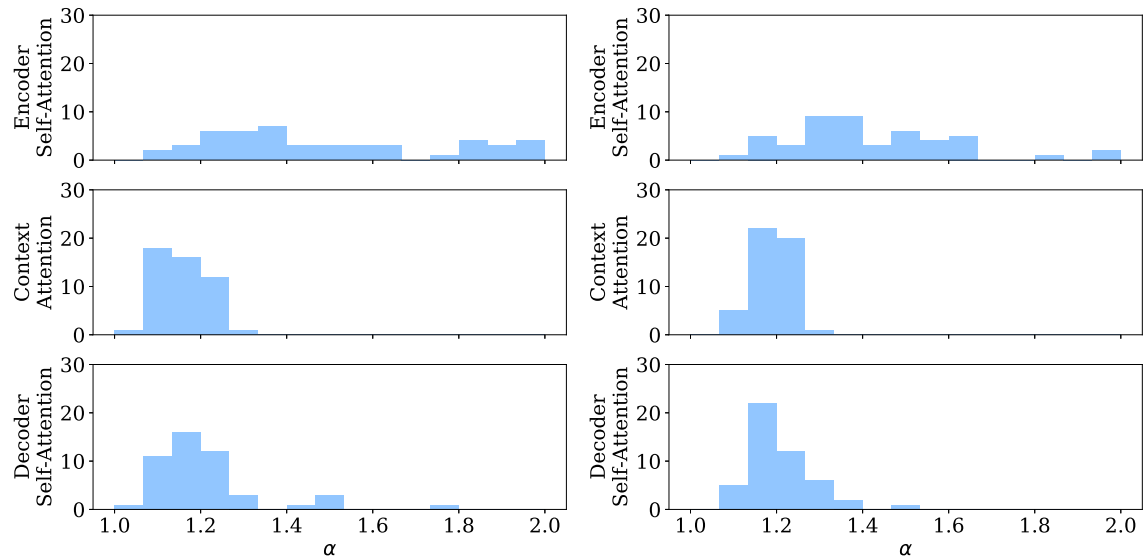


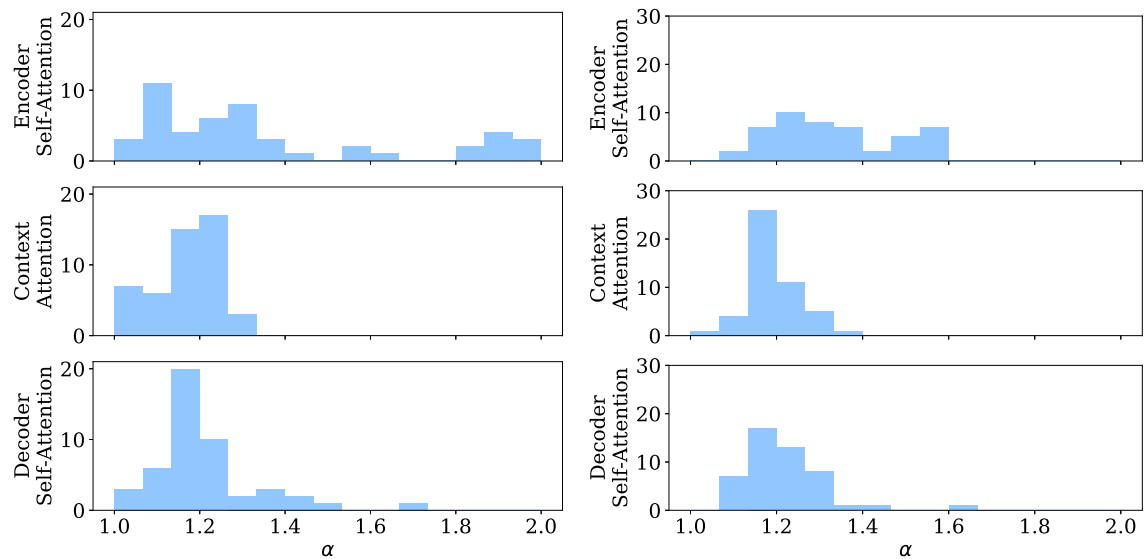
# Supplementary Material

## A High-Level Statistics Analysis of Other Language Pairs



(a) WMT 2016 RO  $\rightarrow$  EN.

