

DEPTH-ADAPTIVE TRANSFORMER

Anonymous authors

Paper under double-blind review

ABSTRACT

State of the art sequence-to-sequence models perform a fixed number of computations for each input sequence regardless of whether it is easy or hard to process. In this paper, we train Transformer models which can make output predictions at different stages of the network and we investigate different ways to predict how much computation is required for a particular sequence. Unlike dynamic computation in Universal Transformers, which applies the same set of layers iteratively, we apply different layers at every step to adjust both the amount of computation as well as the model capacity. Experiments on machine translation benchmarks show that this approach can match the accuracy of a baseline Transformer while using only half the number of decoder layers.

1 INTRODUCTION

Neural sequence models achieve state-of-the-art results in many natural language processing tasks including machine translation (Gehring et al., 2017; Vaswani et al., 2017), language modeling (Baevski & Auli, 2019; Radford et al., 2019), summarization (Celikyilmaz et al., 2018; Wu et al., 2019) and natural language understanding (Radford et al., 2018; Devlin et al., 2019).

These models can count hundreds or even billions of parameters (Radford et al., 2019). For example, the winning entry of the WMT'19 machine translation task in English-German used an ensemble totaling two billion parameters (Ng et al., 2019). While large models are required to do better on hard examples, small models are likely to perform as well on easy examples, *e.g.*, the aforementioned ensemble is probably not required to translate a short phrase such as "Thank you". However, current models apply the same amount of computation regardless of whether the input is easy or hard.

In this paper, we propose Transformers which adapt the number of layers to each input in order to achieve a good speed-accuracy trade off. We extend Graves (2016; ACT) who introduced dynamic computation to recurrent neural networks in several ways: we apply different layers at each stage, we investigate a range of designs and training targets for the halting module and we explicitly supervise through simple oracles to achieve good performance on large-scale tasks.

Universal Transformers (UT) use dynamic computation via ACT but Dehghani et al. (2018) found that fixing the number of computation steps performs better for large-scale machine translation while as our approach demonstrates substantial improvements in speed at no loss in accuracy. UT layers have the size of an entire standard Transformer and they are invoked several times which has a large impact on speed. Our approach does not increase the size of individual layers. We also extend the resource efficient object classification work of Huang et al. (2017) to structured prediction where dynamic computation decisions impact future computation.

We encode the input sequence using a standard Transformer encoder to generate the output sequence with a varying amount of computation in the decoder network. Dynamic computation poses a challenge for self-attention because omitted layers in prior time-steps may be required in the future. We experiment with two approaches to address this and show that a simple approach works well (§2).

Next, we investigate different mechanisms to control the amount of computation in the decoder network, either for the entire sequence or on a per-token basis. This includes multinomial and binomial classifiers supervised by the model likelihood or whether the argmax is already correct as well as simply thresholding the model score (§3). Experiments on IWSLT14 German-English translation (Cettolo et al., 2014) as well as WMT'14 English-French translation show that we can match the baseline performance at substantially less computation (§4).

2 ANYTIME STRUCTURED PREDICTION

We first present a model that can make predictions at different layers. This is known as *anytime prediction* for computer vision models and we extend it to structured prediction (Huang et al., 2017).

2.1 TRANSFORMER WITH MULTIPLE OUTPUT CLASSIFIERS

We base our approach on the Transformer sequence-to-sequence model (Vaswani et al., 2017). Both encoder and decoder networks contain N stacked blocks where each has several sub-blocks surrounded by residual skip-connections. The first sub-block is a multi-head dot-product self-attention and the second a position-wise fully connected feed-forward network. For the decoder, there is an additional sub-block after the self-attention to add source context via another multi-head attention.

Given a pair of source-target sequences (\mathbf{x}, \mathbf{y}) , \mathbf{x} is processed with the encoder to give representations $\mathbf{s} = (s_1, \dots, s_{|\mathbf{x}|})$. Next, the decoder generates \mathbf{y} step-by-step. For every new token \mathbf{y}_t input to the decoder at time t , the N decoder blocks process it to yield hidden states $(h_t^n)_{1 \leq n \leq N}$:

$$h_t^0 = \text{embed}(\mathbf{y}_t), \quad h_t^n = \text{block}_n(h_{<t}^{n-1}, \mathbf{s}),$$

where block_n is the mapping associated with the n^{th} block and embed is a lookup table.

The output distribution for predicting the next token is computed by feeding the activations of the last decoder layer h_t^N into a softmax normalized output classifier W :

$$p(\mathbf{y}_{t+1} | h_t^N) = \text{softmax}(W h_t^N)$$

Standard Transformers have a single output classifier attached to the top of the decoder network. However, for dynamic computation we need to be able to make predictions at different stages of the network. To achieve this, we attach output classifiers \mathcal{C}_n parameterized by W_n to the output h_t^n of each of the N decoder blocks:

$$\forall n, p(\mathbf{y}_{t+1} | h_t^n) = \text{softmax}(W_n h_t^n)$$

The classifiers can be parameterized independently or we can share the weights across the N blocks.

