

# 1 Justification for Construction 3.1

Set  $A^{(1)}$  to be the (scaled) adjacency matrix of  $\mathcal{G}$ , i.e.  $A_{i,j}^{(1)} = \beta_1 \mathbf{1}(j = p(i))$ , and  $A^{(2)} = \beta_2 I_S$ , where  $\beta_1, \beta_2 \rightarrow \infty$ . We will now show that the output of the disentangled transformer approximates  $\hat{\pi}_{s_{1:T}}(\cdot \mid s_T)$ .

**First Attention.** Note that by the construction of  $\tilde{A}^{(1)}$ ,  $\tilde{X} \tilde{A}^{(1)} \tilde{X}^\top = A^{(1)}$ , which is the scaled adjacency matrix of  $\mathcal{G}$ . If  $i$  is not a root node (i.e.  $i \in \overline{\mathcal{R}}$ ,  $p(i) \neq \emptyset$ ), then

$$\mathcal{S}(\tilde{X} \tilde{A}^{(1)} \tilde{X}^\top)_{i,j} = \mathbf{1}(j = p(i)) \quad (1)$$

so  $i$  attends to its parent  $p(i)$ . Therefore, the output of the first attention is the token at position  $p(i)$ , i.e.  $\text{attn}(\tilde{X}; \tilde{A}^{(1)})_i = \tilde{x}_{p(i)}$ . The transformer then appends  $\tilde{x}_{p(i)}$  to the residual stream of token  $i$ .

When  $i$  is a root node (i.e.  $i \in \mathcal{R}$ ,  $p(i) = \emptyset$ ), then for all  $j$ ,  $(\tilde{X} \tilde{A}^{(1)} \tilde{X}^\top)_{ij} = 0$ . Therefore after the softmax,  $i$  will attend equally to all previous tokens:

$$\mathcal{S}(\tilde{X} \tilde{A}^{(1)} \tilde{X}^\top)_{i,j} = \frac{1}{i} \quad \text{for all } j \leq i. \quad (2)$$

Thus the first attention layer averages all of the tokens in the sequence:  $\text{attn}(\tilde{X}; \tilde{A}^{(1)})_i = \frac{1}{i} \sum_{j \leq i} \tilde{x}_j$ . It then copies this average into the residual stream.

**Second Attention.** We next show that the  $T$ th token attends to all prior tokens whose parents tokens are equal to  $s_T$ . It then averages them and copies them into the residual stream.

After the first attention layer, the residual stream is  $h_j^{(1)} = [\tilde{x}_j, \text{attn}(\tilde{X}; \tilde{A}^{(1)})_j]^\top$ . The second attention layer compares the  $T$ th token of the original sequence  $\tilde{x}_T$  to the output of the first attention at all other positions. Explicitly, the attention pattern is equal to:

$$h_T^{(1)\top} \tilde{A}^{(2)} h_j^{(1)} = \beta_2 \cdot \tilde{x}_T^\top \begin{bmatrix} A^{(2)} & 0_{S \times T} \\ 0_{T \times S} & 0_{T \times T} \end{bmatrix} \text{attn}(\tilde{X}; \tilde{A}^{(1)})_j = \beta_2 \cdot \begin{cases} \mathbf{1}(s_{p(i)} = s_T) & i \in \overline{\mathcal{R}} \\ \frac{1}{i} \sum_{j \leq i} \mathbf{1}(s_j = s_T) & i \in \mathcal{R}. \end{cases} \quad (3)$$

As  $\beta_2 \rightarrow \infty$ , the softmax converges to a hard max, and so the  $T$ th token attends equally to all tokens  $i$  such that  $s_{p(i)} = s_T$ . The attention then averages all of these tokens, so the  $T$ th token in the residual stream is equal to  $h_T^{(2)} = [\tilde{x}_T, \frac{1}{T} \sum_{j \leq T} \tilde{x}_j, Z, \tilde{x}_T]$  where

$$Z := \frac{\sum_{s_{p(i)} = s_T} \tilde{x}_i}{|\{i : s_{p(i)} = s_T\}|} \quad (4)$$

is the average of the tokens whose parent is equal to  $s_T$ .

**Output Layer.**  $W_O$  reads from the third block in this stream, which we denoted by  $Z$  in (4) above. It then returns the token embedding of  $Z$  which is equal to:

$$f_{\tilde{\theta}}(s_{1:T}) = \frac{\sum_{s_{p(i)}=s_T} e_{s_i}}{|\{i : s_{p(i)} = s_T\}|} = \hat{\pi}_{s_{1:T}}(\cdot | s_T), \quad (5)$$

as desired.

## 2 Justification for Equation 3

The output of the first attention layer is

$$\text{attn}(\tilde{X}; \tilde{A}^{(1)}) = \mathcal{S}(\text{MASK}(\tilde{X} \tilde{A}^{(1)} \tilde{X}^\top) \tilde{X} = \mathcal{S}(\text{MASK}(A^{(1)})) \tilde{X}.$$

Next, we have that

$$\begin{aligned} h_T^{(1)\top} \tilde{A}^{(2)} h_T^{(1)\top} &= \tilde{x}_T^\top \begin{bmatrix} A^{(2)} & 0_{S \times T} \\ 0_{T \times S} & 0_{T \times T} \end{bmatrix} \text{attn}(\tilde{X}; \tilde{A}^{(1)})^\top \\ &= \tilde{x}_T^\top \begin{bmatrix} A^{(2)} & 0_{S \times T} \\ 0_{T \times S} & 0_{T \times T} \end{bmatrix} \tilde{X}^\top \mathcal{S}(\text{MASK}(A^{(1)}))^\top \\ &= \bar{x}_T^\top A^{(2)} \bar{X}^\top \mathcal{S}(\text{MASK}(A^{(1)}))^\top. \end{aligned}$$

Thus the output of the second attention layer is

$$\begin{aligned} \text{attn}(h^{(1)}; \tilde{A}^{(2)})_T &= h^{(1)\top} \mathcal{S}\left(h^{(1)} \left(\tilde{A}^{(2)}\right)^\top h^{(T)}\right) \\ &= h^{(1)\top} \mathcal{S}\left(\mathcal{S}(\text{MASK}(A^{(1)})) \bar{X} A^{(2)\top} \bar{x}_T\right) \end{aligned}$$

Finally, the output is

$$\begin{aligned} \widetilde{\text{TF}}_{\tilde{\theta}}(s_{1:T}) &= \widetilde{W}_O^\top h_T^{(2)} \\ &= [I_S \quad 0_{S \times T} \mid 0_{S \times d}] \text{attn}(h^{(1)}; \tilde{A}^{(2)})_T \\ &= [I_S \quad 0_{S \times T} \mid 0_{S \times d}] h^{(1)\top} \mathcal{S}\left(\mathcal{S}(\text{MASK}(A^{(1)})) \bar{X} A^{(2)\top} \bar{x}_T\right) \\ &= \bar{X}^\top \mathcal{S}\left(\mathcal{S}(\text{MASK}(A^{(1)})) \bar{X} A^{(2)\top} \bar{x}_T\right), \end{aligned}$$

as desired.