(b) KFTT JA  $\rightarrow$  EN.



(c) WMT 2014 EN  $\rightarrow$  DE.

(d) IWSLT 2017 DE  $\rightarrow$  EN.

Figure 11: Histograms of  $\alpha$  values.

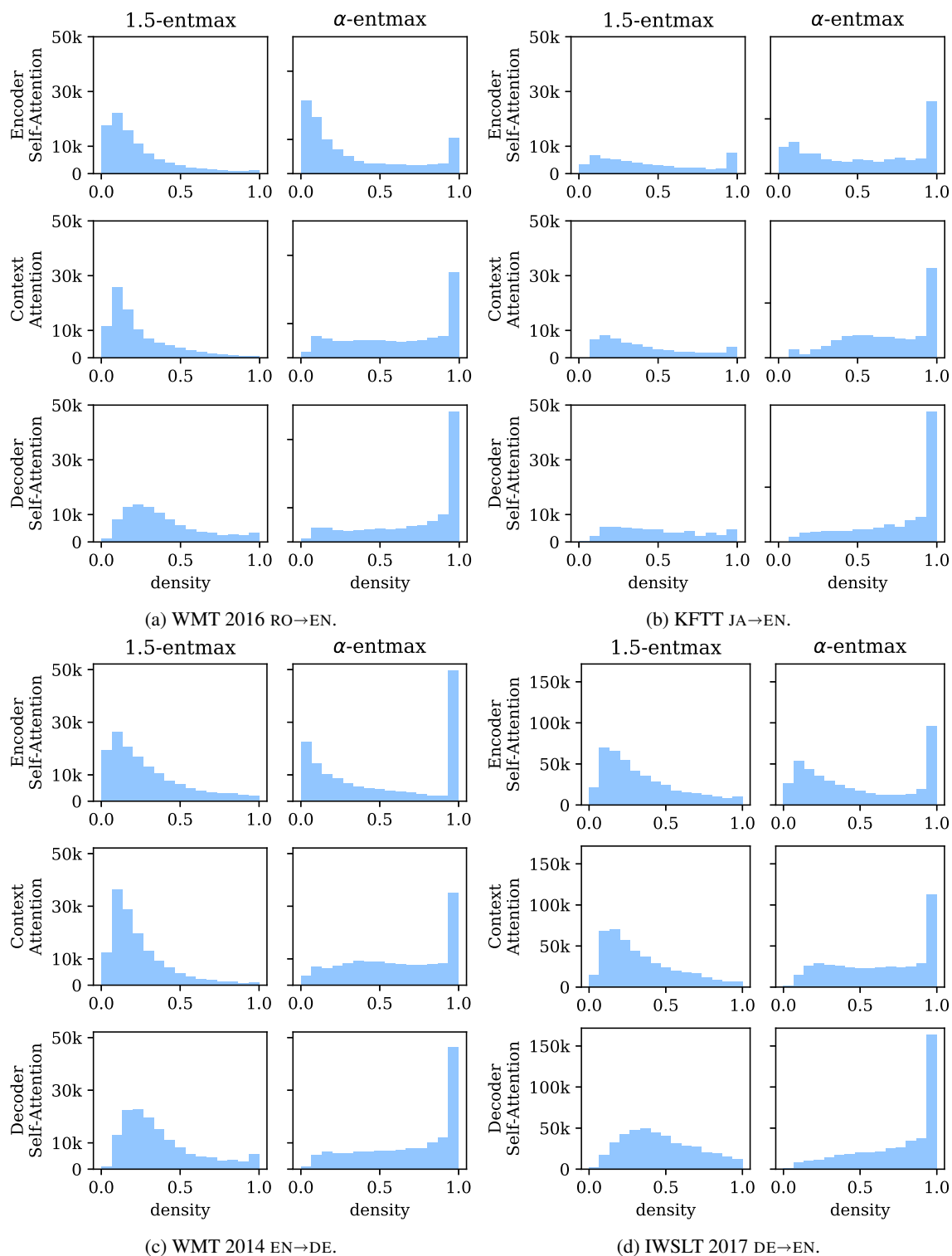
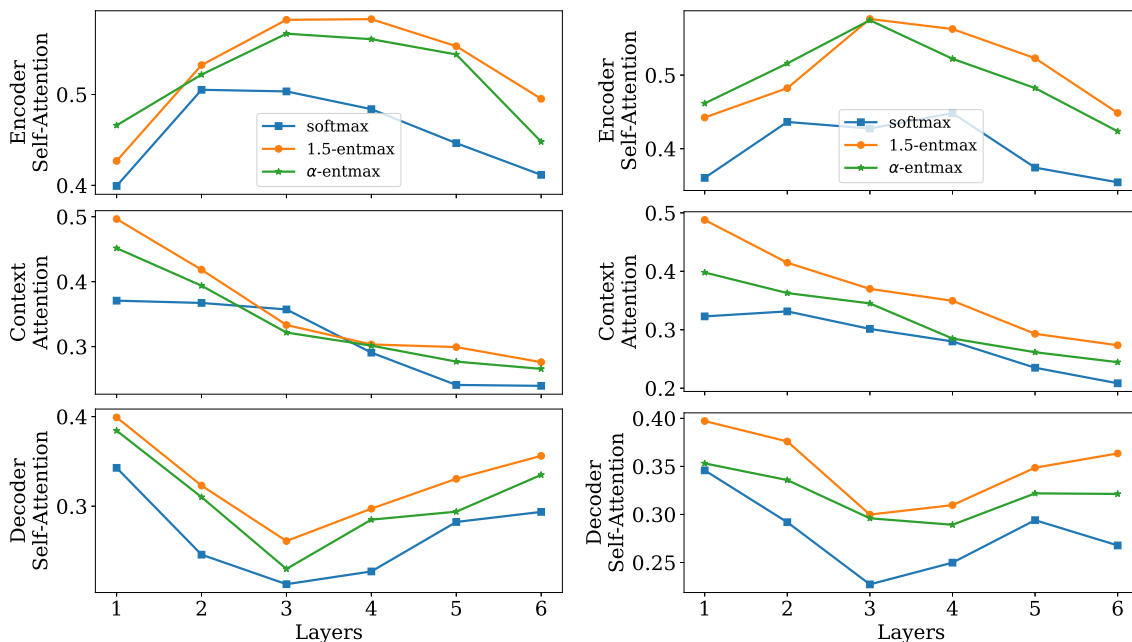
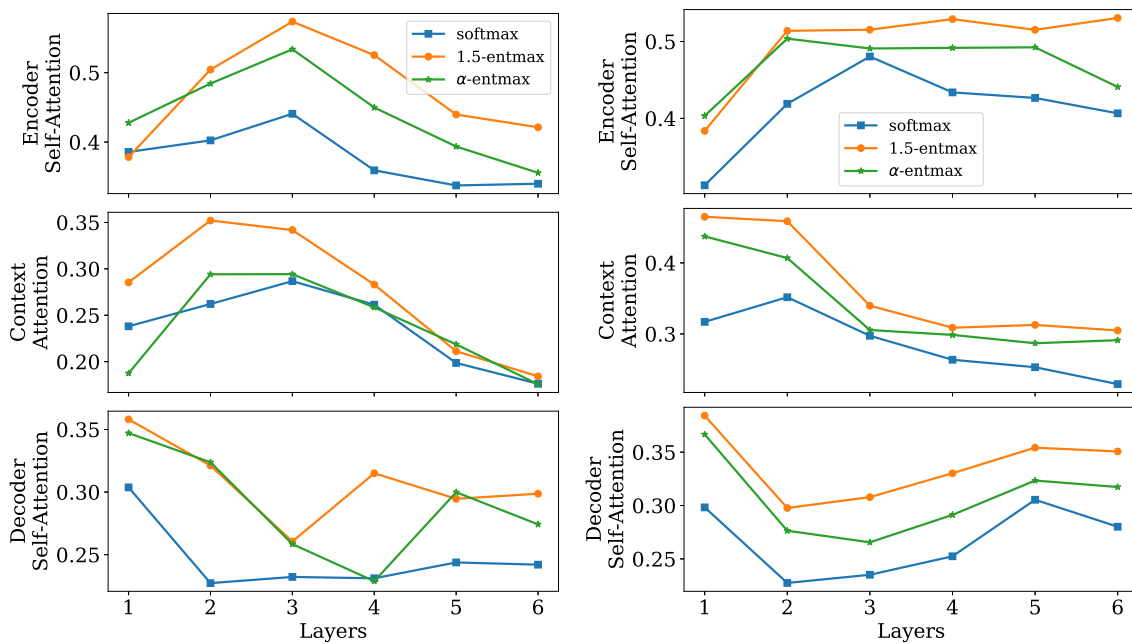


Figure 12: Histograms of head densities.