2.2 TRAINING MULTIPLE OUTPUT CLASSIFIERS

Dynamic computation enables the model to use any of the N exit classifiers instead of just the final one. Some of our models can choose a different output classifier at each time-step which results in an exponential number of possible output classifier combinations in the sequence length.

We consider two possible ways to train the decoder network (Figure 1). *Aligned training* optimizes all classifiers simultaneously and assumes all previous hidden states required by the self-attention are available. However, at test time this is often not the case when we choose a different exit for every token which leads to misaligned states. Instead, *mixed training* samples several sequences of exits for a given sentence and exposes the model to hidden states from different layers.

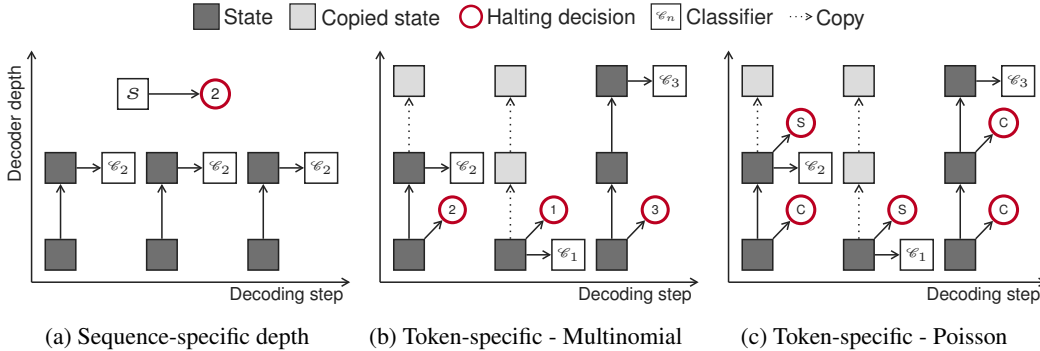
Generally, for a given output sequence \mathbf{y} , we have a sequence of chosen exits $(n_1, \dots, n_{|\mathbf{y}|})$ and we denote the block at which we exit at time t as n_t .

2.2.1 ALIGNED TRAINING

Aligned training assumes all hidden states $h_1^{n-1}, \dots, h_t^{n-1}$ are available in order to compute self-attention and it optimizes N loss terms, one for each exit (Figure 1a):

$$\text{LL}_t^n = \log p(\mathbf{y}_t | h_{t-1}^n), \quad \text{LL}^n = \sum_{t=1}^{|\mathbf{y}|} \text{LL}_t^n, \quad \mathcal{L}_{dec}(\mathbf{x}, \mathbf{y}) = - \frac{1}{\sum_n \omega_n} \sum_{n=1}^N \omega_n \text{LL}^n. \quad (1)$$

The compound loss $\mathcal{L}_{dec}(\mathbf{x}, \mathbf{y})$ is a weighted average of N terms w.r.t. to $(\omega_1, \dots, \omega_N)$. We found that uniform weights achieve better BLEU compared to other weighing schemes (c.f. Appendix A). At inference time, not all time-steps will have hidden states for the current layer since the model exited early. In this case, we simply *copy* the last computed state to all upper layers, similar to mixed training (§2.2.2). However, we do apply layer-specific key and value projections to the copied state.



Multinomial q_t :

$$q_t(n|\mathbf{x}, \mathbf{y}_{<t}) = \text{softmax}(W_h h_t^1 + b_h), \quad (8)$$

Poisson binomial q_t :

$$\forall n \in [1..N-1], \chi_{t,n} = \sigma(W_h h_t^n + b_h), \quad q_t(n|\mathbf{x}, \mathbf{y}_{<t}) = \begin{cases} \chi_t^n \prod_{n' < n} (1 - \chi_t^{n'}), & \text{if } n < N \\ \prod_n \chi_t^n, & \text{otherwise} \end{cases} \quad (9)$$

where W_h and b_h are the weights and biases of the halting mechanism. The two classifiers are trained to minimize the cross-entropy with respect to either one the following oracle distributions:

Likelihood-based: At each time-step t , we choose the block whose exit classifier has the highest likelihood plus a regularization term weighted by λ to encourage lower exits.

$$q_t^*(\mathbf{x}, \mathbf{y}) = \delta(\arg \max_n \text{LL}_t^n - \lambda n). \quad (10)$$

This oracle ignores the impact of the current decision on the future time-steps and we therefore consider a smoothed likelihood which includes information from surrounding tokens:

$$\begin{aligned} \kappa(t, t') &= e^{-\frac{|t-t'|^2}{\sigma}}, \quad \text{smoothLL}_t^n = \sum_{t'} \kappa(t, t') \text{LL}_{t'}^n, \\ q_t^*(\mathbf{x}, \mathbf{y}) &= \delta(\arg \max_n \text{smoothLL}_t^n - \lambda n), \end{aligned} \quad (11)$$

where we control the size of the surrounding context with σ . We refer to this oracle as $\text{LL}(\sigma, \lambda)$ including the case where we only look at the likelihood of the current token with $\sigma = 0$.

