# Scaling Laws for Neural Machine Translation

**Behrooz Ghorbani**
ghorbani@google.com

**Orhan Firat**
orhanf@google.com

**Markus Freitag**
freitag@google.com

**Ankur Bapna**
ankurbpn@google.com

**Maxim Krikun**
krikun@google.com

**Xavier Garcia**
xgarcia@google.com

**Ciprian Chelba**
ciprianchelba@google.com

**Colin Cherry**
colincherry@google.com

## Abstract

We present an empirical study of scaling properties of encoder-decoder Transformer models used in neural machine translation (NMT). We show that cross-entropy loss as a function of model size follows a certain scaling law. Specifically (i) We propose a formula which describes the scaling behavior of cross-entropy loss as a bivariate function of encoder and decoder size, and show that it gives accurate predictions under a variety of scaling approaches and languages; we show that the total number of parameters alone is not sufficient for such purposes. (ii) We observe different power law exponents when scaling the decoder vs scaling the encoder, and provide recommendations for optimal allocation of encoder/decoder capacity based on this observation. (iii) We also report that the scaling behavior of the model is acutely influenced by *composition bias* of the train/test sets, which we define as any deviation from naturally generated text (either via machine generated or human translated text). We observe that natural text on the target side enjoys scaling, which manifests as successful reduction of the cross-entropy loss. (iv) Finally, we investigate the relationship between the cross-entropy loss and the quality of the generated translations. We find two different behaviors, depending on the nature of the test data. For test sets which were originally translated from target language to source language, both loss and BLEU score improve as model size increases. In contrast, for test sets originally translated from source language to target language, the loss improves, but the BLEU score stops improving after a certain threshold. We release generated text from all models used in this study.

## 1 Introduction

Scaling properties of neural networks have long been an intriguing topic of study [2, 3]. Along with the practical success of modern neural networks at scale, theoretical understanding of the factors governing the quality and training dynamics of large neural networks has also being developing [1, 31, 11, 12, 8, 20, 5, 27]. In particular, scaling model sizes, datasets and the total computation budget has been identified as a reliable approach to improve generalization performance on several machine learning tasks. For many of these tasks the scaling behavior of neural networks is highly predictable; model fit or test loss can be described precisely as a function of its number of parameters [18, 21, 16, 17, 31]. Neural machine translation (NMT) has long enjoyed the benefits of scaling [19, 4, 25], but studies investigating the scaling behavior of NMT models are missing. We present

the first large-scale systematic study of scaling laws for encoder-decoder Transformer models applied to NMT [36]. [1] [2]

We start with highlighting the major differences between decoder-only language models, where the majority of the previous work has focused, and encoder-decoder (conditional) language models applied to NMT. The two differ along a few crucial dimensions. The first difference results from the very nature of the separate architectures being used, i.e. decoder-only vs encoder-decoder. The presence of separate architectural components complicates the study of scaling properties due to the increased degree of freedom. Second, contrary to language modeling, the task of machine translation is conditional: the task is predictive rather than fully generative. Furthermore, this prediction task is ambiguous: there is no one right answer for a given source, and translations can vary substantially depending on the translator's incentives. This manifests itself as different scaling benefits for different test sets. To take an extreme example, a test set translated by someone who writes nearly word-for-word translations may benefit less from model scaling than one translated by someone who considers each translation a work of art. In this work, these differences in difficulty coincide with the translation direction of the test set; that is, whether the source was translated into the target (source-original) or vice versa (target-original). Source-original data has translated text on the target side, which contains several artifacts of "translationese" that distinguish it from text originally written in that language, often lacking the diversity and complexity of "natural" text [24], while target-original data requires the prediction of more complex natural text on the target side. Finally, unlike language models, NMT is evaluated on metrics that quantify generation quality against reference translations (for eg. BLEU) [28] instead of evaluating model fit (perplexity) on an evaluation set.

In this paper, we aim to provide empirical answers to the following research questions:

1. **Does the encoder-decoder architecture for NMT share the same scaling law function as the language models?** Contrary to previous work on LM, we show that a univariate law depending on the total number of parameters in the network does not adequately describe the scaling behavior of NMT models. Our scaling laws parameterize the cross entropy loss as a bivariate function of the number of encoder parameters and the number of decoder parameters as separate variables. Our results indicate that the scaling behavior is largely determined by the total capacity of the model, and the capacity allocation between the encoder and the decoder.

2. **How does the naturalness of source/target side data affect scaling behavior?** We study the effect of naturalness of the source and target text, both for training and evaluation. When evaluating with target side natural text, scaling the model capacity continues improving model quality throughout our range of measurements. On the other hand, improvements on cross-entropy saturate (or reaches the irreducible error region) on source side natural evaluation sets even for moderately-sized models.

3. **Do scaling improvements in cross-entropy translate into corresponding improvements in generation quality?** Finally we study the relationship between generation quality and cross-entropy and how their correlation changes as we: (i) Scale different components of the model (encoder vs decoder) and (ii) Evaluate on source-natural or target-natural evaluation sets.

Our results on multiple language pairs and training/test data compositions validate that **model scaling predictably improves the cross-entropy on validation data**. However, our findings also raise several questions regarding the effect of naturalness of training and evaluation text and how cross-entropy eventually relates with generation quality for auto-regressive generative models.

## 2 Effect of Scaling on Cross-Entropy

### 2.1 Experimental setting

**Model Architectures and Training** We train a series of pre-layer norm Transformer networks with varying sizes [39]. Models are trained with per-token cross-entropy loss and Adafactor optimizer [35].

---

[1] An initial version of this study was submitted to NeurIPS 2021.

[2] A few weeks before the publication of this manuscript on Arxiv, [13] appeared on OpenReview. While both papers study scaling laws for NMT, our studies focus on different parameter regimes (393K-56M vs 100M-3.5B).

All models are trained with a fixed batch-size of 500k tokens and dropout rate of 0.1 for residuals, feed-forward activations and attention. All models are trained to near convergence for 500k training steps. Details of the model hyper-parameters are described in Appendix A.

**Model Scaling**  Transformer architecture consists of Transformer Blocks: a cascade of self-attention, cross-attention and feed-forward layers, each having multiple adjustable hyper-parameters (e.g. model-dimension, number of attention heads, attention projection dimension etc.). Considering the combinatorial expansion of the search space for scaling each one [29, 26, 37], in this study we choose to vary only the total number of Transformer Blocks, while keeping the internal hyper-parameters intact across different scales. In other words, we scale the depth of the Transformers while keeping width and other variables fixed. We use GPipe pipeline parallelism for scaling [19] thanks to its flexible API across various depths.

In an encoder-decoder Transformer architecture for NMT, depth scaling can naturally be implemented by varying encoder-decoder blocks independently or symmetrically. Hence, we examine the change in the cross-entropy loss as the number of parameters increase with three depth scaling approaches:

*Encoder Scaling*: vary encoder depth (2 to 64) while decoder depth is fixed (6 layers).

*Decoder Scaling*: vary decoder depth (2 to 64) while encoder depth is fixed (6 layers).

*Symmetric Scaling*: increasing decoder and encoder layers together (from 2 to 64), i.e. the number of Transformer Blocks in the encoder and decoder being equal.

For all experiments, configuration of the individual layers is unchanged: the model dimension, width of the feed-forward layer, and number of attention heads are fixed respectively at 1024, 8192, and 16. [3] Each encoder layer adds approximately 20M parameters to the model while each decoder layer adds around 25M parameters. In this section, we train 95 such models which scale the encoder / decoder size by approximately a factor of 32 (from roughly 40M parameters to 1.5B parameters). Following the convention in this literature, we do not count the parameters in the embedding and softmax layers towards the model size.

**Language Pairs**  We report results on two language pairs, English→German and German→English, using an in-house web-crawled dataset with around 2.2 billion sentence pairs (approximately 55 billion tokens) for both translation directions. This dataset provides a large enough training set to ensure the dataset size is not a bottleneck in the model performance.

**Evaluation Sets**  We use a variety of test sets for evaluation covering different domains: (i) Web-Domain (ii) News-Domain (iii) Wikipedia (iv) Patents. The news-domain test sets come from the WMT2019 [6] evaluation campaign (newstest2019) for all language pairs. The other test sets are internal test sets representing the different domains, ranging from 500 to 5000 sentence pairs. For each domain, we randomly sample sentences in the source language and use professional translators to generate a reference translation in the target language. Throughout the paper, we will refer this type of test sets as *source-original* as the source sentences have been crawled from the web while the reference translations are added later. For most of the domains, we also have a *target-original* counterpart which is generated in the opposite direction: Sentences are crawled in the target language and human translated into the source language. Earlier work [9, 10, 14] showed that it is important to differentiate between the two different kinds of test sets as the style of natural sentences and human (or machine) translations (translationese) is quite different. Cross-entropy loss is evaluated on the different test sets during training. To reduce the variation caused by the parameter fluctuations at the end of the training, we present the median loss over the last 50k steps of the training as the final loss.

## 2.2  Results

Figure 1 shows the empirical evolution of the test loss on the Web-Domain test sets for encoder and decoder scaling for English→German. To compare the empirical results with the scaling laws present in the literature for decoder only models [21, 16], we have fitted a power law of the form

$$\hat{L}(N) = \alpha N^{-p} + L_\infty \tag{1}$$

---

[3]A complete description of the model architecture is provided in Appendix A

to the data. [4] Here, $N$ is the total number of parameters outside of embedding / softmax layers and $\{\alpha, p, L_\infty\}$ are the fitted parameters of the power law. As Figure 1 suggests, scaling the encoder has different effects on the test loss compared to scaling the decoder. As such, simple power-law relations similar to Eq. (1) that only consider the total number of parameters, fail to capture the correct scaling behavior of the model.
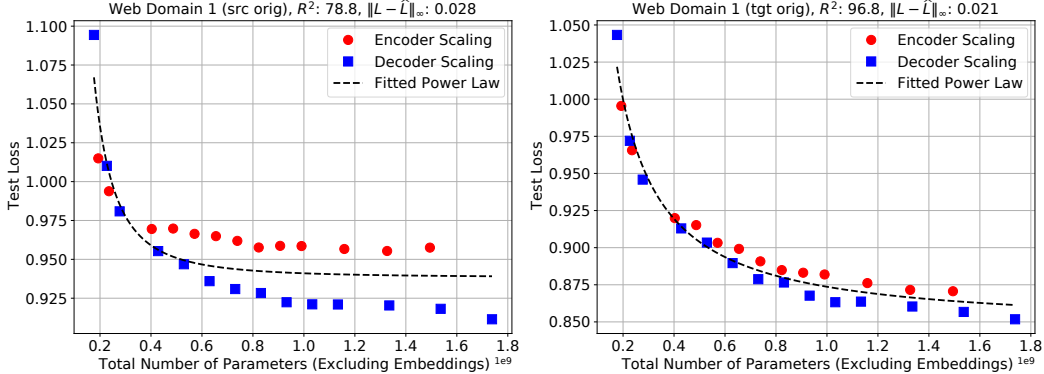


Figure 1: Evolution of the test loss as a function of the total model parameters for English→German. Scaling the encoder has different effects compared to scaling the decoder. As such, traditional power laws of type expressed as in Eq. (1) are unable to capture the correct scaling behavior. R-squared $(100 \times \frac{\text{explained variance}}{\text{total variance}})$ and maximum absolute deviation ($\|\cdot\|_\infty$) are reported for each fit.

**Proposed Scaling Law**     To tackle this issue, we present a new scaling law that reflects the encoder-decoder nature of the architecture as well as the bilingual format of the data. Let $N_e$ and $N_d$ be the number of non-embedding parameters in the encoder and the decoder respectively. Then, our proposed scaling law has the form

$$\hat{L}(N_e, N_d) = \alpha \left( \frac{\bar{N}_e}{N_e} \right)^{p_e} \left( \frac{\bar{N}_d}{N_d} \right)^{p_d} + L_\infty \tag{2}$$

where $\{\alpha, p_e, p_d, L_\infty\}$ are test set specific (fitted) parameters. $\bar{N}_e$ and $\bar{N}_d$ are fixed normalization parameters corresponding to the number of encoder / decoder parameters in our baseline 12-layer encoder-decoder model.[5] In this formulation, $\alpha$ corresponds to the maximum loss reduction (as compared to the baseline model) that one can hope from scaling, while $p_e$ and $p_d$ are the scaling exponents for encoder and decoder respectively. $L_\infty$ corresponds to the irreducible loss of the data.