(a) WMT 2016 RO→EN.

(b) KFTT JA→EN.



(c) WMT 2014 EN→DE.

(d) IWSLT 2017 DE→EN.

Figure 13: Jensen-Shannon divergence over layers.

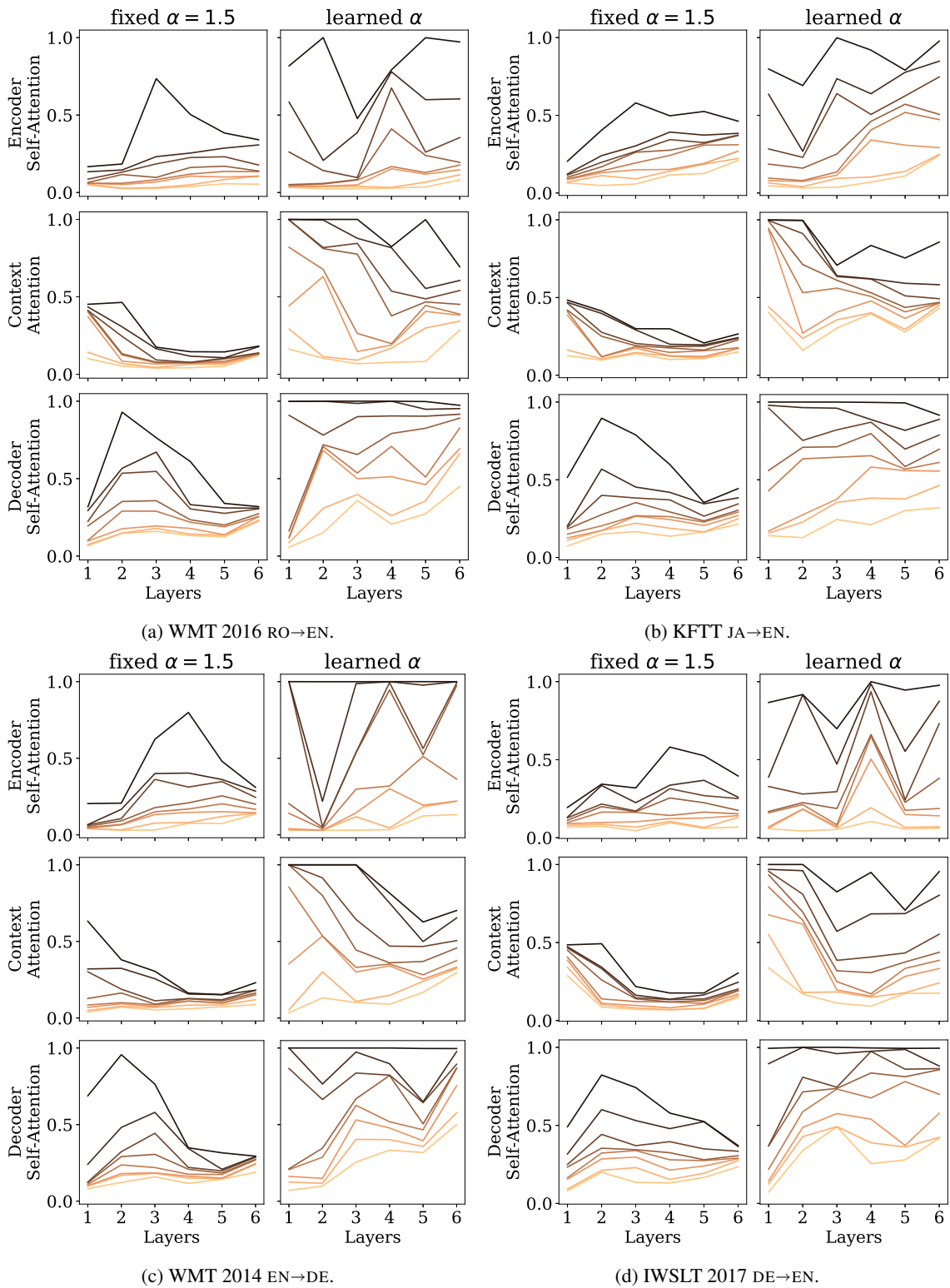


Figure 14: Head densities over layers.