Correctness-based: Similar to the likelihood-based oracle we can look at the correctness of the prediction at time-step t as well as surrounding positions. We define the target q_t^* as follows:

$$\begin{aligned} C_t^n &= y_t = \arg \max_y p(y|h_{t-1}^n), \quad \text{smoothC}_t^n = \sum_{t'} \kappa(t, t') C_{t'}^n, \\ q_t^*(\mathbf{x}, \mathbf{y}) &= \delta(\arg \max_n \text{smoothC}_t^n - \lambda n). \end{aligned} \quad (12)$$

Confidence thresholding Finally, we consider thresholding the model predictions (§2), i.e., exit when the argmax score of the current output classifier $p(y_{t+1}|h_t^n)$ exceeds a threshold θ_n . This does not require training and the thresholds θ are simply tuned on the valid set to maximize BLEU.

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

We evaluate on several benchmarks and measure tokenized BLEU (Papineni et al., 2002):

IWSLT’14 German to English (De-En). We use the setup of Edunov et al. (2018) and train on 160K sentence pairs. We use $N = 6$ blocks, FFN dimension 1024, 4 heads, dropout 0.3, embedding dimension $d_{\text{enc}} = 512$ for the encoder and $d_{\text{dec}} = 256$ for the decoder. Embeddings are untied with 6 different output classifiers. We evaluate with a single checkpoint and a beam of width 5.

WMT’14 English to French (En-Fr). We also experiment on the much larger WMT’14 English-French task comprising 35.5m training sentence pairs. We develop on 26k held out pairs and test on newstest14. The vocabulary consists of 44k joint BPE types (Sennrich et al., 2016). We use a Transformer *big* architecture and tie the embeddings of the encoder, the decoder and the output classifiers ($(W_n)_{1 \leq n \leq 6}$; §2.1). We average the last ten checkpoints and use a beam of width 4.

Models are implemented in fairseq (Ott et al., 2019) and are trained with Adam (Kingma & Ba, 2015). We train for 50k updates on 128 GPUs with a batch size of 460k tokens for WMT’14 En-Fr and on 2 GPUs with 8k tokens per batch for IWSLT’14 De-En. To stabilize training, we re-normalize the gradients if the norm exceeds $g_{\text{clip}} = 3$.

	Uniform	$n = 1$	$n = 2$	$n = 3$	$n = 4$	$n = 5$	$n = 6$	Average
Baseline	-	34.5	35.5	35.8	35.7	35.8	36.0	35.5
Aligned ($\omega_n = 1$)	35.5	34.1	35.5	35.8	36.1	36.1	36.2	35.6
Mixed $M = 1$	34.1	32.9	34.3	34.5	34.5	34.6	34.5	34.2
Mixed $M = 3$	35.1	33.9	35.2	35.4	35.5	35.5	35.5	35.2
Mixed $M = 6$	35.3	34.2	35.4	35.8	35.9	35.8	35.9	35.5

Table 1: Aligned vs. mixed training on IWSLT De-En. We report valid BLEU for a uniformly sampled exit $n \sim \mathcal{U}([1..6])$ at each token, a fixed exit $n \in [1..6]$ for all tokens, as well as the average BLEU over the fixed exits. As baseline we show six standard Transformer models with 1-6 blocks.

For models with adaptive exits, we first train without exit prediction ($\alpha = 0$ in Eq. (4)) using the aligned mode (*c.f.* §2.2.1) for 50k updates and then continue training with $\alpha \neq 0$ until convergence. For WMT’14 En-Fr we freeze the encoder and decoder parameters in the second phase and only optimize the halting mechanism. The exit prediction classifiers are parameterized by a single linear layer (Eq. (5)) with the same input dimension as the embedding dimension, *e.g.*, 1024 for a big Transformer; the output dimension is N for a multinomial classifier or one for Poisson. We exit when $\chi_{t,n} > 0.5$ for Poisson classifiers.

4.2 TRAINING MULTIPLE OUTPUT CLASSIFIERS

We first compare the two training regimes for our model (§2.2). Aligned training performs self-attention on aligned states (§2.2.1) and mixed training exposes self-attention to hidden states from different blocks (§2.2.2).

We compare the two training modes when choosing either a uniformly sampled exit or a fixed exit $n = 1, \dots, 6$ at inference time for every time-step. The sampled exit experiment tests the robustness to mixed hidden states and the fixed exit setup simulates an ideal setting where all previous states are available. As baselines we show six separate standard Transformers with $N \in [1..6]$ decoder blocks. All models are trained with an equal number of updates and mixed training with $M=6$ paths is most comparable to aligned training since the number of losses per sample is identical.

