

## A. Weight Initialization

We derive a weight initialization scheme tailored to the GLU activation function similar to Glorot & Bengio (2010); He et al. (2015b) by focusing on the variance of activations within the network for both forward and backward passes. We also detail how we modify the weight initialization for dropout.

### A.1. Forward Pass

Assuming that the inputs  $\mathbf{x}_l$  of a convolutional layer  $l$  and its weights  $W_l$  are independent and identically distributed (i.i.d.), the variance of its output, computed as  $\mathbf{y}_l = W_l \mathbf{x}_l + \mathbf{b}_l$ , is

$$\text{Var}[\mathbf{y}_l] = n_l \text{Var}[w_l x_l] \quad (3)$$

where  $n_l$  is the number inputs to the layer. For one-dimensional convolutional layers with kernel width  $k$  and input dimension  $c$ , this is  $kc$ . We adopt the notation in (He et al., 2015b), i.e.  $y_l$ ,  $w_l$  and  $x_l$  represent the random variables in  $\mathbf{y}_l$ ,  $W_l$  and  $\mathbf{x}_l$ . With  $w_l$  and  $x_l$  independent from each other and normally distributed with zero mean, this amounts to

$$\text{Var}[\mathbf{y}_l] = n_l \text{Var}[w_l] \text{Var}[x_l]. \quad (4)$$

$x_l$  is the result of the GLU activation function  $\mathbf{y}_{l-1}^a \sigma(\mathbf{y}_{l-1}^b)$  with  $\mathbf{y}_{l-1} = (\mathbf{y}_{l-1}^a, \mathbf{y}_{l-1}^b)$ , and  $\mathbf{y}_{l-1}^a, \mathbf{y}_{l-1}^b$  i.i.d. Next, we formulate upper and lower bounds in order to approximate  $\text{Var}[x_l]$ . If  $\mathbf{y}_{l-1}$  follows a symmetric distribution with mean 0, then

$$\text{Var}[x_l] = \text{Var}[\mathbf{y}_{l-1}^a \sigma(\mathbf{y}_{l-1}^b)] \quad (5)$$

$$= E[(\mathbf{y}_{l-1}^a \sigma(\mathbf{y}_{l-1}^b))^2] - E^2[\mathbf{y}_{l-1}^a \sigma(\mathbf{y}_{l-1}^b)] \quad (6)$$

$$= \text{Var}[\mathbf{y}_{l-1}^a] E[\sigma(\mathbf{y}_{l-1}^b)^2]. \quad (7)$$

A lower bound is given by  $(1/4)\text{Var}[\mathbf{y}_{l-1}^a]$  when expanding (6) with  $E^2[\sigma(\mathbf{y}_{l-1}^b)] = 1/4$ :

$$\text{Var}[x_l] = \text{Var}[\mathbf{y}_{l-1}^a \sigma(\mathbf{y}_{l-1}^b)] \quad (8)$$

$$= \text{Var}[\mathbf{y}_{l-1}^a] E^2[\sigma(\mathbf{y}_{l-1}^b)] + \quad (9)$$

$$\text{Var}[\mathbf{y}_{l-1}^a] \text{Var}[\sigma(\mathbf{y}_{l-1}^b)] \\ = \frac{1}{4} \text{Var}[\mathbf{y}_{l-1}^a] + \text{Var}[\mathbf{y}_{l-1}^a] \text{Var}[\sigma(\mathbf{y}_{l-1}^b)] \quad (10)$$

and  $\text{Var}[\mathbf{y}_{l-1}^a] \text{Var}[\sigma(\mathbf{y}_{l-1}^b)] > 0$ . We utilize the relation  $\sigma(x)^2 \leq (1/16)x^2 - 1/4 + \sigma(x)$  (Appendix B) to provide an upper bound on  $E[\sigma(x)^2]$ :

$$E[\sigma(x)^2] \leq E\left[\frac{1}{16}x^2 - \frac{1}{4} + \sigma(x)\right] \quad (11)$$

$$= \frac{1}{16}E[x^2] - \frac{1}{4} + E[\sigma(x)] \quad (12)$$

With  $x \sim \mathcal{N}(0, \text{std}(x))$ , this yields

$$E[\sigma(x)^2] \leq \frac{1}{16}E[x^2] - \frac{1}{4} + \frac{1}{2} \quad (13)$$

$$= \frac{1}{16}\text{Var}[x] + \frac{1}{4}. \quad (14)$$

With (7) and  $\text{Var}[\mathbf{y}_{l-1}^a] = \text{Var}[\mathbf{y}_{l-1}^b] = \text{Var}[\mathbf{y}_{l-1}]$ , this results in

$$\text{Var}[x_l] \leq \frac{1}{16}\text{Var}[\mathbf{y}_{l-1}]^2 + \frac{1}{4}\text{Var}[\mathbf{y}_{l-1}]. \quad (15)$$