Figure 2 presents the fit achieved by the proposed scaling law on Web-Domain test sets. The dashed lines describe the fit of the scaling law given in Eq. (2) to the empirical (encoder & decoder scaling) data. The plots suggest that our proposed scaling law is able to simultaneously capture both encoder and decoder scaling behaviors.

To validate the (out-of-sample) prediction power of these scaling laws, we compare their predictions with empirical loss values achieved by our symmetric scaling models. Figure 3 presents this comparison. The plots suggest that the predictions of the scaling law match the empirical (out-of-sample) results with remarkable accuracy. These results suggest that the predictions of the scaling law are not sensitive to the scaling approach; the scaling law fitted on encoder / decoder scaling data is able to almost perfectly predict the scaling behavior of symmetric models. Notice that the parameter range of symmetric scaling models is much larger than either of the encoder or decoder scaling models. Nevertheless, our fitted scaling laws are able to extrapolate effectively to models with sizes beyond the ones used for fitting them.

Figures 2 & 3 suggest that the functional form proposed in Eq. (2) captures the scaling behavior of English→German models accurately. To verify the robustness of our proposed scaling law, we

---

[4]Details of the curve fitting procedure are presented in Appendix E.
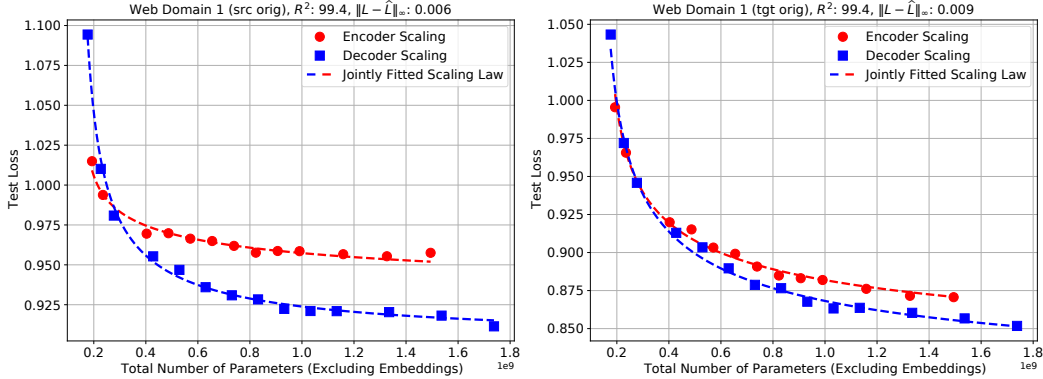[5]Corresponds to 6-layer encoder - 6-layer decoder.

Figure 2: Evolution of log-perplexity as a function of the model size for English→German models. Eq. (2) is jointly fitted to the empirical loss values from encoder scaling and decoder scaling experiments. Our proposed scaling law is able to capture more than $99\%$ of the variation in the data. We anticipate some fluctuations around the predicted trend (with estimated standard deviation of $0.003$) caused by the randomness in the training pipeline (see Appendix C). We observe similar results for our other test sets (see Figures 12 & 13 of the appendix).
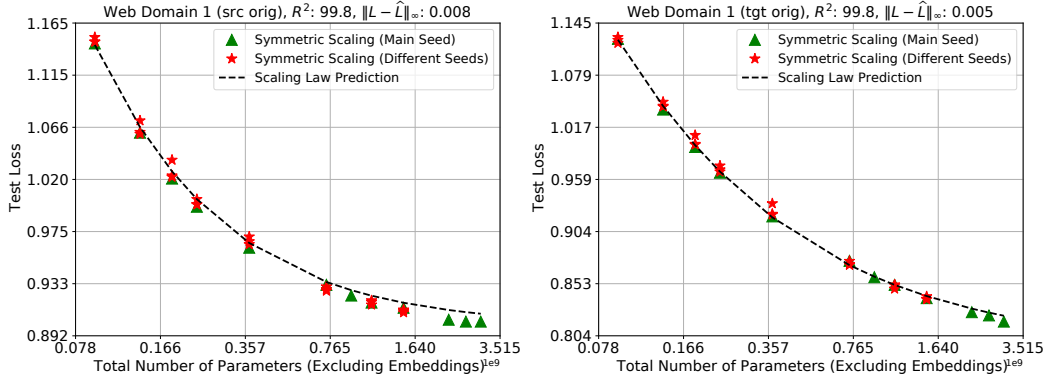


Figure 3: Comparison of the (out-of-sample) predictions of the scaling law with the empirical test loss values from symmetric scaling English→German models. Eq. (2) is fitted only using the encoder / decoder scaling data and then just evaluated on the symmetric scaling model parameters. Our proposed scaling law is able to almost fully capture the variation in the data ($R^2 = 99.8\%$) even though it has not been fitted on it. To examine the randomness in the results, we have repeated a subset of training runs with $4$ different random seeds (see Appendix C for more details). We observe similar results for our other test sets (see Figure 14 of the appendix).

evaluate it on an additional translation task namely, German→English (De→En). Figure 4 depicts the corresponding scaling behavior on Web-Domain test set. Similar to the En→De case, our proposed functional form is able to closely capture the scaling behavior of the models.

## 2.3 Analysis

The above results suggest that scaling law formalized in Eq. (2) captures the scaling behavior of the Transformer NMT models in multiple language pairs. As such, we can study the fitted coefficients to fully understand the scaling properties of these models. Figure 5 presents the fitted coefficients for all of the test sets under consideration. Several observations are in order:

**Decoder vs Encoder Scaling:** On all our test sets, the decoder exponents were observed to be larger than the encoder exponents, $p_d > p_e$. As a result, when improving the test loss is concerned, it is much more effective to scale the decoder rather than the encoder. This is contrary to the usual
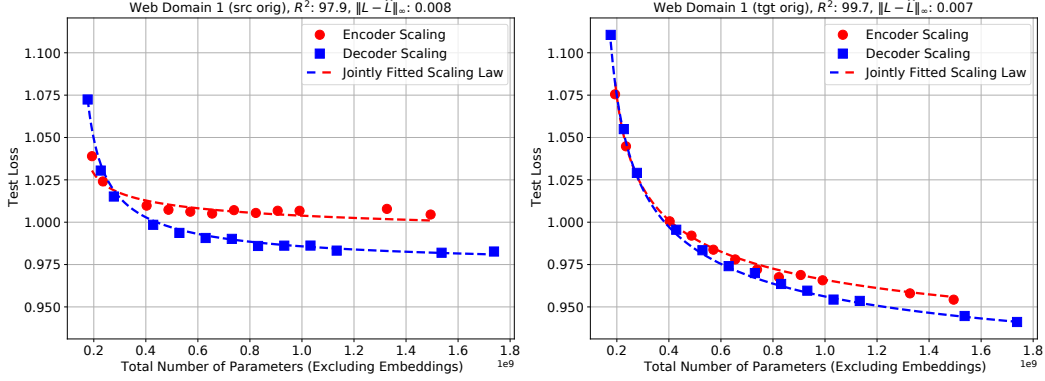
Figure 4: Fitted scaling laws for De→En translation task. The scaling law (2) is jointly fitted to the empirical loss value from encoder scaling and decoder scaling experiments. Similar to En→De case, the law is able to describe the empirical scaling behavior of the models with high accuracy. See Figures 15 & 16 in the appendix for the fit on other test sets.

practice; due to latency considerations, many practitioners train NMT models with deep encoders and shallow decoders [22]. Our results suggest this practice could be sub-optimal in terms of loss reduction. Proposition 1 below leverages Eq. (2) to provide guidance on how to allocate parameters in between the encoder and decoder optimally. The proof is presented in Appendix D.

**Proposition 1** (Optimal Scaling). *Assume the loss performance of the model is described by Eq. (2). Let $B$ denote the budget for total number of parameters. Then, the optimal encoder / decoder sizes (denoted respectively by $N_e^*$ and $N_d^*$) are:*

$$N_e^* = \frac{p_e}{p_e + p_d}B, \qquad N_d^* = \frac{p_d}{p_e + p_d}B. \qquad (3)$$

*In addition, when optimally scaling the model, the scaling law reduces to:*

$$\hat{L}_{opt}(B) = \alpha^* B^{-(p_d + p_e)} + L_\infty, \qquad \alpha^* \equiv \alpha\left(\frac{\bar{N}_e(p_e + p_d)}{p_e}\right)^{p_e}\left(\frac{\bar{N}_d(p_e + p_d)}{p_d}\right)^{p_d}. \qquad (4)$$


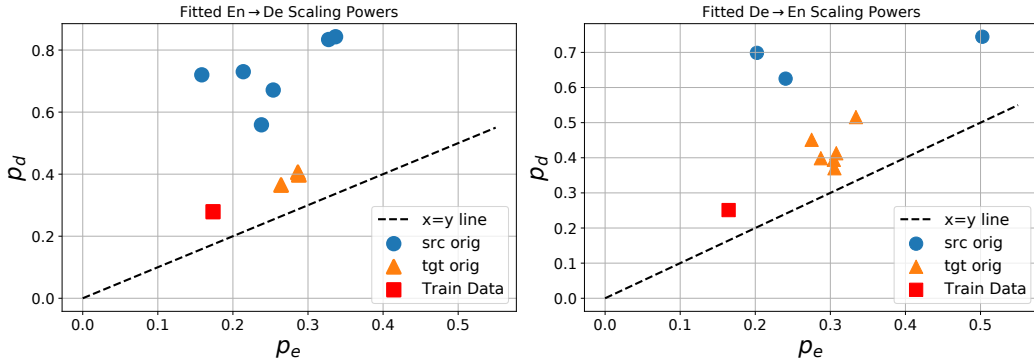
Figure 5: Fitted scaling law exponents for all our test sets. Across all the test sets under consideration, we observe $p_d > p_e$.

Proposition 1 suggests that when $\frac{N_e}{N_d} = \frac{p_e}{p_d}$, one can achieve the best possible scaling behavior for the task. Inspection of the functional form of Eq. (2) suggests that as long as $N_d/N_e$ is fixed as the model scales (i.e. the encoder and decoder grow proportionally together), the optimal scaling exponent, $(p_e + p_d)$, can be achieved, albeit with a potentially sub-optimal multiplicative constant, $\alpha^\#$. To examine how significant this sub-optimality can be, in Figure 6, we compare the multiplicative constants resulting from proportional scaling of the encoder and decoder with different values of

$N_d/N_e$. The results suggest that as long as the parameter allocation is not extremely far from $(N_e^*, N_d^*)$, the scaling scheme is approximately optimal. In particular, symmetrically scaling the encoder and decoder layers, which yields $N_d/N \approx 0.55$, is barely distinguishable from the optimal scaling scheme described in Proposition 1. In contrast, lopsided models which heavily favor the encoder or the decoder achieve a much worse multiplicative constant.
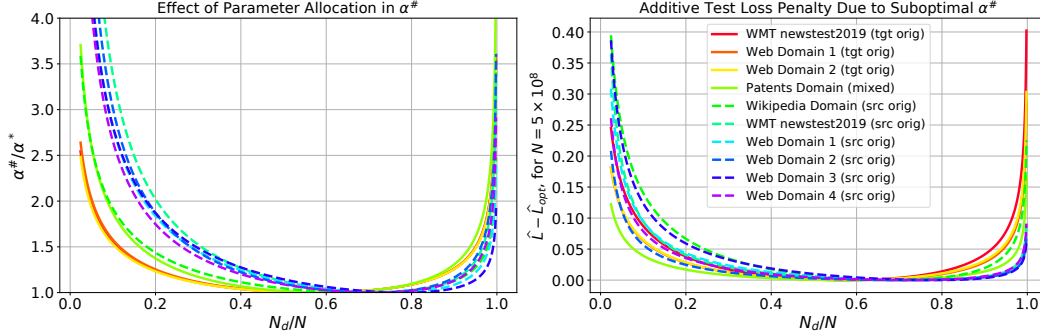


Figure 6: We use our fitted scaling laws to evaluate the effect of encoder / decoder parameter allocation ratio when proportionally scaling the encoder and the decoder. Left: $\alpha^\# / \alpha*$ for different parameter allocation schemes. Right: The predicted additive loss penalty, $(\widehat{L} - \widehat{L}_{opt})$, for a model with $5 \times 10^8$ total (non-embedding) parameters. Each line corresponds to a different test set.