Table 1 shows that aligned training outperforms mixed training both for fixed exits as well as for randomly sampled exits. The latter is surprising since aligned training never exposes the self-attention mechanism to hidden states from other blocks. We suspect that this is due to the residual connections which *copy* features from lower blocks to subsequent layers and which are ubiquitous in Transformer models (§2). Aligned training also performs very competitively to the individual baseline models.

Aligned training is conceptually simple and fast. We can process a training example with N exits in a single forward/backward pass while M passes are needed for mixed training. In the remaining paper, we use the aligned mode to train our models. Appendix A reports experiments with weighing the various output classifiers differently but we found that a uniform weighting scheme worked well.

4.3 ADAPTIVE DEPTH ESTIMATION

Next, we train models with aligned states and compare adaptive depth classifiers in terms of BLEU as well as computational effort. We measure the latter as the average exit per output token (AE).

As baselines we use again six separate standard Transformers with $N \in [1..6]$ with a single output classifier. We also measure the performance of the aligned mode trained model for fixed exits $n \in [1..6]$. For the adaptive depth token-specific models (Tok), we train four combinations: likelihood-based oracle (LL) + Poisson, likelihood-based oracle (LL) + multinomial, correctness based oracle (C) + Poisson and correctness-based oracle (C) + multinomial. Sequence-specific models (Seq) are trained with the correctness oracle (C) and the likelihood oracle (LL) with different values for the regularization weight λ . All parameters are tuned on the valid set and we report results on the test set for a range of average exits.

Figure 3 shows that the aligned model (blue line) can match the accuracy of a standard 6-block Transformer (black line) at half the number of layers ($n = 3$) by always exiting at the third block. On

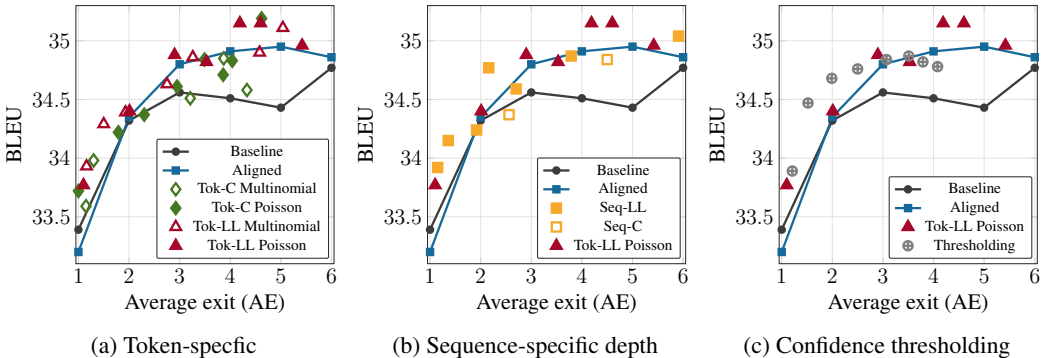


Figure 3: Trade-off between speed (average exit or AE) and accuracy (BLEU) for depth-adaptive methods on the IWSLT14 De-En test set.

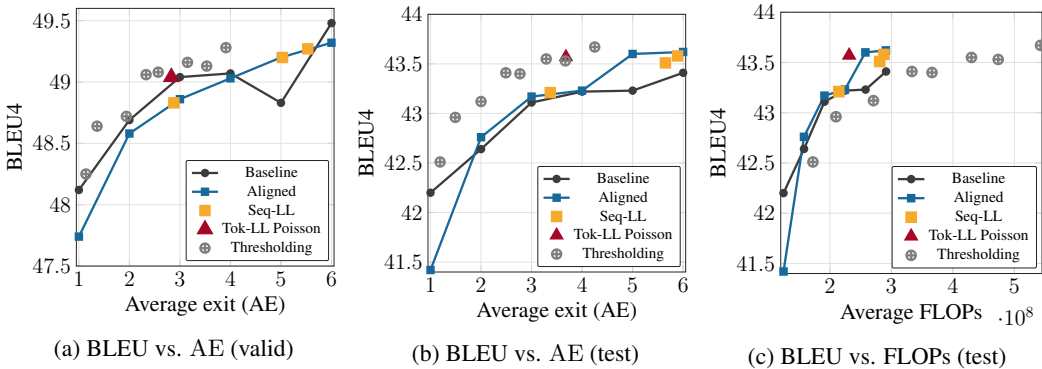


Figure 4: Speed and accuracy on the WMT’14 English-French benchmark (*c.f.* Figure 3).

the valid set, parity was already achieved at two blocks (*c.f.* Table 1). The aligned model outperforms the baseline for $n = 2, \dots, 6$. This is likely due the regularization effect of jointly training six outputs classifiers, one on top of each block, leading to a strong training signal for the early blocks.