We initialize the embedding matrices in our network with small variances (around 0.01), which allows us to dismiss the quadratic term and approximate the GLU output variance with

$$\text{Var}[x_l] \approx \frac{1}{4}\text{Var}[\mathbf{y}_{l-1}]. \quad (16)$$

If  $L$  network layers of equal size and with GLU activations are combined, the variance of the final output  $\mathbf{y}_L$  is given by

$$\text{Var}[\mathbf{y}_L] \approx \text{Var}[\mathbf{y}_1] \left( \prod_{l=2}^L \frac{1}{4} n_l \text{Var}[w_l] \right). \quad (17)$$

Following (He et al., 2015b), we aim to satisfy the condition

$$\frac{1}{4} n_l \text{Var}[w_l] = 1, \forall l \quad (18)$$

so that the activations in a network are neither exponentially magnified nor reduced. This is achieved by initializing  $W_l$  from  $\mathcal{N}(0, \sqrt{4/n_l})$ .

### A.2. Backward Pass

The gradient of a convolutional layer is computed via back-propagation as  $\Delta \mathbf{x}_l = \hat{W}_l \mathbf{y}_l$ . Considering separate gradients  $\Delta \mathbf{y}_l^a$  and  $\Delta \mathbf{y}_l^b$  for GLU, the gradient of  $\mathbf{x}$  is given by

$$\Delta \mathbf{x}_l = \hat{W}_l^a \Delta \mathbf{y}_l^a + \hat{W}_l^b \Delta \mathbf{y}_l^b. \quad (19)$$

$\hat{W}$  corresponds to  $W$  with re-arranged weights to enable back-propagation. Analogously to the forward pass,  $\Delta x_l$ ,  $\hat{w}_l$  and  $\Delta y_l$  represent the random variables for the values in  $\Delta \mathbf{x}_l$ ,  $\hat{W}_l$  and  $\Delta \mathbf{y}_l$ , respectively. Note that  $W$  and  $\hat{W}$  contain the same values, i.e.  $\hat{w} = w$ . Similar to (3), the variance of  $\Delta x_l$  is

$$\text{Var}[\Delta x_l] = \hat{n}_l \left( \text{Var}[w_l^a] \text{Var}[\Delta y_l^a] + \text{Var}[w_l^b] \text{Var}[\Delta y_l^b] \right). \quad (20)$$

Here,  $\hat{n}_l$  is the number of inputs to layer  $l+1$ . The gradients for the GLU inputs are:

$$\Delta \mathbf{y}_l^a = \Delta \mathbf{x}_{l+1} \sigma(\mathbf{y}_l^b) \quad \text{and} \quad (21)$$

$$\Delta \mathbf{y}_l^b = \Delta \mathbf{x}_{l+1} \mathbf{y}_l^a \sigma'(\mathbf{y}_l^b). \quad (22)$$

The approximation for the forward pass can be used for  $\text{Var}[\Delta y_l^a]$ , and for estimating  $\text{Var}[\Delta y_l^b]$  we assume an upper bound on  $E[\sigma'(y_l^b)^2]$  of  $1/16$  since  $\sigma'(y_l^b) \in [0, \frac{1}{4}]$ . Hence,

$$\text{Var}[\Delta y_l^a] - \frac{1}{4} \text{Var}[\Delta x_{l+1}] \leq \frac{1}{16} \text{Var}[\Delta x_{l+1}] \text{Var}[y_l^b] \quad (23)$$

$$\text{Var}[\Delta y_l^b] \leq \frac{1}{16} \Delta \text{Var}[\Delta x_{l+1}] \text{Var}[y_l^a] \quad (24)$$

We observe relatively small gradients in our network, typically around 0.001 at the start of training. Therefore, we approximate by discarding the quadratic terms above, i.e.

$$\text{Var}[\Delta y_l^a] \approx \frac{1}{4} \text{Var}[\Delta x_{l+1}] \quad (25)$$

$$\text{Var}[\Delta y_l^b] \approx 0 \quad (26)$$

$$\text{Var}[\Delta x_l] \approx \frac{1}{4} \hat{n}_l \text{Var}[w_l^a] \text{Var}[\Delta x_{l+1}] \quad (27)$$

As for the forward pass, the above result can be generalized to backpropagation through many successive layers, resulting in

$$\text{Var}[\Delta x_2] \approx \text{Var}[\Delta x_{L+1}] \left( \prod_{l=2}^L \frac{1}{4} \hat{n}_l \text{Var}[w_l^a] \right) \quad (28)$$

and a similar condition, i.e.  $(1/4)\hat{n}_l \text{Var}[w_l^a] = 1$ . In the networks we consider, successions of convolutional layers usually operate on the same number of inputs so that most cases  $n_l = \hat{n}_l$ . Note that  $W_l^b$  is discarded in the approximation; however, for the sake of consistency we use the same initialization for  $W_l^a$  and  $W_l^b$ .