## 3   Effect of Dataset Composition Bias on Scaling Behavior

Translation deals with the problem of mapping a sequence in one language into another language. A good translation should not only be adequate and fluent, but should ideally also adopt the style of a sentence naturally written in the target language. This necessitates MT models to make sense of natural looking inputs and generate natural looking outputs. As mentioned in Section 2, the examples used to train or test NMT models carry a critical bias, which we refer to as *composition bias*. Composition bias is introduced because of the unavailability of source-target examples (pairs) that are both natural in the accessible data generating distribution. For any given naturally generated text in a language, the corresponding text in the other language is either translated by humans, introducing *translationese* bias or translated by other machine translation systems, introducing *MT* bias. We consider both biases affecting the problem from a similar angle, hence we bundle them and call it composition bias. While machine translation by design has composition bias in the training/test sets employed [10, 30], its effect on model scaling is unknown. In this section we investigate the role of composition bias in scaling and identify critical factors playing role.

We caution the reader to not take the composition bias as a problem specific to NMT. In fact as most training corpora in NMT are web-crawled, they can contain machine translation output on either the source or target side. Considering the growth of generated content in the web by machine learning models [6] [7], it is not improbable that a proportion of the content collected and used by machine learning models is going to be biased by other models that are continuously generating content.

**The Effect of Test Set Construction:**   We will first take a look at the impact of composition bias on the test sets used in this study and then investigate the influence on the training set. Figure 5 shows the fitted scaling law coefficient for all of our test sets. The coefficients suggests that the scaling powers for source-original test sets are drastically different from those of target-original test sets. This behavior is in direct contrast with language modeling setting [21] where it was observed that the evaluation on different test sets merely acted as a scaling penalty that only changed the multiplicative constants of the scaling law.

To elucidate this phenomenon further, in Figure 7, we compare the scaling trends for different source and target original test sets. To factor out the effect of the data domain, we present one source original

---

[6]https://openai.com/blog/gpt-3-apps/
[7]https://blog.google/products/translate/one-billion-installs/

and one target original test set for each domain. Several observations are in order: Test sets with a similar composition approach (source or target original) have a qualitatively similar scaling behavior. However, scaling behavior is vastly different between the two composition approaches. Reducible loss quickly decays to zero for source original test sets. In fact, we observe that scaling our baseline 6L-6L model by a factor of $2.5$ is sufficient for ensuring that reducible loss is below $0.05$ for all source original test sets. In contrast, on target original test sets, the loss decays much more slowly with model size. For comparison, to ensure that reducible loss is below $0.05$ for all target original test sets, we estimate that the baseline model has to be scaled up by a factor of $11$.
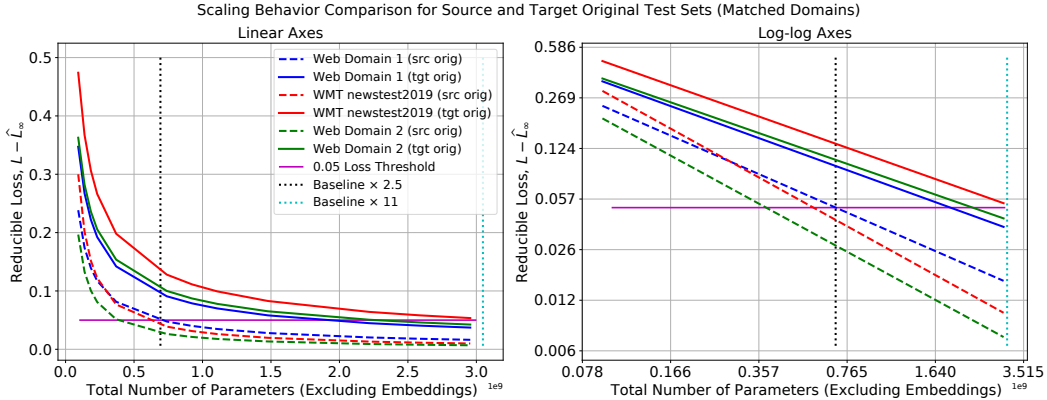


Figure 7: A comparison of scaling behavior across source and target original test sets. We use our fitted scaling laws to estimate the evolution of reducible loss for each test set. All scaling trends correspond to symmetrically scaling the encoder and decoder layers.

Because of this behavior, the value of larger models in NMT is closely tied to their evaluation sets: On source original test sets, due to larger scaling exponents, even moderate increases in model size are sufficient for pushing the reducible loss close to zero. Hence, beyond a few hundred million parameters, there is no benefit in increasing the model size. In contrast, for target original test sets, which generally have smaller scaling exponents, large models are needed to push the reducible loss to zero.

**The Effect of Training Set Construction:** The results of the previous section suggest that the construction of the test data plays a key role in the scaling behavior of the model. Now, we briefly examine the role of underlined training data construction on the scaling behavior. To do this, we generate two En→De datasets, that were not used in the previous experiments. One fully target original and another completely source original.

To generate the target original dataset, we compile a set of German documents from the web. Documents are screened to ensure the data is not machine generated. We use a Hybrid model (with 380M parameters) [7] to back-translate (BT) these documents to English. Similarly, for the source original data, we collect human generated English documents and (forward) translate them to German using a hybrid model (with approximately 327M parameters). Both datasets provide us with approximately 2.2 billion training examples. We mimic the experimental setup of Section 2.

Note that even though these datasets are not human generated, they reflect important aspects of training large NMT models. Many modern NMT datasets are harvested from the web and as a result, are contaminated with machine generated data. Moreover, many popular data augmentation algorithms such as Back Translation [34], sequence level distillation [23] and self training [15] purposefully add machine generated data into the training pipeline in order to take advantage of monolingual data.

Figure 8 describes the scaling behavior for models trained on target-original data. We observe that even though larger models are successful in reducing the training loss, they are unable to improve the test loss after roughly $400M$ parameters. Once this size threshold is exceeded, models overfit the training data and the test loss starts to deteriorate across all of our test sets. We hypothesize that this size threshold corresponds to the capacity of the original back-translation model. This assertion suggests that in order for back-translation to be beneficial for training large models, it has to be

performed with a models with comparable capacity or higher. Although quite intriguing, we leave the verification of this hypothesis to future work.
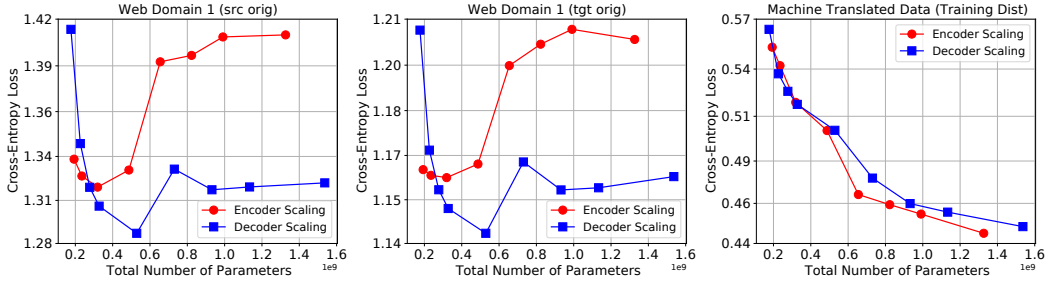


Figure 8: Scaling behavior of models trained on back-translated data. Right: Increasing the model size successfully reduces the loss on the training distribution. However, on the test data (left and center) the loss increases after approximately $400M$ parameters.

Figure 9 paints another interesting picture for the models trained on the source-original data only, implying the target side having the composition bias, expected to be simpler, dull and not rich in its content, in short - not natural looking. As experiments suggest, even our smallest models are able to achieve extremely low loss values (roughly $0.16$), with an apparent overfitting pattern. We believe the same phenomenon is also related to the "data simplification" effect sought by non-autoregressive models in NMT [40].
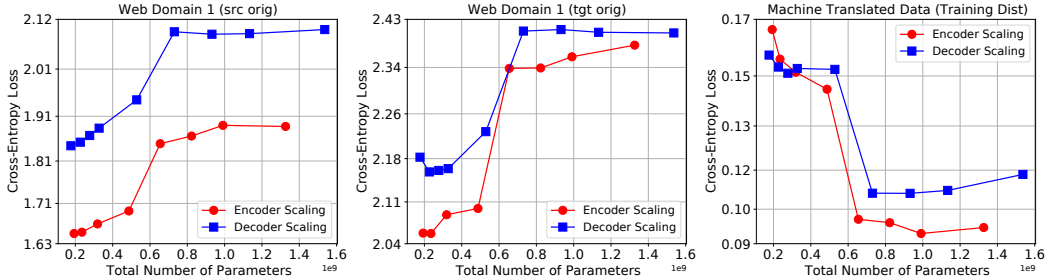


Figure 9: Scaling behavior of models trained on forward translated data. Left / center: early stopping test loss on Web-Domain. Right: loss at the end of the training for a subset of the training data.

## 4 Evolution of Generation Quality

We examine the effects of scaling on the output quality as measured by BLEU score [8]. For the analysis of this section, we focus on output generated via beam search [38]. For tractability purposes, we do not attempt to tune the (many) hyper-parameters of beam-search for each model. Instead, we use the configuration optimized for the baseline model (listed in Appendix F) in all the decoding tasks.

Figure 10 presents the co-evolution of BLEU score and cross-entropy loss throughout the training for all of our models. Depending on the construction of the test sets, two different empirical behaviors emerge. On target-original test sets, larger models are able to improve (lower) the test loss. These improvements in the loss are accompanied with consistent improvements (increases) in BLEU score. In fact, we observe that a simple power law of the form

$$\text{BLEU} = c_B L^{-p_B}, \qquad c_B, \ p_B > 0. \tag{5}$$

can capture the relationship between BLEU score and cross-entropy loss for high-quality models. [9]

---

[8]We computed the BLEU scores using an internal reimplementation of Moses scorer: `mteval-v13a.pl`.

[9]We observe certain deviations from this trend for smaller models and for early checkpoints. We document these deviations in Appendix F.

In contrast, on source-original test sets, this relationship is absent; larger models consistently achieve better test losses, however, beyond a certain threshold, BLEU scores begin to deteriorate. Figures 21 and 22 exhibit that this phenomenon is not due to over-training; the BLEU score gap between large and small models is persistent throughout the training.

To ensure that this observation truly reflects the generation quality of the models (as opposed to potential biases of BLEU score), we repeat our analysis with BLEURT score [32, 33]. The results are presented in Figure 11. As the figure suggests, BLEURT scores closely mirror the behavior of BLEU scores with respect to model scaling.