## B Background

### B.1 Regularized Fenchel-Young prediction functions

**Definition 1** (Blondel et al. 2019). Let  $\Omega: \Delta^d \rightarrow \mathbb{R} \cup \{\infty\}$  be a strictly convex regularization function. We define the prediction function  $\pi_\Omega$  as

$$\pi_\Omega(\mathbf{z}) = \operatorname{argmax}_{\mathbf{p} \in \Delta^d} (\mathbf{p}^\top \mathbf{z} - \Omega(\mathbf{p})) \quad (12)$$

### B.2 Characterizing the $\alpha$ -entmax mapping

**Lemma 1** (Peters et al. 2019). For any  $\mathbf{z}$ , there exists a unique  $\tau^*$  such that

$$\alpha\text{-entmax}(\mathbf{z}) = [(\alpha - 1)\mathbf{z} - \tau^*\mathbf{1}]_+^{1/\alpha-1}. \quad (13)$$

*Proof:* From the definition of  $\alpha$ -entmax,

$$\alpha\text{-entmax}(\mathbf{z}) := \operatorname{argmax}_{\mathbf{p} \in \Delta^d} \mathbf{p}^\top \mathbf{z} + H_\alpha^\top(\mathbf{p}), \quad (14)$$

we may easily identify it with a regularized prediction function (Def. 1):

$$\alpha\text{-entmax}(\mathbf{z}) \equiv \pi_{-H_\alpha^\top}(\mathbf{z}).$$

We first note that for all  $\mathbf{p} \in \Delta^d$ ,

$$-(\alpha - 1)H_\alpha^\top(\mathbf{p}) = \frac{1}{\alpha} \sum_{i=1}^d p_i^\alpha + \text{const}. \quad (15)$$

From the constant invariance and scaling properties of  $\pi_\Omega$  (Blondel et al., 2019, Proposition 1, items 4–5),

$$\pi_{-H_\alpha^\top}(\mathbf{z}) = \pi_\Omega((\alpha - 1)\mathbf{z}), \quad \text{with} \quad \Omega(\mathbf{p}) = \sum_{j=1}^d g(p_j), \quad g(t) = \frac{t^\alpha}{\alpha}.$$

Using (Blondel et al., 2019, Proposition 5), noting that  $g'(t) = t^{\alpha-1}$  and  $(g')^{-1}(u) = u^{1/\alpha-1}$ , yields

$$\pi_\Omega(\mathbf{z}) = [\mathbf{z} - \tau^*\mathbf{1}]_+^{1/\alpha-1}, \quad \text{and therefore} \quad \alpha\text{-entmax}(\mathbf{z}) = [(\alpha - 1)\mathbf{z} - \tau^*\mathbf{1}]_+^{1/\alpha-1}. \quad (16)$$

Since  $H_\alpha^\top$  is strictly convex on the simplex,  $\alpha$ -entmax has a unique solution  $\mathbf{p}^*$ . Equation 16 implicitly defines a one-to-one mapping between  $\mathbf{p}^*$  and  $\tau^*$  as long as  $\mathbf{p}^* \in \Delta$ , therefore  $\tau^*$  is also unique. ■

### B.3 Connections to softmax and sparsemax

The Euclidean projection onto the simplex, sometimes referred to, in the context of neural attention, as sparsemax (Martins and Astudillo, 2016), is defined as

$$\text{sparsemax}(\mathbf{z}) := \operatorname{argmin}_{\mathbf{p} \in \Delta} \|\mathbf{p} - \mathbf{z}\|_2^2. \quad (17)$$

The solution can be characterized through the unique threshold  $\tau$  such that  $\sum_i \text{sparsemax}(\mathbf{z})_i = 1$  and (Held et al., 1974)

$$\text{sparsemax}(\mathbf{z}) = [\mathbf{z} - \tau\mathbf{1}]_+. \quad (18)$$

Thus, each coordinate of the sparsemax solution is a piecewise-linear function. Visibly, this expression is recovered when setting  $\alpha = 2$  in the  $\alpha$ -entmax expression (Equation 21); for other values of  $\alpha$ , the exponent induces curvature.