For arbitrarily large variances of network inputs and activations, our approximations are invalid; in that case, the initial values for  $W_l^a$  and  $W_l^b$  would have to be balanced for the input distribution to be retained. Alternatively, methods that explicitly control the variance in the network, e.g. batch normalization (Ioffe & Szegedy, 2015) or layer normalization (Ba et al., 2016) could be employed.

### A.3. Dropout

Dropout retains activations in a neural network with a probability  $p$  and sets them to zero otherwise (Srivastava et al., 2014). It is common practice to scale the retained activations by  $1/p$  during training so that the weights of the network do not have to be modified at test time when  $p$  is set to 1. In this case, dropout amounts to multiplying activations  $x$  by a Bernoulli random variable  $r$  where  $\Pr[r=1/p] = p$  and  $\Pr[r=0] = 1-p$  (Srivastava et al., 2014). It holds that  $E[r] = 1$  and  $\text{Var}[r] = (1-p)/p$ . If  $x$  is independent of  $r$  and  $E[x] = 0$ , the variance after dropout is

$$\text{Var}[xr] = E[r]^2 \text{Var}[x] + \text{Var}[r] \text{Var}[x] \quad (29)$$

$$= \left(1 + \frac{1-p}{p}\right) \text{Var}[x] \quad (30)$$

$$= \frac{1}{p} \text{Var}[x] \quad (31)$$

Assuming that a the *input* of a convolutional layer has been subject to dropout with a retain probability  $p$ , the variations of the forward and backward activations from §A.1 and §A.2 can

now be approximated with

$$\text{Var}[x_{l+1}] \approx \frac{1}{4p} n_l \text{Var}[w_l] \text{Var}[x_l] \quad \text{and} \quad (32)$$

$$\text{Var}[\Delta x_l] \approx \frac{1}{4p} n_l \text{Var}[w_l^a] \text{Var}[\Delta x_{l+1}]. \quad (33)$$

This amounts to a modified initialization of  $W_l$  from a normal distribution with zero mean and a standard deviation of  $\sqrt{4p/n}$ . For layers without a succeeding GLU activation function, we initialize weights from  $\mathcal{N}(0, \sqrt{p/n})$  to calibrate for any immediately preceding dropout application.

## B. Upper Bound on Squared Sigmoid

The sigmoid function  $\sigma(x)$  can be expressed as a hyperbolic tangent by using the identity  $\tanh(x) = 2\sigma(2x) - 1$ . The derivative of  $\tanh$  is  $\tanh'(x) = 1 - \tanh^2(x)$ , and with  $\tanh(x) \in [0, 1], x \geq 0$  it holds that

$$\tanh'(x) \leq 1, x \geq 0 \quad (34)$$

$$\int_0^x \tanh'(x) dx \leq \int_0^x 1 dx \quad (35)$$

$$\tanh(x) \leq x, x \geq 0 \quad (36)$$

We can express this relation with  $\sigma(x)$  as follows:

$$2\sigma(x) - 1 \leq \frac{1}{2}x, x \geq 0 \quad (37)$$

Both terms of this inequality have rotational symmetry w.r.t 0, and thus

$$(2\sigma(x) - 1)^2 \leq \left(\frac{1}{2}x\right)^2 \quad \forall x \quad (38)$$

$$\Leftrightarrow \sigma(x)^2 \leq \frac{1}{16}x^2 - \frac{1}{4}\sigma(x). \quad (39)$$

## C. Attention Visualization

Figure 3 shows attention scores for a generated sentence from the WMT'14 English-German task. The model used for this plot has 8 decoder layers and a 80K BPE vocabulary. The attention passes in different decoder layers capture different portions of the source sentence. Layer 1, 3 and 6 exhibit a linear alignment. The first layer shows the clearest alignment, although it is slightly off and frequently attends to the corresponding source word of the previously generated target word. Layer 2 and 8 lack a clear structure and are presumably collecting information about the whole source sentence. The fourth layer shows high alignment scores on nouns such as "festival", "way" and "work" for both the generated target nouns as well as their preceding words. Note that in German, those preceding words depend on gender and object relationship of the respective noun. Finally, the attention scores in layer 5 and 7 focus on "built", which is reordered in the German translation and is moved from the beginning to the very end of

the sentence. One interpretation for this is that as generation progresses, the model repeatedly tries to perform the re-ordering. “aufgebaut” can be generated after a noun or pronoun only, which is reflected in the higher scores at positions 2, 5, 8, 11 and 13.