A careful look at the left-subplots of Figures 10 & 11 brings up another interesting trend. At similar values of the test loss, encoder-scaled models result in better generation quality compared to decoder-scaled models. This findings agrees with previous work that relied on encoder-scaling when optimizing for BLEU and inference latency [22]. Whether these differences in the effects of encoder-scaling and decoder-scaling are caused by insufficient search algorithms, or just different model fits from different architectural priors is left to future work.
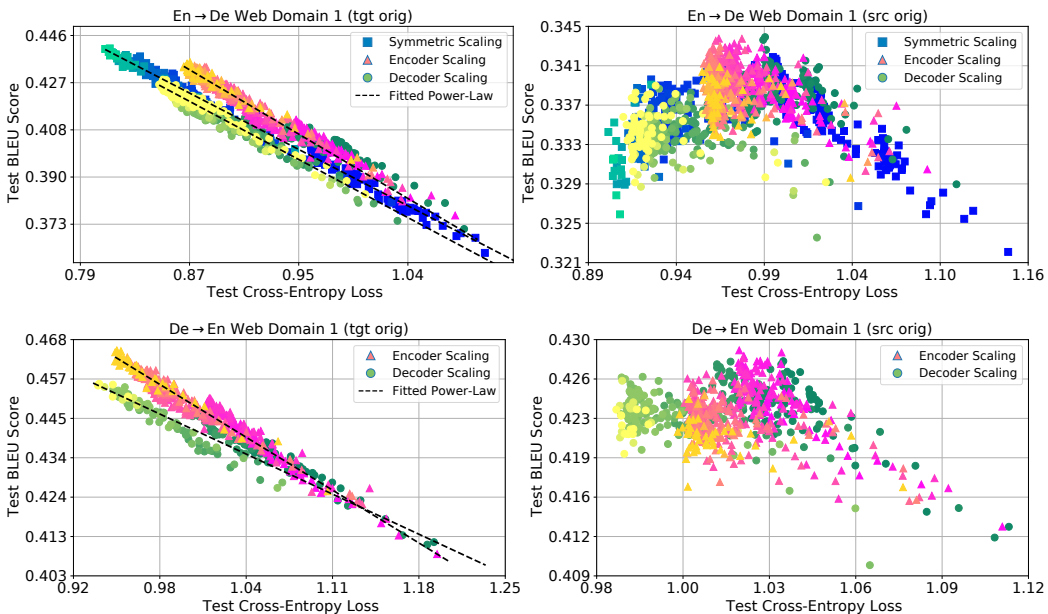


Figure 10: Log-log plot of the evolution of BLEU score as a function of cross-entropy loss for different models. For each scaling approach, warmer colors represent larger models. Each individual color represents different checkpoints of a single model during training. On target original data (left column), improvements to cross-entropy loss lead to consistent improvements in BLEU score. Dashed lines correspond to fit achieved by Eq. (5). The relationship breaks down for source original data (right column). More examples of this phenomenon are presented in Appendix F.

## 5  Conclusion and Limitation

In this work we have attempted to quantify the evolution of model quality as a function of model capacity for encoder-decoder NMT models.

While a univariate scaling law describing the cross-entropy as a function of the total number of parameters in the model is insufficient, a bivariate law treating the number of encoder and decoder parameters as separate variables adequately describes the scaling behavior of these models under various scaling strategies. We validate this behavior on 2 language pairs and on a variety of evaluation sets with different compositions. Whether this behavior is intrinsic to the encoder-decoder architecture, or arising from the nature of the NMT task, requires further study.

Next, we demonstrate that this scaling behavior is highly dependent on the composition of the evaluation data, specifically on whether the source or the target sentences are "original". Our findings
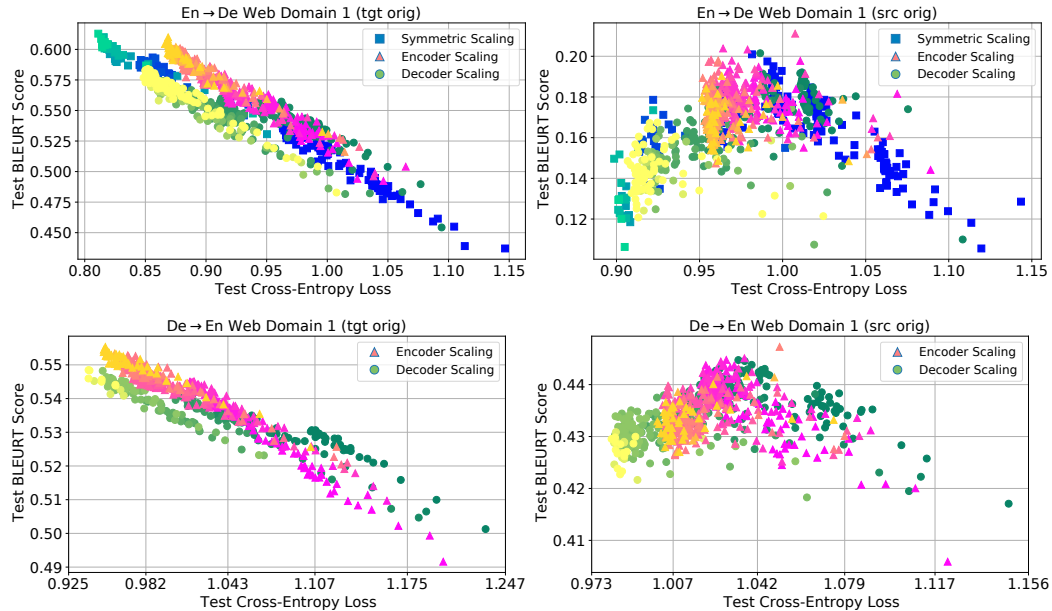
Figure 11: The evolution of BLEURT score as a function of cross-entropy loss for different models. For each scaling approach, warmer colors represent larger models. Each individual color represents different checkpoints of a single model during training. On target original data (left column), improvements to cross-entropy loss lead to consistent improvements in BLEURT score. The relationship breaks down for source original data (right column). More examples are provided in Appendix F.

indicate that target-original evaluation sets continue benefiting from model scaling throughout our range of measurements, while the reducible error on source-original evaluation sets quickly saturates to 0. This could be an artifact of the lack of diversity in translated text; a simpler target distribution doesn't require much capacity to model while generating fluent or natural-looking text could benefit much more from scale.

We also study how the composition of training data affects the scaling behavior of models. When training on target-original (back-translated) text, model quality keeps improving until a point after which the trend saturates. In our study the capacity where saturation manifests first is perilously close to the capacity of the model used for back-translation, indicating that the capacity of the generative model used to generate synthetic text might have a role to play, but this requires further investigation. When training on source-original text, even low-capacity models are sufficient to reach the irreducible loss region, painting a gloomy picture for synthetic data. While we have explored these ideas in the context of machine translation, given the proliferation of generative models this problem will likely be a challenge for future practitioners training on web-scraped monolingual datasets as well. For low-resource languages, the proliferation of machine translated text is already a problem given that a significant portion of web text in these languages is machine translated.

Finally, we attempt to understand how generation quality evolves with the improvements in cross-entropy resulting from model scaling. As with our previous findings, dataset composition plays a major role in determining the trends. For source-original evaluation sets, the correlation between cross-entropy and generation quality breaks down. On target-original evaluation, we observe an inverse correlation between cross-entropy and BLEU/BLEURT, suggesting that improved model fit results in a corresponding improvement in generation quality. The slope of this relationship is different for encoder-scaling and decoder-scaling, with encoder-scaled models performing better on BLEU/BLEURT than decoder-scaled models, at the same level of cross-entropy loss. Whether this is an artifact of our search strategy (beam search, tuned to a 6L-encoder 6L-decoder model) or the difference in architectural priors is something that requires further investigation.

Our findings suggest that scaling behavior of encoder-decoder NMT models is predictable, but the exact formulation of scaling laws might vary depending on the particular architecture or task being

11

studied. Our empirical findings also raise concerns regarding the effect of synthetic data on model scaling and evaluation, and how proliferation of machine generated text might hamper the quality of future models trained on web-text.

## Acknowledgments

## References

[1] Madhu S. Advani and Andrew M. Saxe. High-dimensional dynamics of generalization error in neural networks, 2017.

[2] Subutai Ahmad and Gerald Tesauro. Scaling and generalization in neural networks: a case study. Advances in Neural Information Processing Systems, 1:160–168, 1988.

[3] S. Amari, Naotake Fujita, and S. Shinomoto. Four types of learning curves. Neural Computation, 4:605–618, 1992.

[4] Naveen Arivazhagan, Ankur Bapna, Orhan Firat, Dmitry Lepikhin, Melvin Johnson, Maxim Krikun, Mia Xu Chen, Yuan Cao, George F. Foster, Colin Cherry, Wolfgang Macherey, Zhifeng Chen, and Yonghui Wu. Massively multilingual neural machine translation in the wild: Findings and challenges. CoRR, abs/1907.05019, 2019.

[5] Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural scaling laws, 2021.

[6] Loïc Barrault, Ondřej Bojar, Marta R Costa-Jussa, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, et al. 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, 2019.

[7] Mia Xu Chen, Orhan Firat, Ankur Bapna, Melvin Johnson, Wolfgang Macherey, George Foster, Llion Jones, Niki Parmar, Mike Schuster, Zhifeng Chen, Yonghui Wu, and Macduff Hughes. The best of both worlds: Combining recent advances in neural machine translation, 2018.

[8] Kyunghyun Cho. Scaling laws of recovering bernoulli. Blog Post, 2020.

[9] Markus Freitag, Isaac Caswell, and Scott Roy. APE at Scale and Its Implications on MT Evaluation Biases. In Proceedings of the Fourth Conference on Machine Translation, pages 34–44, Florence, Italy, August 2019. Association for Computational Linguistics.

[10] Markus Freitag, David Grangier, and Isaac Caswell. BLEU might be guilty but references are not innocent. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 61–71, Online, November 2020. Association for Computational Linguistics.

[11] Mario Geiger, Arthur Jacot, Stefano Spigler, Franck Gabriel, Levent Sagun, Stéphane d' Ascoli, Giulio Biroli, Clément Hongler, and Matthieu Wyart. Scaling description of generalization with number of parameters in deep learning. Journal of Statistical Mechanics: Theory and Experiment, 2020(2):023401, Feb 2020.

[12] Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Linearized two-layers neural networks in high dimension, 2020.

[13] Mitchell A Gordon, Kevin Duh, and Jared Kaplan. Data and parameter scaling laws for neural machine translation. In ACL Rolling Review - May 2021, 2021.

[14] Yvette Graham, Barry Haddow, and Philipp Koehn. Statistical power and translationese in machine translation evaluation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 72–81, 2020.

[15] Junxian He, Jiatao Gu, Jiajun Shen, and Marc'Aurelio Ranzato. Revisiting self-training for neural sequence generation, 2020.

[16] Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative modeling. arXiv preprint arXiv:2010.14701, 2020.

[17] Danny Hernandez, Jared Kaplan, Tom Henighan, and Sam McCandlish. Scaling laws for transfer, 2021.

[18] Joel Hestness, Sharan Narang, Newsha Ardalani, Gregory Diamos, Heewoo Jun, Hassan Kianinejad, Md Patwary, Mostofa Ali, Yang Yang, and Yanqi Zhou. Deep learning scaling is predictable, empirically. arXiv preprint arXiv:1712.00409, 2017.

[19] Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Xu Chen, HyoukJoong Lee, Jiquan Ngiam, Quoc V Le, Yonghui Wu, et al. Gpipe: Efficient training of giant neural networks using pipeline parallelism. In NeurIPS, 2019.

[20] Marcus Hutter. Learning curve theory. CoRR, abs/2102.04074, 2021.

[21] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. arXiv preprint arXiv:2001.08361, 2020.

[22] Jungo Kasai, Nikolaos Pappas, Hao Peng, James Cross, and Noah A. Smith. Deep encoder, shallow decoder: Reevaluating non-autoregressive machine translation, 2021.

[23] Yoon Kim and Alexander M. Rush. Sequence-level knowledge distillation, 2016.

[24] Moshe Koppel and Noam Ordan. Translationese and Its Dialects. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies - Volume 1, pages 1318–1326, 2011.

[25] Dmitry Lepikhin, HyoukJoong Lee, Yuanzhong Xu, Dehao Chen, Orhan Firat, Yanping Huang, Maxim Krikun, Noam Shazeer, and Zhifeng Chen. {GS}hard: Scaling giant models with conditional computation and automatic sharding. In International Conference on Learning Representations, 2021.

[26] Yoav Levine, Noam Wies, Or Sharir, Hofit Bata, and Amnon Shashua. Limits to depth efficiencies of self-attention. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.

[27] Bruno Loureiro, Cédric Gerbelot, Hugo Cui, Sebastian Goldt, Florent Krzakala, Marc Mézard, and Lenka Zdeborová. Capturing the learning curves of generic features maps for realistic data sets with a teacher-student model, 2021.

[28] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. BLEU: a Method for Automatic Evaluation of Machine Translation. In Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, pages 311–318. Association for Computational Linguistics, 2002.

[29] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. Journal of Machine Learning Research, 21(140):1–67, 2020.

[30] Parker Riley, Isaac Caswell, Markus Freitag, and David Grangier. Translationese as a language in "multilingual" nmt. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7737–7746, 2020.