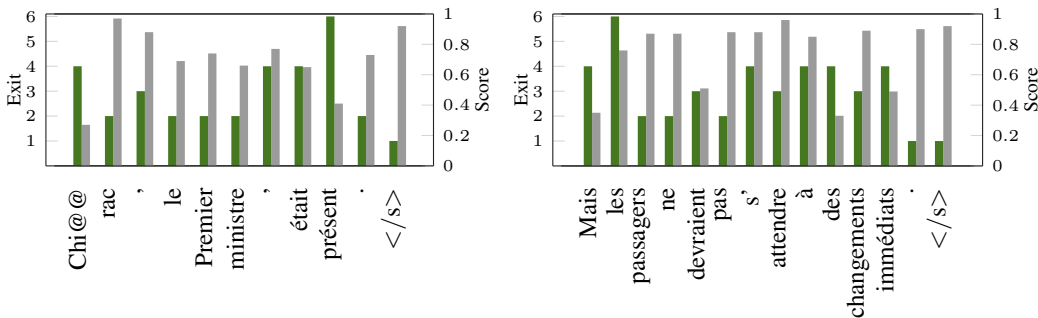
For token specific halting mechanisms (Figure 3a) the likelihood oracle tends to outperforms the correctness oracle and the Poisson classifier achieves a better speed-accuracy trade-off than the multinomial. At the sequence-level, likelihood is also the better oracle (Figure 3b). Tok-LL Poisson with $(\sigma, \lambda) = (2, 0.2)$ achieves $AE = 2.9$ and 34.9 BLEU which corresponds to similar accuracy as the $N = 6$ baseline at 52% fewer decoding blocks. Seq-LL using $\lambda = 0.02$ with $AE = 2.1$ and BLEU 34.9 has slightly higher accuracy than the $N = 6$ baseline at 65% fewer blocks.

Confidence thresholding (Figure 3c) performs very well for $AE < 3$ while Tok-LL Poisson only achieves a marginal improvement over the aligned model with fixed exits in this AE range; this may be due to the large depth regularization (high λ ; see Eq. (11)) required to achieve very low AE.

4.4 SCALING THE ADAPTIVE-DEPTH MODELS

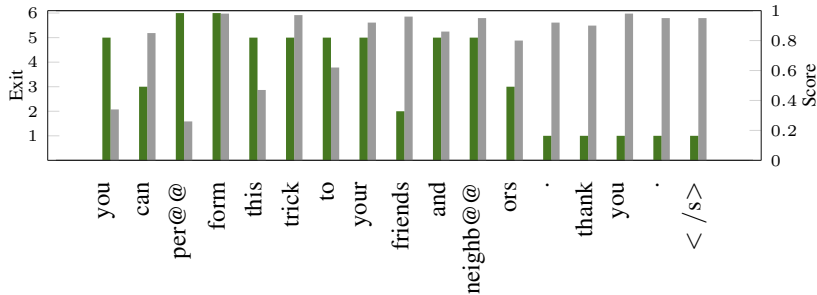
Finally, we take the best performing models form the IWSLT benchmark and test them on the large WMT’14 English-French benchmark. Results on the valid set (Figure 4a) and the test set (Figure 4b) show that confidence thresholding works very well and that the sequence-specific depth approach does not improve over the baseline. Tok-LL Poisson achieves 43.6 BLEU which slightly improves over the $N=6$ baseline (43.4 BLEU) at $AE = 3.7$ or 40% fewer blocks.

Confidence thresholding matches the accuracy of the $N=6$ baseline with $AE 2.5$ or 59% fewer decoding blocks. However, confidence thresholding requires computing the output classifier at each block to determine whether to halt or continue. This is a large overhead since output classifiers



(a) **Src:** Chi@@rac , the Prime Minister , was there . **Ref:** Chi@@rac , Premier ministre , est là .
 (b) **Src:** But passengers shoul@@dn't expect changes to happen immediatly . **Ref:** Mais les passagers ne devraient pas s' attendre à des changements immédiats .

Figure 5: Examples from the WMT'14 En-Fr test set for Tok-LL Poisson. Token exits are in black and model scores (probabilities) are in gray. The '@@' are due to BPE or subword tokenization.



(a) **Src:** diesen trick können sie ihren freunden und nachbarn vor@@führen . danke . **Ref:** there is a trick you can do for your friends and neighb@@ors . thanks .

Figure 6: Example from the IWSLT De-En test set using Tok-LL Poisson. See Figure 5.

predict 44k types for this benchmark (§4.1). To better account for this, we measure the average number of FLOPs per output token (details in Appendix B). Figure 4c shows that the Tok-LL Poisson approach provides a better trade-off when the overhead of the output classifiers is considered.