On the other hand, the well-known softmax is usually defined through the expression

$$\text{softmax}(\mathbf{z})_i := \frac{\exp(z_i)}{\sum_j \exp(z_j)}, \quad (19)$$

which can be shown to be the unique solution of the optimization problem

$$\text{softmax}(\mathbf{z})_i = \operatorname{argmax}_{\mathbf{p} \in \Delta} \mathbf{p}^\top \mathbf{z} + \mathbf{H}^S(\mathbf{p}), \quad (20)$$

where  $\mathbf{H}^S(\mathbf{p}) := -\sum_i p_i \log p_i$  is the Shannon entropy. Indeed, setting the gradient to 0 yields the condition  $\log p_i = z_i - \nu_i - \tau - 1$ , where  $\tau$  and  $\nu_i > 0$  are Lagrange multipliers for the simplex constraints  $\sum_i p_i = 1$  and  $p_i \geq 0$ , respectively. Since the *l.h.s.* is only finite for  $p_i > 0$ , we must have  $\nu_i = 0$  for all  $i$ , by complementary slackness. Thus, the solution must have the form  $p_i = \exp(z_i)/Z$ , yielding Equation 19.

## C Jacobian of $\alpha$ -entmax *w.r.t.* the shape parameter $\alpha$ : Proof of Proposition 1

Recall that the entmax transformation is defined as:

$$\alpha\text{-entmax}(\mathbf{z}) := \operatorname{argmax}_{\mathbf{p} \in \Delta^d} \mathbf{p}^\top \mathbf{z} + \mathbf{H}_\alpha^T(\mathbf{p}), \quad (21)$$

where  $\alpha \geq 1$  and  $\mathbf{H}_\alpha^T$  is the Tsallis entropy,

$$\mathbf{H}_\alpha^T(\mathbf{p}) := \begin{cases} \frac{1}{\alpha(\alpha-1)} \sum_j (p_j - p_j^\alpha), & \alpha \neq 1, \\ \mathbf{H}^S(\mathbf{p}), & \alpha = 1, \end{cases} \quad (22)$$

and  $\mathbf{H}^S(\mathbf{p}) := -\sum_j p_j \log p_j$  is the Shannon entropy.

In this section, we derive the Jacobian of entmax with respect to the scalar parameter  $\alpha$ .

### C.1 General case of $\alpha > 1$

From the KKT conditions associated with the optimization problem in Eq. 21, we have that the solution  $\mathbf{p}^*$  has the following form, coordinate-wise:

$$p_i^* = [(\alpha - 1)(z_i - \tau^*)]_+^{1/(\alpha-1)}, \quad (23)$$

where  $\tau^*$  is a scalar Lagrange multiplier that ensures that  $\mathbf{p}^*$  normalizes to 1, *i.e.*, it is defined implicitly by the condition:

$$\sum_i [(\alpha - 1)(z_i - \tau^*)]_+^{1/(\alpha-1)} = 1. \quad (24)$$

For general values of  $\alpha$ , Eq. 24 lacks a closed form solution. This makes the computation of the Jacobian

$$\frac{\partial \alpha\text{-entmax}(\mathbf{z})}{\partial \alpha} \quad (25)$$

non-trivial. Fortunately, we can use the technique of implicit differentiation to obtain this Jacobian.

The Jacobian exists almost everywhere, and the expressions we derive yield a generalized Jacobian (Clarke, 1990) at any non-differentiable points that may occur for certain  $(\alpha, \mathbf{z})$  pairs. We begin

by noting that  $\frac{\partial p_i^*}{\partial \alpha} = 0$  if  $p_i^* = 0$ , because increasing  $\alpha$  keeps sparse coordinates sparse.<sup>4</sup> Therefore we need to worry only about coordinates that are in the support of  $\mathbf{p}^*$ . We will assume hereafter that the  $i^{\text{th}}$  coordinate of  $\mathbf{p}^*$  is non-zero. We have:

$$\begin{aligned}
\frac{\partial p_i^*}{\partial \alpha} &= \frac{\partial}{\partial \alpha} [(\alpha - 1)(z_i - \tau^*)]^{\frac{1}{\alpha-1}} \\
&= \frac{\partial}{\partial \alpha} \exp \left[ \frac{1}{\alpha-1} \log [(\alpha - 1)(z_i - \tau^*)] \right] \\
&= p_i^* \frac{\partial}{\partial \alpha} \left[ \frac{1}{\alpha-1} \log [(\alpha - 1)(z_i - \tau^*)] \right] \\
&= \frac{p_i^*}{(\alpha - 1)^2} \left[ \frac{\frac{\partial}{\partial \alpha} [(\alpha - 1)(z_i - \tau^*)]}{z_i - \tau^*} - \log [(\alpha - 1)(z_i - \tau^*)] \right] \\
&= \frac{p_i^*}{(\alpha - 1)^2} \left[ \frac{z_i - \tau^* - (\alpha - 1) \frac{\partial \tau^*}{\partial \alpha}}{z_i - \tau^*} - \log [(\alpha - 1)(z_i - \tau^*)] \right] \\
&= \frac{p_i^*}{(\alpha - 1)^2} \left[ 1 - \frac{\alpha - 1}{z_i - \tau^*} \frac{\partial \tau^*}{\partial \alpha} - \log [(\alpha - 1)(z_i - \tau^*)] \right]. \tag{26}
\end{aligned}$$

We can see that this Jacobian depends on  $\frac{\partial \tau^*}{\partial \alpha}$ , which we now compute using implicit differentiation.

Let  $\mathcal{S} = \{i : p_i^* > 0\}$ . By differentiating both sides of Eq. 24, re-using some of the steps in Eq. 26, and recalling Eq. 23, we get

$$\begin{aligned}
0 &= \sum_{i \in \mathcal{S}} \frac{\partial}{\partial \alpha} [(\alpha - 1)(z_i - \tau^*)]^{1/(\alpha-1)} \\
&= \sum_{i \in \mathcal{S}} \frac{p_i^*}{(\alpha - 1)^2} \left[ 1 - \frac{\alpha - 1}{z_i - \tau^*} \frac{\partial \tau^*}{\partial \alpha} - \log [(\alpha - 1)(z_i - \tau^*)] \right] \\
&= \frac{1}{(\alpha - 1)^2} - \frac{\partial \tau^*}{\partial \alpha} \sum_{i \in \mathcal{S}} \frac{p_i^*}{(\alpha - 1)(z_i - \tau^*)} - \sum_{i \in \mathcal{S}} \frac{p_i^*}{(\alpha - 1)^2} \log [(\alpha - 1)(z_i - \tau^*)] \\
&= \frac{1}{(\alpha - 1)^2} - \frac{\partial \tau^*}{\partial \alpha} \sum_i (p_i^*)^{2-\alpha} - \sum_i \frac{p_i^*}{\alpha - 1} \log p_i^* \\
&= \frac{1}{(\alpha - 1)^2} - \frac{\partial \tau^*}{\partial \alpha} \sum_i (p_i^*)^{2-\alpha} + \frac{\mathbf{H}^{\mathcal{S}}(\mathbf{p}^*)}{\alpha - 1}, \tag{27}
\end{aligned}$$

from which we obtain:

$$\frac{\partial \tau^*}{\partial \alpha} = \frac{\frac{1}{(\alpha-1)^2} + \frac{\mathbf{H}^{\mathcal{S}}(\mathbf{p}^*)}{\alpha-1}}{\sum_i (p_i^*)^{2-\alpha}}. \tag{28}$$

Finally, plugging Eq. 28 into Eq. 26, we get:

$$\begin{aligned}
\frac{\partial p_i^*}{\partial \alpha} &= \frac{p_i^*}{(\alpha - 1)^2} \left[ 1 - \frac{1}{(p_i^*)^{\alpha-1}} \frac{\partial \tau^*}{\partial \alpha} - (\alpha - 1) \log p_i^* \right] \\
&= \frac{p_i^*}{(\alpha - 1)^2} \left[ 1 - \frac{1}{(p_i^*)^{\alpha-1}} \frac{\frac{1}{(\alpha-1)^2} + \frac{\mathbf{H}^{\mathcal{S}}(\mathbf{p}^*)}{\alpha-1}}{\sum_i (p_i^*)^{2-\alpha}} - (\alpha - 1) \log p_i^* \right] \\
&= \frac{p_i^* - \tilde{p}_i(\alpha)}{(\alpha - 1)^2} - \frac{p_i^* \log p_i^* + \tilde{p}_i(\alpha) \mathbf{H}^{\mathcal{S}}(\mathbf{p}^*)}{\alpha - 1}, \tag{29}
\end{aligned}$$