[31] Jonathan S Rosenfeld, Amir Rosenfeld, Yonatan Belinkov, and Nir Shavit. A constructive prediction of the generalization error across scales. In International Conference on Learning Representations, 2019.

[32] Thibault Sellam, Dipanjan Das, and Ankur P Parikh. Bleurt: Learning robust metrics for text generation. arXiv preprint arXiv:2004.04696, 2020.

[33] Thibault Sellam, Amy Pu, Hyung Won Chung, Sebastian Gehrmann, Qijun Tan, Markus Freitag, Dipanjan Das, and Ankur P Parikh. Learning to evaluate translation beyond english: Bleurt submissions to the wmt metrics 2020 shared task. arXiv preprint arXiv:2010.04297, 2020.

[34] Rico Sennrich, Barry Haddow, and Alexandra Birch. Improving neural machine translation models with monolingual data. In ACL (1), 2016.

[35] Noam Shazeer and Mitchell Stern. Adafactor: Adaptive learning rates with sublinear memory cost. In International Conference on Machine Learning, pages 4596–4604. PMLR, 2018.

[36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. CoRR, abs/1706.03762, 2017.

[37] Noam Wies, Yoav Levine, Daniel Jannai, and Amnon Shashua. Which transformer architecture fits my data? a vocabulary bottleneck in self-attention, 2021.

[38] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. Google's neural machine translation system: Bridging the gap between human and machine translation. arXiv preprint arXiv:1609.08144, 2016.

[39] Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tieyan Liu. On layer normalization in the transformer architecture. In International Conference on Machine Learning, pages 10524–10533. PMLR, 2020.

[40] Chunting Zhou, Graham Neubig, and Jiatao Gu. Understanding knowledge distillation in non-autoregressive machine translation, 2021.

# A   Architecture and Hyper-Parameter Details

As described in Section 2, all our models use a similar configuration for their Transformer Blocks. In particular, we fix model dimension to $1024$, feed-forward layer dimension to $8192$, number of attention heads to $16$, and attention hidden dimension to $1024$. Our models use a sentence-piece vocabulary of size $32000$.

**Regularization:**   We use a dropout of $0.1$ for residuals, feed-forward activations and attention. Models are trained with label smoothing of magnitude $0.1$. To improve the training stability, all models use logit clipping of $10$.

**Optimizer:**   We use Adafactor [35] optimizer for training our models. We use $40$k linear warm-up steps and an inverse square root learning rate schedule. For Adafactor we used momentum with $0.9$ and factored second moment to save memory.

Table 1 (and Table 2 resp.) describes the parameter decomposition of the encoder scaling (decoder scaling) models. Table 3 describes the parameter counts for the symmetric scaling models. The largest model we used (64L-64L) has more than 3 billion parameters while the smallest model we used (2L-2L) has only 92M non-embedding parameters.

Table 1: Parameter decomposition of the encoder scaling models. The total number of parameters includes 98M parameters representing the softmax and embedding layers.

| Encoder | | Decoder | | |
|---|---|---|---|---|
| Layers | Parameters | Layers | Parameters | Total Parameters |
| 2 | 42M | 6 | 151M | 291M |
| 4 | 84M | 6 | 151M | 333M |
| 8 | 168M | 6 | 151M | 417M |
| 12 | 252M | 6 | 151M | 501M |
| 16 | 336M | 6 | 151M | 585M |
| 20 | 420M | 6 | 151M | 669M |
| 24 | 504M | 6 | 151M | 753M |
| 28 | 588M | 6 | 151M | 837M |
| 32 | 672M | 6 | 151M | 921M |
| 36 | 756M | 6 | 151M | 1005M |
| 40 | 840M | 6 | 151M | 1089M |
| 48 | 1007M | 6 | 151M | 1257M |
| 56 | 1175M | 6 | 151M | 1425M |
| 64 | 1343M | 6 | 151M | 1593M |

Table 2: Parameter decomposition of the decoder scaling models. The total number of parameters includes 98M parameters representing the softmax and embedding layers. Note that the 6L-6L model is the baseline model we used for hyper-parameter tuning.

| Encoder | | Decoder | | |
|---|---|---|---|---|
| Layers | Parameters | Layers | Parameters | Total Parameters |
| 6 | 126M | 2 | 50M | 275M |
| 6 | 126M | 4 | 101M | 325M |
| 6 | 126M | 6 | 151M | 375M |
| 6 | 126M | 8 | 202M | 426M |
| 6 | 126M | 12 | 302M | 527M |
| 6 | 126M | 16 | 403M | 627M |
| 6 | 126M | 20 | 504M | 728M |
| 6 | 126M | 24 | 605M | 829M |
| 6 | 126M | 28 | 705M | 930M |
| 6 | 126M | 32 | 806M | 1030M |
| 6 | 126M | 36 | 907M | 1131M |
| 6 | 126M | 40 | 1008M | 1232M |
| 6 | 126M | 48 | 1209M | 1433M |
| 6 | 126M | 56 | 1411M | 1635M |
| 6 | 126M | 64 | 1612M | 1836M |

Table 3: Parameter decomposition of the symmetric scaling models trained for English→German translation task. The total number of parameters includes 98M parameters representing the softmax and embedding layers.

| Encoder | | Decoder | | |
|---|---|---|---|---|
| Layers | Parameters | Layers | Parameters | Total Parameters |
| 2 | 42M | 2 | 50M | 191M |
| 3 | 63M | 3 | 76M | 237M |
| 4 | 84M | 4 | 101M | 283M |
| 5 | 105M | 5 | 126M | 329M |
| 8 | 168M | 8 | 202M | 468M |
| 16 | 336M | 16 | 403M | 837M |
| 20 | 420M | 20 | 504M | 1022M |
| 24 | 504M | 24 | 605M | 1207M |
| 32 | 672M | 32 | 806M | 1576M |
| 48 | 1007M | 48 | 1209M | 2315M |
| 56 | 1175M | 56 | 1411M | 2684M |
| 64 | 1343M | 64 | 1612M | 3054M |

# B   Scaling Laws for Other Test Sets

In order to keep the discussion in the main text focused, we only presented scaling laws for Web Domain 1 test sets. These test sets were chosen as they had a domain similar to the training data (i.e. web). In this appendix, we repeat the same analysis for our other test sets. The details of these test sets are described in Section 2.

Figure 12 demonstrates how well the scaling law in Eq. (2) fits the empirical scaling behavior of our models on all our test sets. For each test set, we have fitted the law jointly on the final test loss achieved by the encoder and decoder scaling models. We measure the final test loss by the median test loss over steps 450K to 500K. Details of the fitting procedure are provided in Appendix E.

The results suggest that Eq. (2) is closely capturing the scaling behavior of the model for all the test sets / domains. In the last row, we also demonstrate the fit for the training data (left column) and the training loss (cross entropy on training data plus regularization, right column). We observe that the scaling law is almost perfectly fitting the empirical data in these cases.

To examine the fit more closely, in Figure 13, we have plotted the same data but with several modifications:

1. Instead of plotting the final loss, we plot the <u>reducible</u> component of the final loss ($L - L_\infty$). As the true value of $L_\infty$ is unknown, we use the value given by the fit of the scaling law.

2. For encoder scaling models, we plot the (reducible) loss against the number of encoder parameters (as opposed to the total number of parameters). Similarly, for decoder scaling models, we plot the loss against the number of decoder parameters.

3. We use log-log scaling on the axes.

4. We use the results of Table 4 to provide a confidence region around our predictions. This confidence region quantifies our expected uncertainty caused by randomness in the initialization and training pipeline.

Eq. (2) predicts that the relationship between the empirical final loss values from the encoder (decoder) scaling models and the number of encoder (decoder) parameters should appear linear on these plots. Figure 13 suggests that the empirical scaling behavior of these models conforms closely to these predictions.