4.5 QUALITATIVE RESULTS

Figure 5 shows outputs for examples of the WMT'14 En-Fr test set together with the exit and model probability for each token. Less computation is used at the end of the sentence since periods and end of sentence markers (</s>) are easy to predict. The amount of computation increases when the model is less confident *e.g.* in Figure 5a, predicting 'présent' (meaning 'present') is hard. A straightforward translation is 'était là' but the model chooses 'present' which is also appropriate. In Figure 5b, the model uses more computation to predict the definite article 'les' since the source has omitted the article for 'passengers'. In the IWSLT'14 De-En example (Figure 6) the model expends very little computation on translating the simple phrase 'danke' (meaning 'thank you').

5 CONCLUSION

We extended anytime prediction to the structured prediction setting and introduced simple but highly effective methods to equip sequence models to make predictions at different points in the network. We compared a number of different mechanisms to predict the required network depth and find that a simple likelihood based Poisson classifier obtains the best trade-off between speed and accuracy. Our results show that the number of decoder layers can be vastly reduced at no loss in accuracy.

REFERENCES

- Alexei Baevski and Michael Auli. Adaptive input representations for neural language modeling. In *Proc. of ICLR*, 2019.
- Asli Celikyilmaz, Antoine Bosselut, Xiaodong He, and Yejin Choi. Deep communicating agents for abstractive summarization. In *Proc. of ICLR*, 2018.
- M. Cettolo, J. Niehues, S. Stüker, L. Bentivogli, and M. Federico. Report on the 11th iwslt evaluation campaign. In *IWSLT*, 2014.
- Mostafa Dehghani, Stephan Gouws, Oriol Vinyals, Jakob Uszkoreit, and Lukasz Kaiser. Universal transformers. In *Proc. of ICLR*, 2018.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *Proc. of NAACL*, 2019.
- Sergey Edunov, Myle Ott, Michael Auli, David Grangier, and Marc’Aurelio Ranzato. Classical structured prediction losses for sequence to sequence learning. In *Proc. of NAACL*, 2018.
- Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. Convolutional sequence to sequence learning. In *Proc. of ICML*, 2017.
- Alex Graves. Adaptive computation time for recurrent neural networks. In *ArXiv preprint*, 2016.
- Gao Huang, Danlu Chen, Tianhong Li, Felix Wu, Laurens van der Maaten, and Kilian Q Weinberger. Multi-scale dense networks for resource efficient image classification. In *Proc. of ICLR*, 2017.
- D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *Proc. of ICLR*, 2015.
- Nathan Ng, Kyra Yee, Alexei Baevski, Myle Ott, Michael Auli, and Sergey Edunov. Facebook fair’s wmt19 news translation task submission. In *Proc. of WMT*, 2019.
- Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. Fairseq: A fast, extensible toolkit for sequence modeling. In *Proc. of NAACL*, 2019.
- K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. BLEU: a method for automatic evaluation of machine translation. In *Proc. of ACL*, 2002.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. In *preprint*, 2018.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. In *Technical report, OpenAI.*, 2019.
- R. Sennrich, B. Haddow, and A. Birch. Neural machine translation of rare words with subword units. In *Proc. of ACL*, 2016.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. Gomez, L. Kaiser, and I. Polosukhin. Attention is all you need. In *Proc. of NeurIPS*, 2017.
- Felix Wu, Angela Fan, Alexei Baevski, Yann N. Dauphin, and Michael Auli. Pay less attention with lightweight and dynamic convolutions. In *Proc. of ICLR*, 2019.

APPENDIX A LOSS SCALING

In this section we experiment with different weights for scaling the output classifier losses. Instead of uniform weighting, we bias towards specific output classifiers by assigning higher weights to their losses. Table 2 shows that weighing the classifiers equally provides good results.

	Uniform	$n = 1$	$n = 2$	$n = 3$	$n = 4$	$n = 5$	$n = 6$	Average
Baseline	-	34.5	35.5	35.8	35.7	35.8	36.0	35.5
$\omega_n = 1$	35.5	34.1	35.5	35.8	36.1	36.1	36.2	35.6
$\omega_n = n$	35.3	32.2	35.0	35.8	36.0	36.2	36.3	35.2
$\omega_n = \sqrt{n}$	35.4	33.3	35.2	35.8	35.9	36.1	36.1	35.4
$\omega_n = 1/\sqrt{n}$	35.6	34.5	35.4	35.7	35.8	35.8	35.9	35.5
$\omega_n = 1/n$	35.3	34.7	35.3	35.5	35.7	35.8	35.8	35.5

(a) IWSLT De-En - Valid

	Uniform	$n = 1$	$n = 2$	$n = 3$	$n = 4$	$n = 5$	$n = 6$	Average
Baseline	-	33.7	34.6	34.6	34.6	34.6	34.8	34.5
$\omega_n = 1$	34.4	33.2	34.4	34.8	34.9	35.0	34.9	34.5
$\omega_n = n$	34.2	31.4	33.8	34.7	34.8	34.8	34.9	34.1
$\omega_n = \sqrt{n}$	34.4	32.5	34.1	34.8	34.9	35.0	35.1	34.4
$\omega_n = 1/\sqrt{n}$	34.6	33.7	34.3	34.6	34.8	34.8	34.9	34.5
$\omega_n = 1/n$	34.2	33.8	34.3	34.5	34.6	34.7	34.7	34.4

(b) IWSLT De-En - Test

Table 2: Aligned training with different weights (ω_n) on IWSLT De-En. For each model we report BLEU on the dev set evaluated with a uniformly sampled exit $n \sim \mathcal{U}([1..6])$ for each token and a fixed exit $n \in [1..6]$ throughout the sequence. The average corresponds to the average BLEU over the fixed exits.

