
705	Jiatao Gu, James Bradbury, Caiming Xiong, Victor O.K. Li, and Richard Socher. 2018. Non-autoregressive neural machine translation . In <i>International Conference on Learning Representations</i> .		756
706			757
707			758
708			759
709	Jiatao Gu and Xiang Kong. 2021. Fully non-autoregressive neural machine translation: Tricks of the trade . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 120–133, Online. Association for Computational Linguistics.		760
710			761
711			762
712			763
713			764
714			765
715	Jiatao Gu, Changhan Wang, and Junbo Zhao. 2019. Levenshtein transformer . In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, <i>Advances in Neural Information Processing Systems 32</i> , pages 11179–11189. Curran Associates, Inc.		766
716			767
717			768
718			769
719			770
720			771
721	Danny Hernandez, Jared Kaplan, Tom Henighan, and Sam McCandlish. 2021. Scaling laws for transfer. <i>arXiv preprint arXiv:2102.01293</i> .		772
722			773
723			774
724	Melvin Johnson, Mike Schuster, Quoc V. Le, Maxim Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat, Fernanda Viégas, Martin Wattenberg, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2017. Google’s multilingual neural machine translation system: Enabling zero-shot translation . <i>Transactions of the Association for Computational Linguistics</i> , 5:339–351.		775
725			776
726			777
727			778
728			779
729			780
730			781
731	Jared Kaplan, Sam McCandlish, T. J. Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeff Wu, and Dario Amodei. 2020. Scaling laws for neural language models. <i>ArXiv</i> , abs/2001.08361.		782
732			783
733			784
734			785
735			786
736	Jungo Kasai, Nikolaos Pappas, Hao Peng, James Cross, and Noah Smith. 2020. Deep encoder, shallow decoder: Reevaluating non-autoregressive machine translation. In <i>International Conference on Learning Representations</i> .		787
737			788
738			789
739			790
740			791
741	Yoon Kim and Alexander M. Rush. 2016. Sequence-level knowledge distillation . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 1317–1327, Austin, Texas. Association for Computational Linguistics.		792
742			793
743			794
744			795
745			796
746	Julia Kreutzer, George Foster, and Colin Cherry. 2020. Inference strategies for machine translation with conditional masking . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1098–1108.		797
747			798
748			799
			800
			801
			802
			803
			804
			805

806 [neural machine translation](#). In *Proceedings of the*
807 *58th Annual Meeting of the Association for Computa-*
808 *tational Linguistics*, pages 2827–2835, Online. As-
809 sociation for Computational Linguistics.

810 Haipeng Sun, Rui Wang, Kehai Chen, Masao Utiyama,
811 Eiichiro Sumita, and Tiejun Zhao. 2020. [Knowledge](#)
812 [distillation for multilingual unsupervised neural ma-](#)
813 [chine translation](#). In *Proceedings of the 58th Annual*
814 *Meeting of the Association for Computational Lin-*
815 *guistics*, pages 3525–3535, Online. Association for
816 Computational Linguistics.

817 David Talbot, Hideto Kazawa, Hiroshi Ichikawa, Jason
818 Katz-Brown, Masakazu Seno, and Franz Josef Och.
819 2011. A lightweight evaluation framework for ma-
820 chine translation reordering. In *Proceedings of the*
821 *Sixth Workshop on Statistical Machine Translation*,
822 pages 12–21.

823 Xu Tan, Jiale Chen, Di He, Yingce Xia, Tao Qin, and
824 Tie-Yan Liu. 2019. Multilingual neural machine
825 translation with language clustering. In *Proceedings*
826 *of the 2019 Conference on Empirical Methods in*
827 *Natural Language Processing and the 9th Interna-*
828 *tional Joint Conference on Natural Language Pro-*
829 *cessing (EMNLP-IJCNLP)*, pages 963–973.

830 Brian Thompson and Matt Post. 2020. [Paraphrase gen-](#)
831 [eration as zero-shot multilingual translation: Disen-](#)
832 [tangling semantic similarity from lexical and syntac-](#)
833 [tic diversity](#). In *Proceedings of the Fifth Conference*
834 *on Machine Translation*, pages 561–570, Online. As-
835 sociation for Computational Linguistics.

836 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob
837 Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
838 Kaiser, and Illia Polosukhin. 2017. Attention is all
839 you need. In *Advances in neural information pro-*
840 *cessing systems*, pages 5998–6008.

841 Weijia Xu and Marine Carpuat. 2021. Editor: An edit-
842 based transformer with repositioning for neural ma-
843 chine translation with soft lexical constraints. *Trans-*
844 *actions of the Association for Computational Lin-*
845 *guistics*, 9:311–328.

846 Weijia Xu, Shuming Ma, Dongdong Zhang, and Ma-
847 rine Carpuat. 2021. [How does distilled data com-](#)
848 [plexity impact the quality and confidence of non-](#)
849 [autoregressive machine translation?](#) In *Findings of*
850 *the Association for Computational Linguistics: ACL-*
851 *IJCNLP 2021*, pages 4392–4400, Online. Associa-
852 tion for Computational Linguistics.

853 Biao Zhang, Philip Williams, Ivan Titov, and Rico Sen-
854 nrich. 2020. [Improving massively multilingual neu-](#)
855 [ral machine translation and zero-shot translation](#). In
856 *Proceedings of the 58th Annual Meeting of the Asso-*
857 *ciation for Computational Linguistics*, pages 1628–
858 1639, Online. Association for Computational Lin-
859 guistics.

860 Chunting Zhou, Jiatao Gu, and Graham Neubig.
861 2019. Understanding knowledge distillation in non-
862 autoregressive machine translation. In *International*
863 *Conference on Learning Representations*.