To see if Eq. (2) continues to capture the scaling behavior of the models out-of-sample, we compare the predictions of the scaling law with the empirical loss values for symmetric scaling models. In other words, we examine how well the scaling laws fitted only using encoder / decoder scaling models predict the final test loss achieved by symmetric scaling models of different sizes. Figure 14 shows this comparison for all of our test sets. We observe a remarkable match between the predictions of the scaling law and the empirical loss values across the board. These observations confirm that Eq. (2) is able to capture the scaling behavior of the model regardless of the scaling approach.
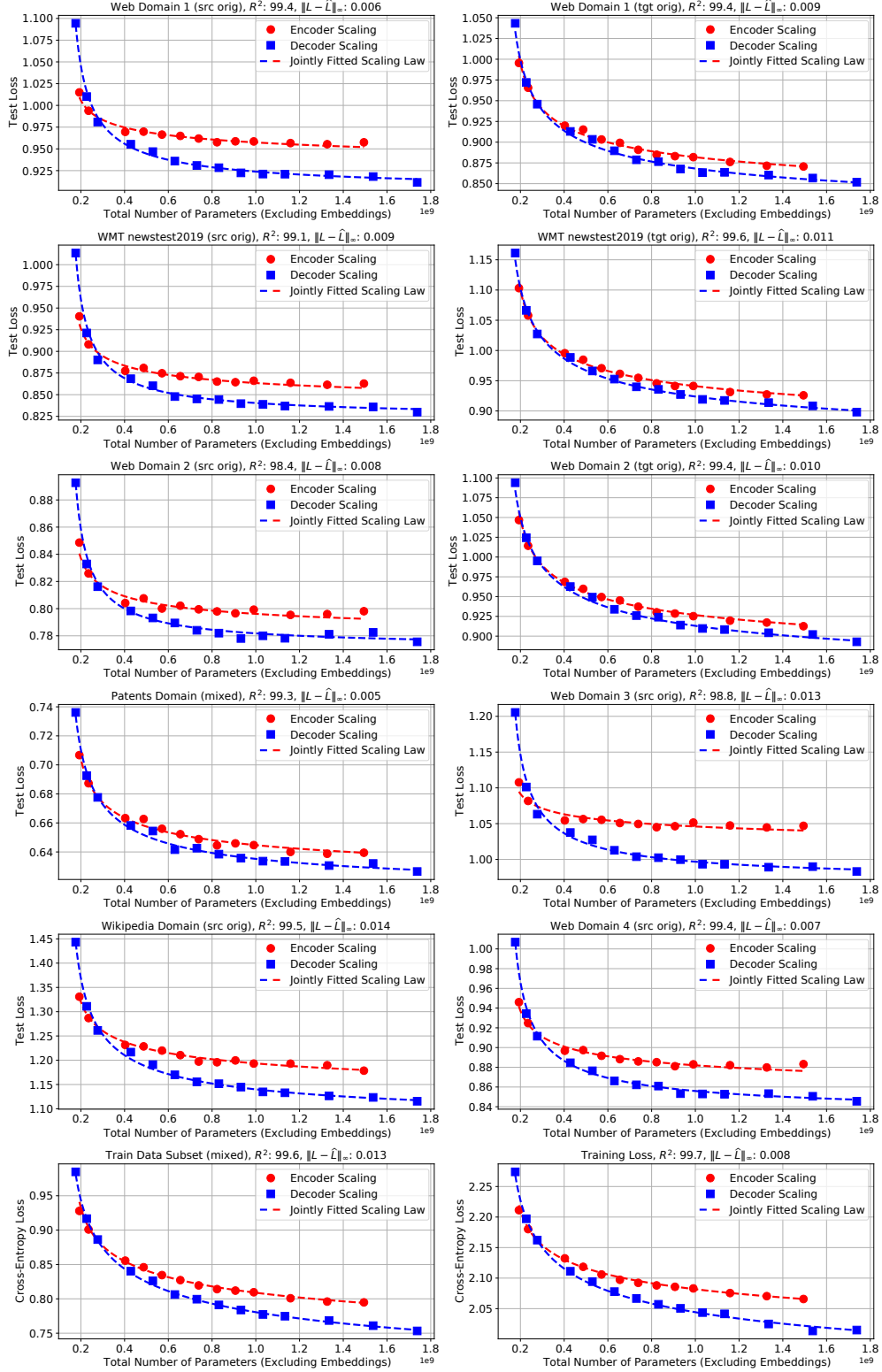
Figure 12: Fitted scaling law (Eq. (2)) for English→German translation task. The scaling law captures the scaling behavior of the models over a diverse collection of test sets and domains. The last row describes the evolution of the cross-entropy loss on the training data (with and without regularization effect).
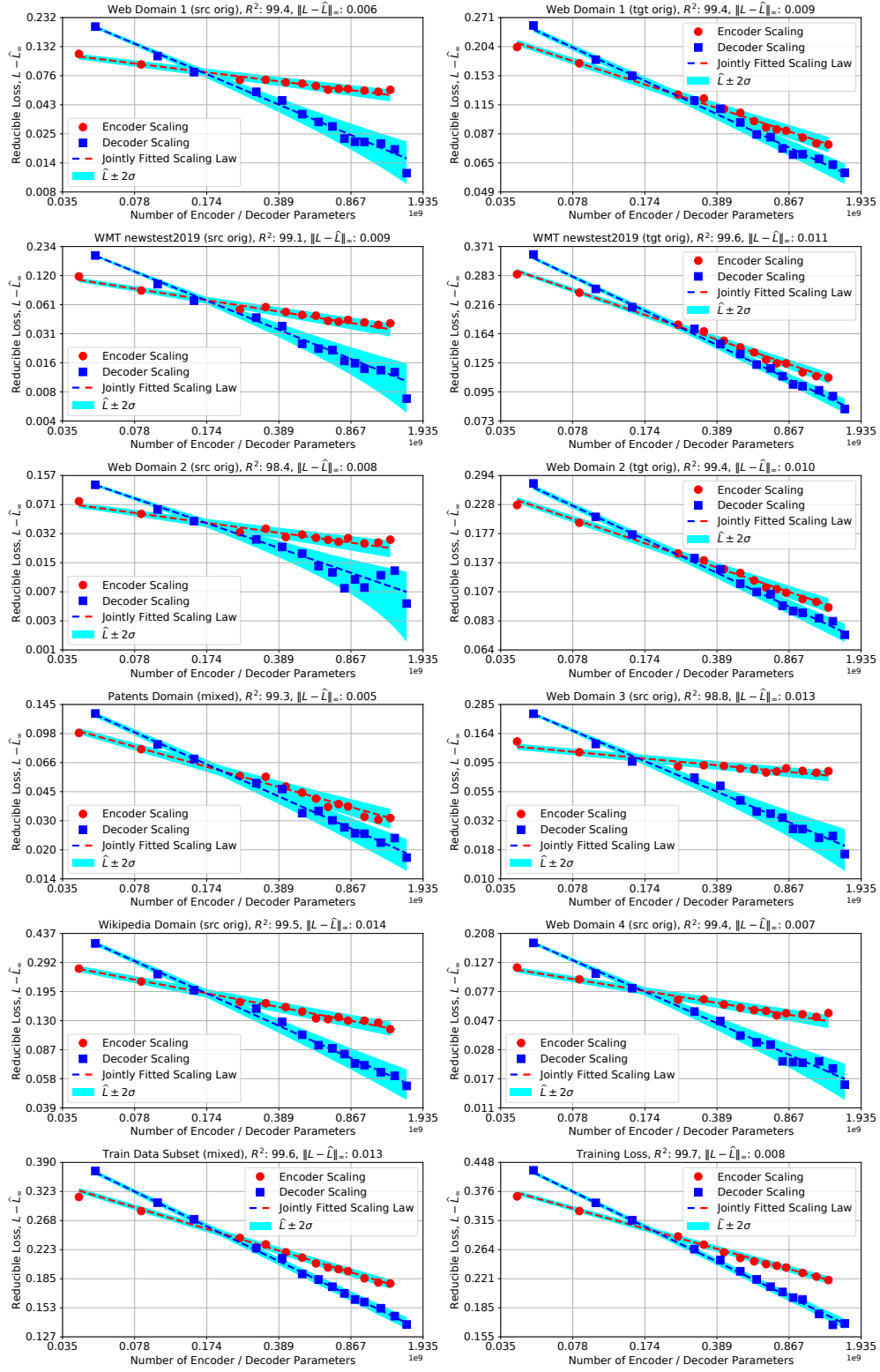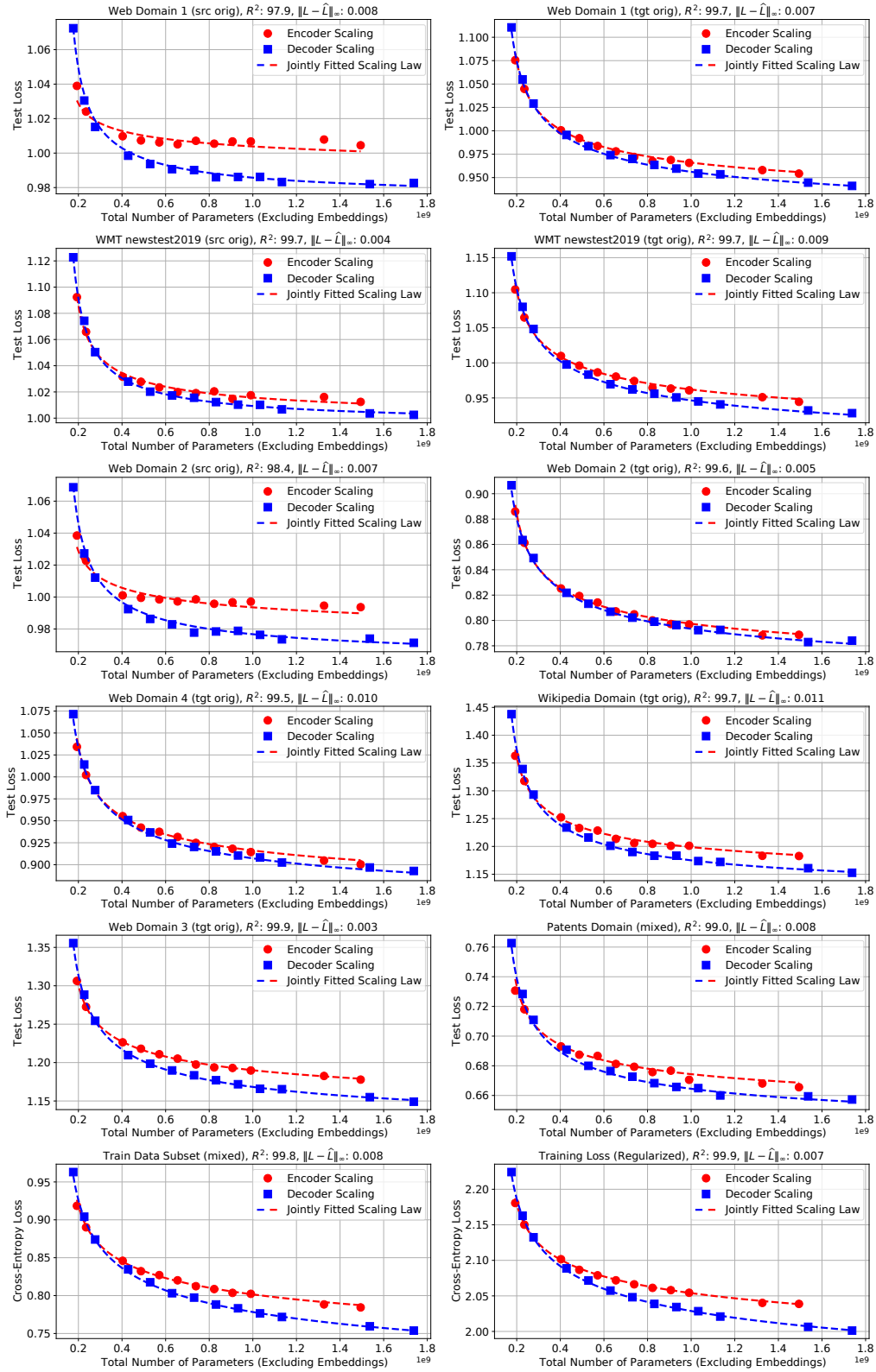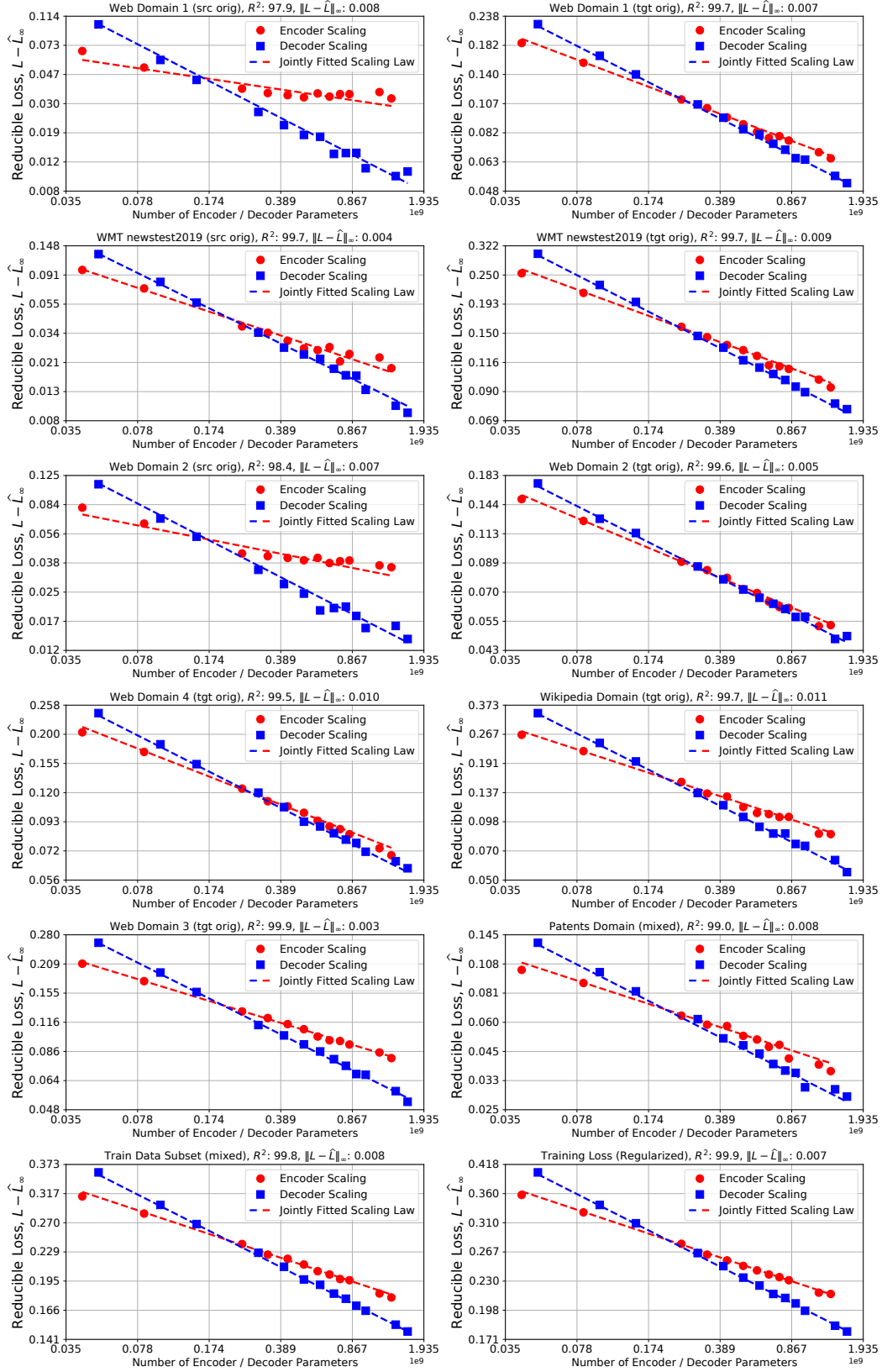
18

Figure 13: Fitted scaling law (Eq. (2)) for English→German translation task. Here, we use a log-log plot in order to inspect the fit more closely. Shaded cyan regions correspond to the uncertainty region given by $\pm 2 \times$standard deviation. Per test set standard deviations are provided in Table 4.