Gradient scaling Adding intermediate supervision at different levels of the decoder results in richer gradients for lower blocks compared to upper blocks. This is because earlier layers affect more loss terms in the compound loss of Eq. (1). To balance the gradients of each block in the decoder, we scale up the gradients of each loss term ($-\text{LL}_n$) when it is updating the parameters of its associated block (block_n with parameters θ_n) and revert it back to its normal scale before back-propagating it to the previous blocks. Figure 7 illustrates this gradient scaling procedure. The θ_n are updated with γ_n -amplified gradients from the block’s supervision and $(N-n)$ gradients from the subsequent blocks. We choose $\gamma_n = \gamma(N-n)$ to control the ratio $\gamma:1$ as the ratio of the block supervision to the subsequent blocks’ supervisions.

Table 3 shows that gradient scaling can benefit the lowest layer at the expense of higher layers. However, no scaling generally works very well.

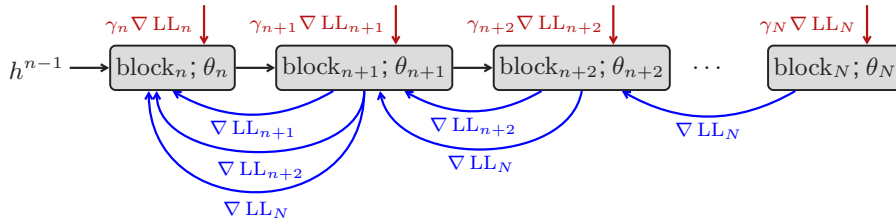


Figure 7: Illustration of gradient scaling.

	Uniform	$n = 1$	$n = 2$	$n = 3$	$n = 4$	$n = 5$	$n = 6$	Average
Baseline	-	34.5	35.5	35.8	35.7	35.8	36.0	35.5
\emptyset	35.5	34.1	35.5	35.8	36.1	36.1	36.2	35.6
$\gamma = 0.3$	35.1	33.7	34.7	35.3	35.7	35.8	36.0	35.2
$\gamma = 0.5$	35.4	34.8	35.4	35.6	35.6	35.7	35.6	35.4
$\gamma = 0.7$	34.9	34.6	35.1	35.1	35.2	35.4	35.3	35.1
$\gamma = 0.9$	34.9	34.8	35.3	35.3	35.3	35.4	35.5	35.3
$\gamma = 1.1$	35.1	34.9	35.2	35.3	35.3	35.3	35.3	35.2

(a) IWSLT De-En - Valid

	Uniform	$n = 1$	$n = 2$	$n = 3$	$n = 4$	$n = 5$	$n = 6$	Average
Baseline	-	33.7	34.6	34.6	34.6	34.6	34.8	34.5
\emptyset	34.4	33.2	34.4	34.8	34.9	35.0	34.9	34.5
$\gamma = 0.3$	34.2	32.8	33.9	34.3	34.6	34.8	35.0	34.2
$\gamma = 0.5$	34.5	33.8	34.2	34.6	34.5	34.7	34.7	34.6
$\gamma = 0.7$	34.0	33.7	34.2	34.3	34.3	34.3	34.3	34.2
$\gamma = 0.9$	34.1	34.0	34.2	34.3	34.4	34.4	34.4	34.3
$\gamma = 1.1$	34.2	34.0	34.3	34.3	34.3	34.3	34.2	34.2

(b) IWSLT De-En - Test

Table 3: Aligned training with different gradient scaling ratios $\gamma : 1$ on IWSLT’14 De-En. For each model we report the BLEU4 score evaluated with a uniformly sampled exit $n \sim \mathcal{U}([1..6])$ for each token and a fixed exit $n \in [1..6]$. The average corresponds to the average BLEU4 of all fixed exits.

APPENDIX B FLOPS APPROXIMATION

This section details the computation of the FLOPS we report. The per token FLOPS are for the decoder network only since we use an encoder of the same size for all models. We breakdown the FLOPS of every operation in Algorithm 1 (blue front of the algorithmic statement). We omit non-linearities, normalizations and residual connections. The main operations we account for are dot-products and by extension matrix-vector products since those represent the vast majority of FLOPS (we assume batch size one to simplify the calculation).