# Convolutional Sequence to Sequence Learning

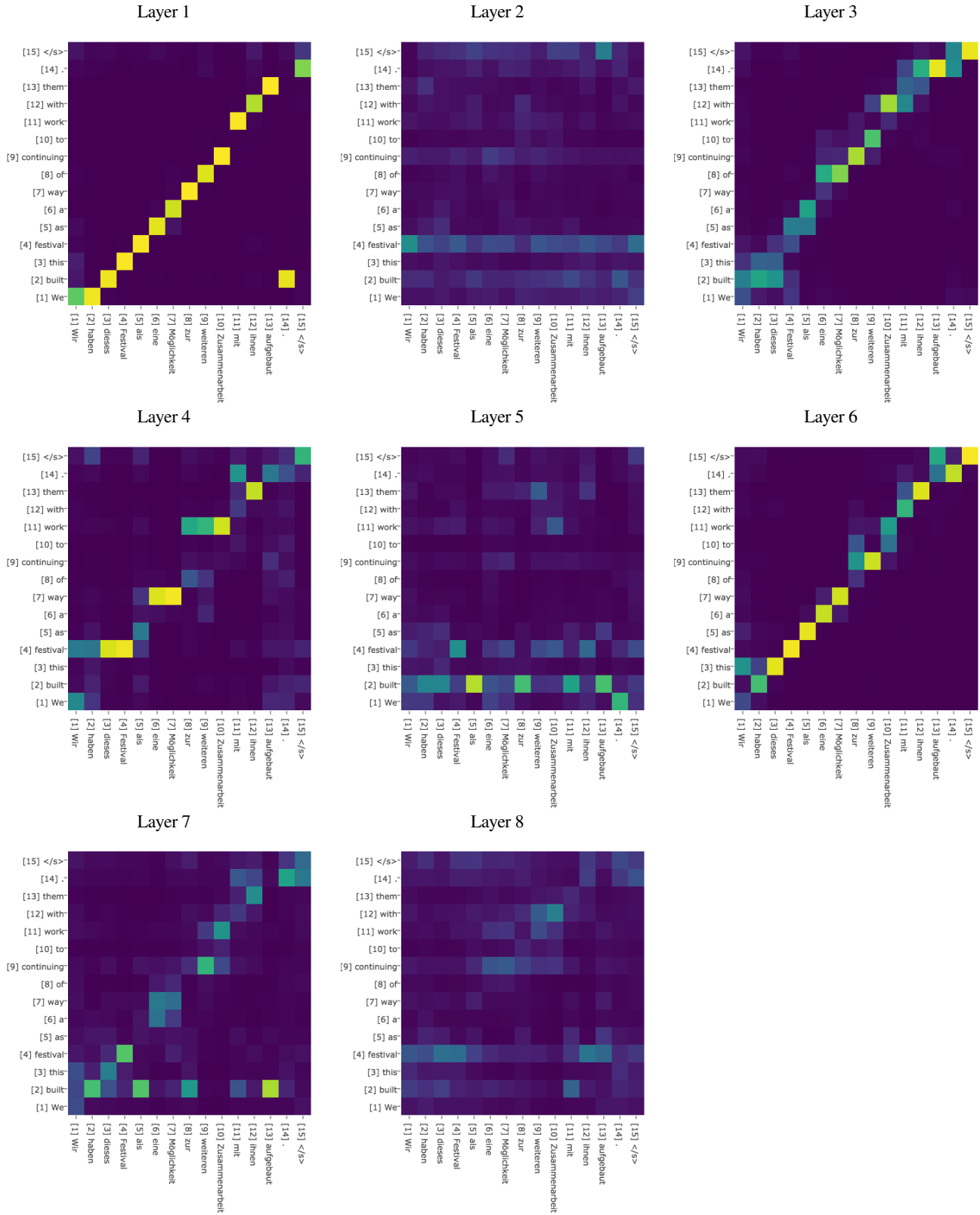


Figure 3. Attention scores for different decoder layers for a sentence translated from English (y-axis) to German (x-axis). This model uses 8 decoder layers and a 80k BPE vocabulary.