Figure 14: Out-of-sample prediction accuracy of English→German scaling laws on symmetric scaling models. Scaling laws are fitted only using the encoder and decoder scaling models. Nevertheless, they accurately predict the scaling behavior of symmetric scaling models.

Figure 15: Fitted scaling law (Eq. (2)) for German→English translation task. The scaling law captures the scaling behavior of the models over a diverse collection of test sets and domains.

Figure 16: German→English scaling law fits on log-log scale.

# C    Quantifying the Random Variations in the Results

Note that the final test loss achieved by the model is a random quantity. Randomness is incorporated into the training pipeline through the initialization step, data order, and hardware failures / preemptions. To quantify the magnitude of the fluctuations caused by this randomness, we retrain a subset of our models (2L-2L, 3L-3L, 4L-4L, 5L-5L, 8L-8L, 16L-16L, 24L-24L, and 32L-32L) with 4 different seeds. Figure 17 presents standard deviation (left) and maximum difference (right) of the final test loss values observed for each model.



Figure 17: Variability of the final test loss across four different seeds.

Table 4: Variability of final test loss for each test dataset (averaged over all models).

| Dataset | Average Standard Deviation | Average Maximum Deviation |
|---|---|---|
| Web Domain 1 (src orig) | 0.0030 | 0.0078 |
| Web Domain 1 (tgt orig) | 0.0030 | 0.0077 |
| WMT newstest2019 (src orig) | 0.0027 | 0.0067 |
| WMT newstest2019 (tgt orig) | 0.0028 | 0.0071 |
| Web Domain 2 (src orig) | 0.0024 | 0.0062 |
| Web Domain 2 (tgt orig) | 0.0030 | 0.0076 |
| Patents Domain (mixed) | 0.0021 | 0.0053 |
| Web Domain 3 (src orig) | 0.0037 | 0.0094 |
| Wikipedia Domain (src orig) | 0.0057 | 0.0146 |
| Web Domain 4 (src orig) | 0.0026 | 0.0066 |
| Train Data Subset (mixed) | 0.0026 | 0.0067 |
| Training Loss | 0.0025 | 0.0064 |

# D Proofs

## D.1 Proof of proposition 1

*Proof.* Let $\beta \equiv \alpha \bar{N}_e^{p_e} \bar{N}_d^{p_d}$. Then the optimal encoder / decoder sizes are optimal parameters of the following optimization problem:

$$\begin{aligned} \text{minimize}_{N_e, N_d} \quad & \beta N_e^{-p_e} N_d^{-p_d} \\ \text{s.t.} \quad & N_e + N_d \leq B \end{aligned} \quad . \tag{6}$$

To convert the problem to a convex problem, we instead consider the log of the objective and adopt the following change of variables:

$$u \equiv \log(N_e), \qquad v \equiv \log(N_d). \tag{7}$$

The transformed optimization problem is of the form:

$$\begin{aligned} \text{minimize}_{u,v} \quad & -p_e u - p_d v \\ \text{s.t.} \quad & \exp(v) + \exp(u) \leq B \end{aligned} \quad . \tag{8}$$

Note that (8) is now convex and therefore, we can use KKT conditions to solve for the optimum. The Lagrangian has the form:

$$\mathcal{L}(u, v, \lambda) = -p_e u - p_d v + \lambda \left( B - \exp(v) - \exp(u) \right). \tag{9}$$

Solving for the first-order conditions yield:

$$-p_e = \lambda \exp(u^*) \tag{10}$$
$$-p_d = \lambda \exp(v^*). \tag{11}$$

Since the constraint is binding, $\lambda \neq 0$. Therefore, we can divide both sides of the equations above which yields:

$$\frac{p_e}{p_d} = \frac{\exp(u^*)}{\exp(v^*)} = \frac{N_e^*}{N_d^*}. \tag{12}$$

Substituting (12) in the constraint yields:

$$N_e^* = \frac{p_e}{p_e + p_d} B, \qquad N_d^* = \frac{p_d}{p_e + p_d} B. \tag{13}$$

Finally, we substitute (13) in the scaling law which yields:

$$\hat{L}_{opt}(B) = \alpha \left( \frac{\bar{N}_e (p_e + p_d)}{p_e B} \right)^{p_e} \left( \frac{\bar{N}_d (p_e + p_d)}{p_d B} \right)^{p_d} + L_\infty \tag{14}$$

$$= \alpha \left( \frac{\bar{N}_e (p_e + p_d)}{p_e} \right)^{p_e} \left( \frac{\bar{N}_d (p_e + p_d)}{p_d} \right)^{p_d} B^{-(p_e + p_d)} + L_\infty \tag{15}$$

$\square$

# E    Curve Fitting Details

We use `scipy.optimize.least_squares` function for curve fitting throughout this paper [10]. To have some robustness to outliers, we use the `loss='soft_l1'` option which is a popular option for robust regression. The code snippet below shows the exact arguments we use for fitting the scaling laws:

```python
def func(p, x, y):
  """Fitting a bivariate scaling law.

  p: A 1-D array of dim 4, corresponding to alpha, p_e, p_d, c.
  x: A matrix of dimension n \times 2. First column encoder params,
    second col decoder params.
  y: A 1-D array of log-pplx of dim n."""
  x_e = NE_bar / x[:, 0]
  x_d = ND_bar / x[:, 1]
  return p[0] * np.power(x_e , p[1]) * np.power(x_d , p[2]) + p[3] - y

def fit_model(x, y, f_scale):
  X = x.to_numpy().copy()
  y = y.to_numpy().copy()
  if np.isnan(X).any() or np.isnan(y).any():
    raise ValueError('Data contains NaNs')
  if len(y.shape) > 1 or y.shape[0] != X.shape[0]:
    raise ValueError('Error in shapes')

  p0 = np.zeros((4,))
  p0[0] = 0.2 # alpha
  p0[1] = 0.4 # p_e
  p0[2] = 0.6 # p_d
  p0[3] = 1.0 # c
  fit = least_squares(func, p0, loss='soft_l1', f_scale=f_scale,
                      args=(X, y), max_nfev=10000, bounds=(0, 10))
  return fit
```

The `'soft_l1'` loss chosen above applies $\ell_2$ penalty on small residuals and a $\ell_1$-like penalty on outlier residuals. The argument `f_scale` determines the boundary where the transition between the two different behaviors occur. For the results presented in this paper, we choose `f_scale` from the grid given by `np.geomspace(0.001, 0.025, num=25)`. Choosing `f_scale=0.025` effectively yields a least-squares regression while smaller values add more robustness to outliers.

---

[10] `https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.least_squares.html`

# F  Analysis of the Generation Quality

**Decoding:**  As described in Section 4, we use beam-search for decoding [38]. To keep the experiments tractable, we did not attempt to tune the hyper-parameters of beam-search for each model. Instead, we use the same hyper-parameters (optimized for the baseline model) for all our decoding jobs. In particular, we fix the length normalization parameter to $1.0$ and number of beams to $4$.

**BLEU-Cross Entropy Loss Co-Evolution:**  Figure 18 presents the relationship between BLEU score and cross-entropy loss for various test datasets. The results closely mimic the phenomenon observed in Figure 10: On target original data, improvements to cross-entropy loss are accompanied with improvements in BLEU score. On source original data however, beyond a certain point, cross-entropy loss and BLEU score exhibit diverging behaviors.

We observe that in large well-trained models, the relationship between BLEU and cross-entropy loss on target-original data is well captured by the power law presented in Eq. (5). The fit achieved by this power law is plotted in our figures. We observe that fitted power laws for encoder scaling models consistently attain larger exponents compared to decoder or symmetrically scaled models. This reflects the fact that encoder scaling models are more successful in improving the generation quality (as measured by BLEU).

Finally, we observe a number of deviations from the predictions of Eq. (5). In particular, models with shallow decoders (6L2L, 6L4L, 6L6L) seem to outperform the trend (Figure 19). Moreover, we observe that models in the beginning of the training process tend to deviate from the overall trend (Figure 20). We postpone an in-depth analysis of these phenomena to future work.
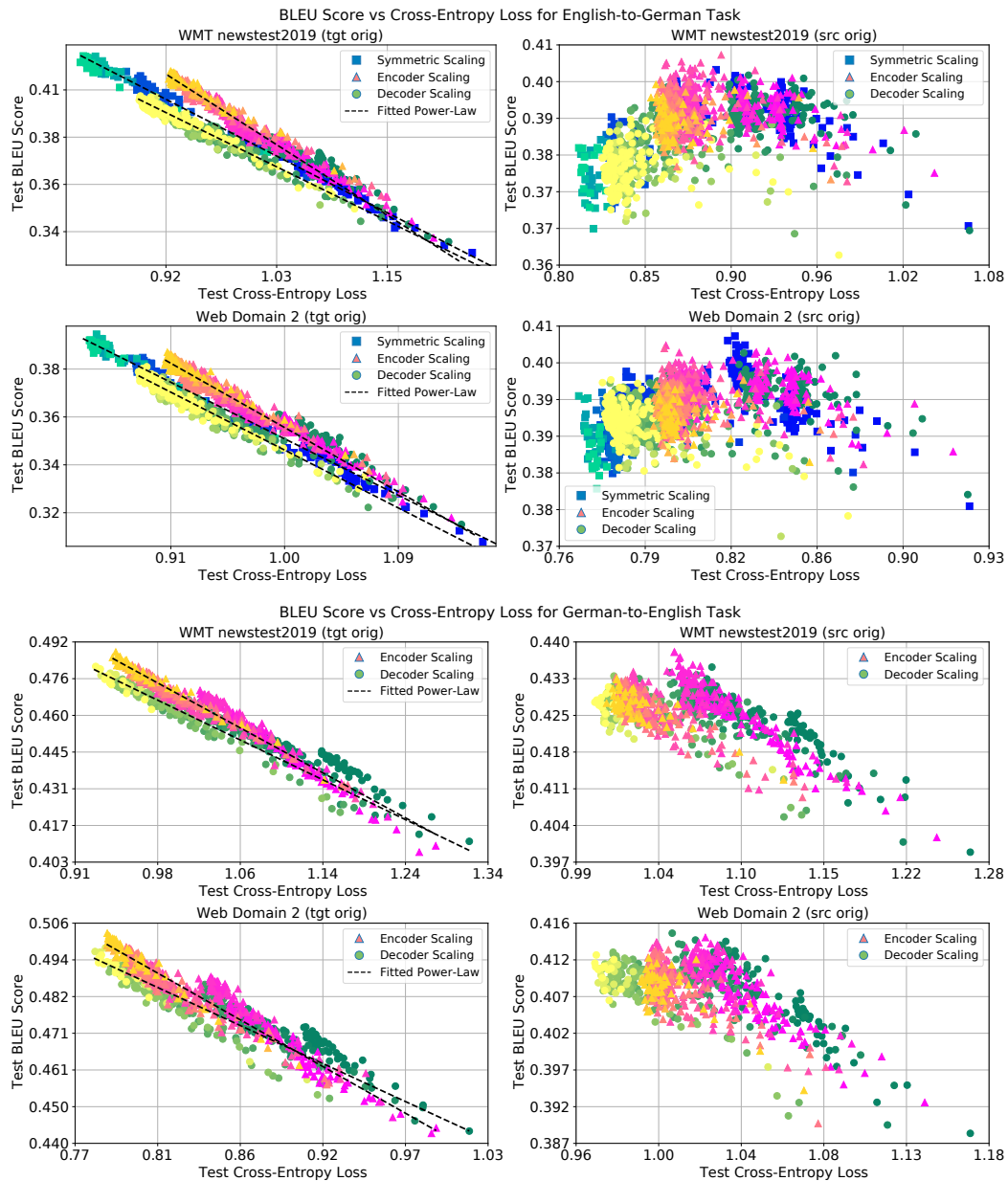
Figure 18: Log-log plot of the evolution of BLEU score as a function of cross-entropy loss for different models. For each scaling approach, warmer colors represent larger models. Each individual color represents different checkpoints of a single model during training. On target original data (left column), improvements to cross-entropy loss lead to consistent improvements in BLEU score. Dashed lines correspond to fit achieved by Eq. (5). The relationship breaks down for source original data (right column).