Parameters		Operation	FLOPS
d_d	decoder embedding dimension.	Dot-product (d)	$2d - 1$
d_e	encoder embedding dimension.	Linear $d_{in} \rightarrow d_{out}$	$2d_{in}d_{out}$
d_f	ffn dimension.		
$ \mathbf{x} $	source length.		
t	Current time-estep ($t \geq 1$).		
V	output vocabulary size.		

Table 4: FLOPS of basic operations, key parameters and variables for the FLOPS estimation.

With this breakdown, the total computational cost at time-step t of a decoder block that we actually go through is:

$$\text{FC}(\mathbf{x}, t) = 12d_d^2 + 4d_f d_d + 4t d_d + 4|\mathbf{x}|d_d + 4[\text{FirstCall}]|\mathbf{x}|d_d d_e,$$

where the cost of mapping the source’ keys and values is incurred the first time the block is called (flagged with FirstCall). This occurs at $t = 1$ for the baseline model but it is input-dependent with depth adaptive estimation and may never occur if all tokens exit early.

If skipped, a block still has to compute the keys and value of the self-attention block so the self-attention of future time-steps can function. This costs $\text{FS} = 4d_d^2$.

Depending on the halting mechanism, an exit prediction cost (FP) is added:

Sequence-specific depth:	$FP(t, q(t)) = 2\llbracket t = 1 \rrbracket N d_d$
Token-specific Multinomial:	$FP(t, q(t)) = 2N d_d$
Token-specific Poisson binomial:	$FP(t, q(t)) = 2d_d q(t)$
Confidence thresholding:	$FP(t, q(t)) = 2q(t) V d_d$

For a set of source sequences $\{\mathbf{x}^{(i)}\}_{i \in \mathcal{I}}$ and generated hypotheses $\{\mathbf{y}^{(i)}\}_{i \in \mathcal{I}}$, the average flops per token is:

$$\begin{aligned} \text{Baseline (} N \text{ blocks):} & \quad \frac{1}{\sum_i |\mathbf{y}^{(i)}|} \sum_i \sum_{t=1}^{|\mathbf{y}^{(i)}|} (N \text{ FC}(\mathbf{x}^{(i)}, t) + 2V d_d) \\ \text{Adaptive depth:} & \quad \frac{1}{\sum_i |\mathbf{y}^{(i)}|} \sum_i \sum_{t=1}^{|\mathbf{y}^{(i)}|} (q(t) \text{FC}(\mathbf{x}^{(i)}, t) + (N - q(t)) \text{FS} + \text{FP}(t, q(t)) + 2V d_d) \end{aligned}$$

In the case of confidence thresholding the final output prediction cost ($2V d_d$) is already accounted for in the exit prediction cost FP.

Algorithm 1 Adaptive decoding with Tok-Poisson

```

1: Input: source codes  $s$ , incremental state
2: Initialization:  $t = 1, y_1 = \langle s \rangle$ 
3: for  $n \in 1 \dots N$  do
4:   FirstCall[ $n$ ] = True.  ▷ A flag signaling if the source' keys and values should be evaluated.
5: end for
6: while  $y_t \neq \langle /s \rangle$  do
7:   Embed the last output token  $y_t$ .
8:   for  $n \in 1 \dots N$  do
9:     ▷ Self-attention.
10:    - Map the input into a key ( $k$ ) and value ( $v$ ). FLOPS= $4d_d^2$ 
11:    - Map the input into a query  $q$ . FLOPS= $2d_d^2$ 
12:    - Score the memory keys with  $q$  to get the attention weights  $\alpha$ . FLOPS= $4td_d$ 
13:    - Map the attention output. FLOPS= $2d_d^2$ 
14:    ▷ Encoder-Decoder interaction.
15:    if FirstCall[ $n$ ] then
16:      Map the source states into keys and values for the  $n$ th block. FLOPS= $4|x|d_e d_d$ 
17:      FirstCall[ $n$ ] = False
18:    end if
19:    - Map the input into a query  $q$ . FLOPS= $2d_d^2$ 
20:    - Score the memory keys with  $q$  to get the attention weights  $\alpha$ . FLOPS= $4|x|d_d$ 
21:    - Map the attention output. FLOPS= $2d_d^2$ 
22:    Feed-forward network. FLOPS= $4d_d d_f$ 
23:    Estimate the halting probability  $\chi_{t,n}$ . FLOPS= $2d_d$ 
24:    if  $\chi_{t,n} > 0.5$  then
25:      Exit the loop (Line 8)
26:    end if
27:  end for
28:  if  $n < N$  then
29:    ▷ Skipped blocks.
30:    for  $n_s \in n + 1 \dots N$  do
31:      Copy and map the copied state into a key ( $k$ ) and value ( $v$ ). FLOPS= $4d_d^2$ 
32:    end for
33:  end if
34:  Project the final state and sample a new output token. FLOPS= $2V d_d$ 
35:   $t++$ 
36: end while

```