Figure 19: Models with shallow decoders tend to outperform predictions of Eq. (5). Points with dark green color represent different checkpoints of 6L2L, 6L4L, and 6L6L models.
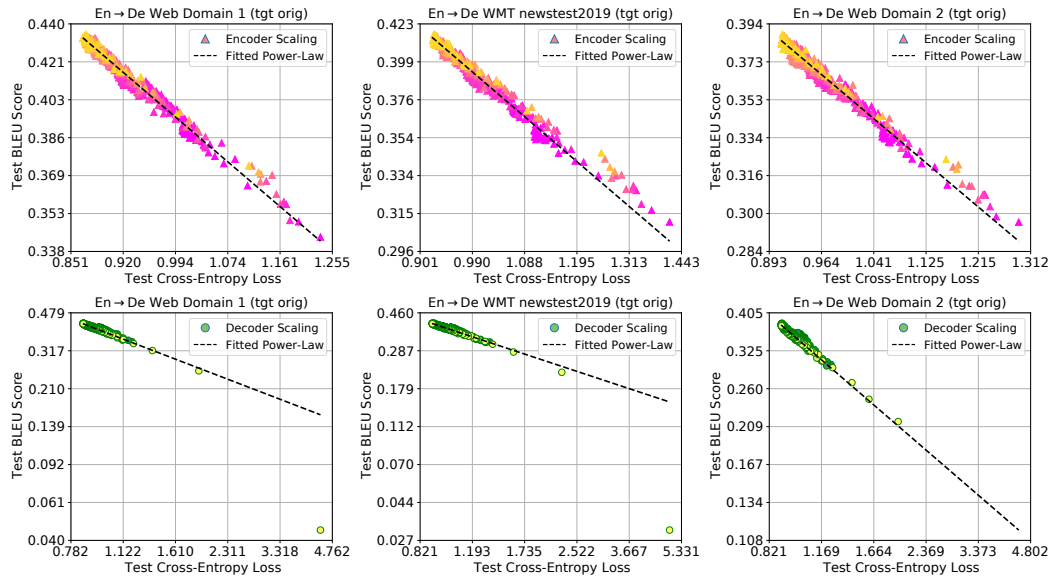


Figure 20: On some of the test sets, data points corresponding to early training checkpoints exhibit deviations from the overall trend.
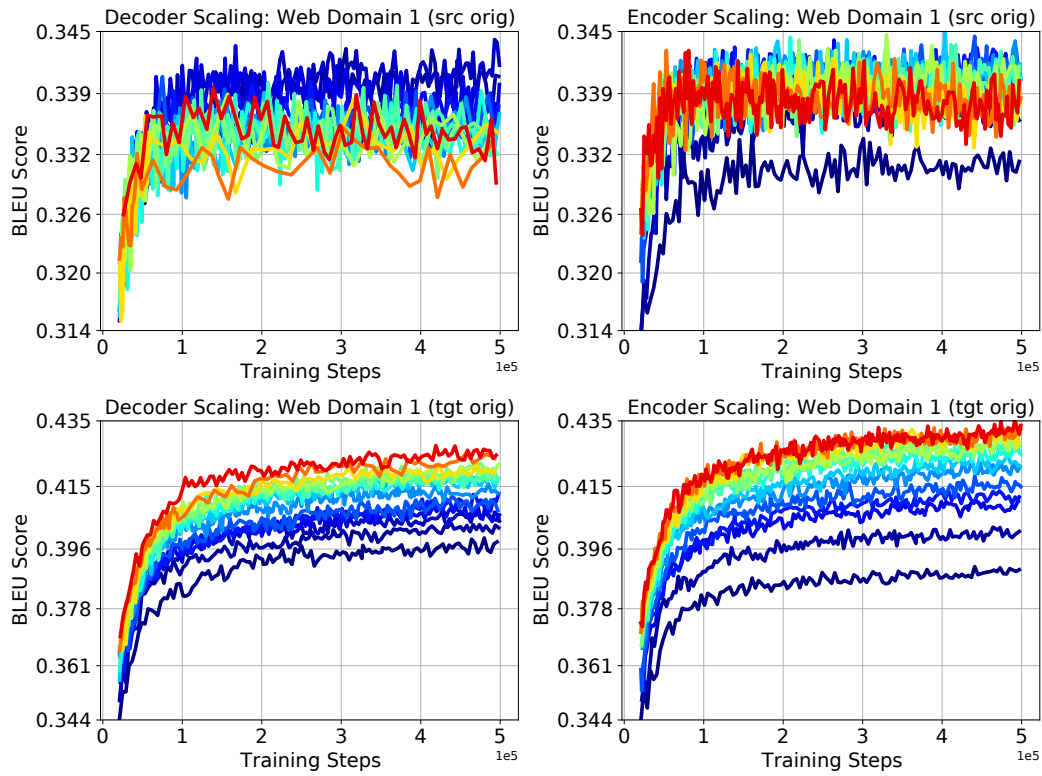
Figure 21: The evolution of BLEU score during the training for English-to-German Web Domain test sets. Warmer colors correspond to larger models. Top row: On source original test data, our largest models achieve lower BLEU scores compared to mid-sized models throughout the training. Bottom row: On target original test data, increasing the model size yields consistent improvements in BLEU score throughout the training.
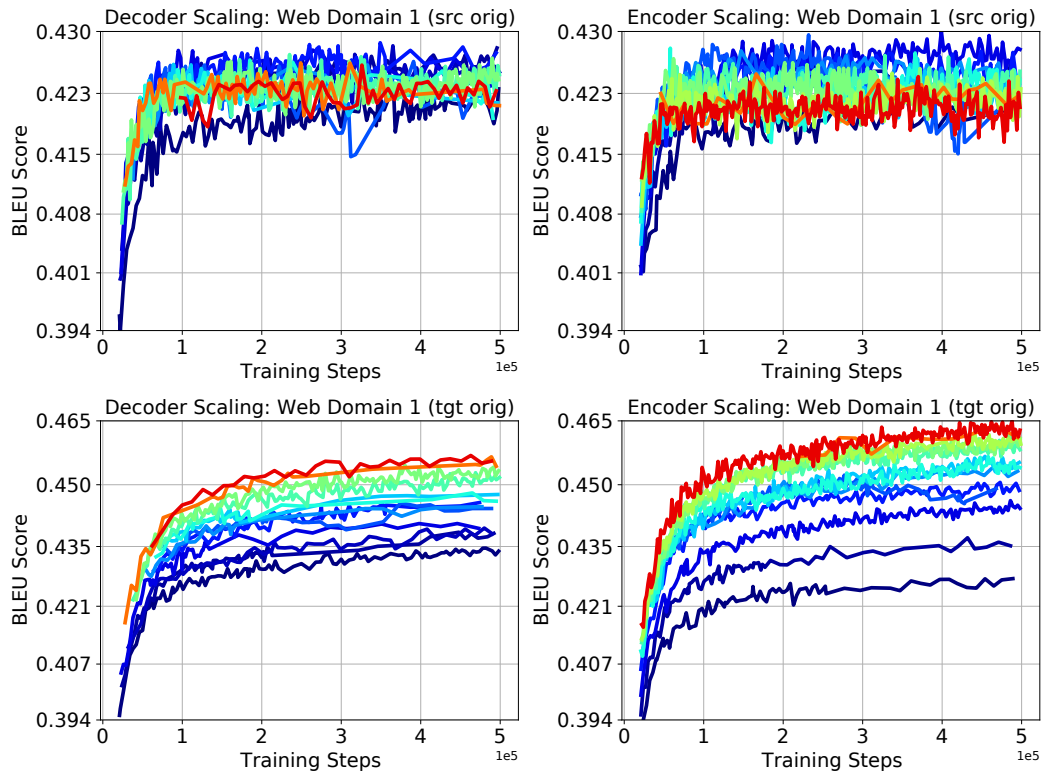
Figure 22: The evolution of BLEU score during the training for German-to-English Web Domain test sets. Warmer colors correspond to larger models. Top row: On source original test data, our largest models achieve lower BLEU scores compared to mid-sized models throughout the training. Bottom row: On target original test data, increasing the model size yields consistent improvements in BLEU score throughout the training.
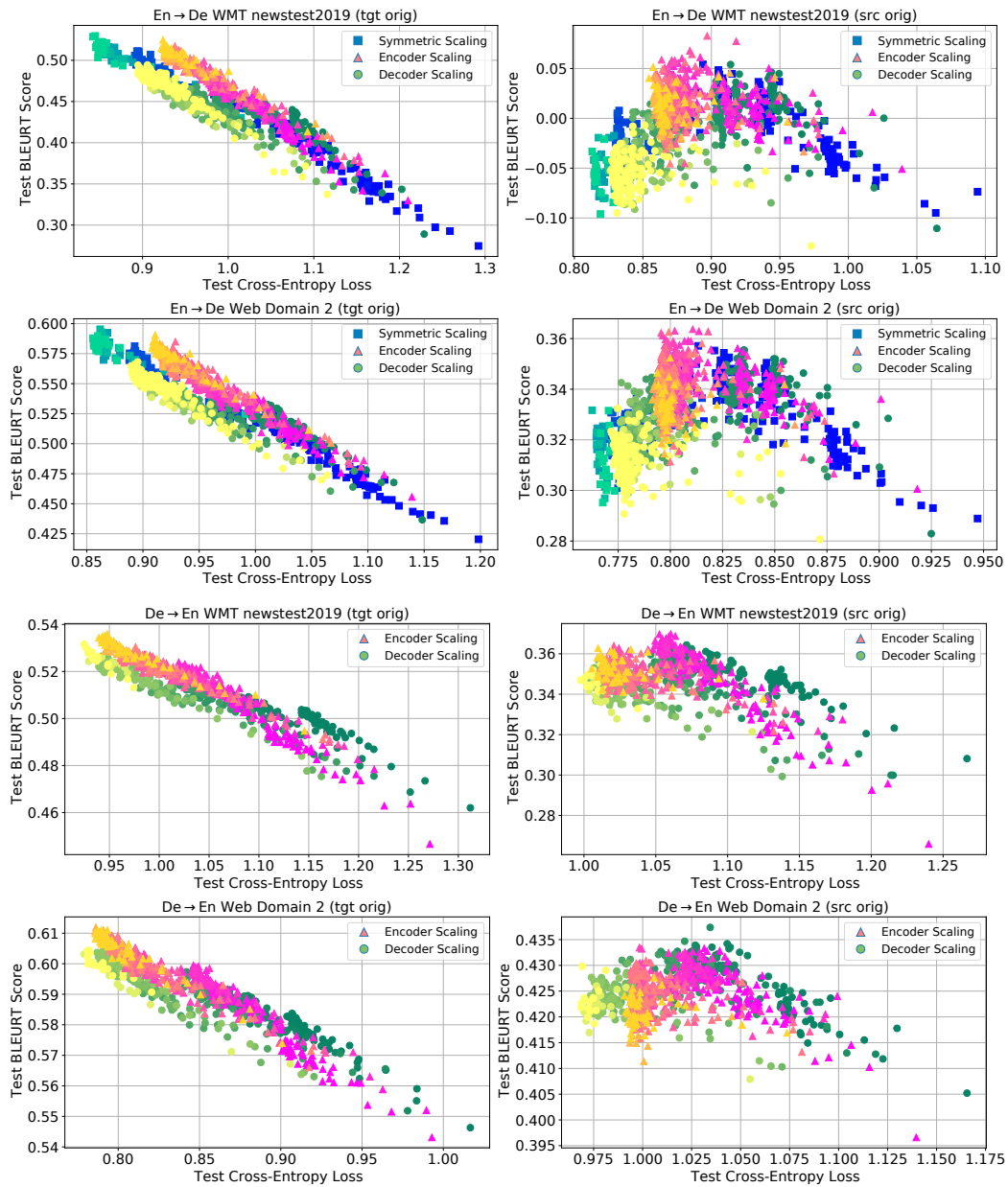
Figure 23: The evolution of BLEURT score as a function of cross-entropy loss for different models. For each scaling approach, warmer colors represent larger models. Each individual color represents different checkpoints of a single model during training. On target original data (left column), improvements to cross-entropy loss lead to consistent improvements in BLEURT score. This relationship breaks down for source original data (right column).