<sup>4</sup>This follows from the margin property of  $\mathbf{H}_\alpha^{\mathcal{T}}$  (Blondel et al., 2019).

where we denote by

$$\tilde{p}_i(\alpha) = \frac{(p_i^*)^{2-\alpha}}{\sum_j (p_j^*)^{2-\alpha}}. \quad (30)$$

The distribution  $\tilde{\mathbf{p}}(\alpha)$  can be interpreted as a “skewed” distribution obtained from  $\mathbf{p}^*$ , which appears in the Jacobian of  $\alpha$ -entmax( $\mathbf{z}$ ) *w.r.t.*  $\mathbf{z}$  as well (Peters et al., 2019).

## C.2 Solving the indetermination for $\alpha = 1$

We can write Eq. 29 as

$$\frac{\partial p_i^*}{\partial \alpha} = \frac{p_i^* - \tilde{p}_i(\alpha) - (\alpha - 1)(p_i^* \log p_i^* + \tilde{p}_i(\alpha) \mathbf{H}^S(\mathbf{p}^*))}{(\alpha - 1)^2}. \quad (31)$$

When  $\alpha \rightarrow 1^+$ , we have  $\tilde{\mathbf{p}}(\alpha) \rightarrow \mathbf{p}^*$ , which leads to a  $\frac{0}{0}$  indetermination.

To solve this indetermination, we will need to apply L’Hôpital’s rule twice. Let us first compute the derivative of  $\tilde{p}_i(\alpha)$  with respect to  $\alpha$ . We have

$$\frac{\partial}{\partial \alpha} (p_i^*)^{2-\alpha} = -(p_i^*)^{2-\alpha} \log p_i^*, \quad (32)$$

therefore

$$\begin{aligned} \frac{\partial}{\partial \alpha} \tilde{p}_i(\alpha) &= \frac{\partial}{\partial \alpha} \frac{(p_i^*)^{2-\alpha}}{\sum_j (p_j^*)^{2-\alpha}} \\ &= \frac{-(p_i^*)^{2-\alpha} \log p_i^* \sum_j (p_j^*)^{2-\alpha} + (p_i^*)^{2-\alpha} \sum_j (p_j^*)^{2-\alpha} \log p_j^*}{\left(\sum_j (p_j^*)^{2-\alpha}\right)^2} \\ &= -\tilde{p}_i(\alpha) \log p_i^* + \tilde{p}_i(\alpha) \sum_j \tilde{p}_j(\alpha) \log p_j^*. \end{aligned} \quad (33)$$

Differentiating the numerator and denominator in Eq. 31, we get:

$$\begin{aligned} \frac{\partial p_i^*}{\partial \alpha} &= \lim_{\alpha \rightarrow 1^+} \frac{(1 + (\alpha - 1) \mathbf{H}^S(\mathbf{p}^*)) \tilde{p}_i(\alpha) (\log p_i^* - \sum_j \tilde{p}_j(\alpha) \log p_j^*) - p_i^* \log p_i^* - \tilde{p}_i(\alpha) \mathbf{H}^S(\mathbf{p}^*)}{2(\alpha - 1)} \\ &= A + B, \end{aligned} \quad (34)$$

with

$$\begin{aligned} A &= \lim_{\alpha \rightarrow 1^+} \frac{\mathbf{H}^S(\mathbf{p}^*) \tilde{p}_i(\alpha) (\log p_i^* - \sum_j \tilde{p}_j(\alpha) \log p_j^*) \mathbf{H}^S(\mathbf{p}^*)}{2} \\ &= \frac{\mathbf{H}^S(\mathbf{p}^*) p_i^* \log p_i^* + p_i^* (\mathbf{H}^S(\mathbf{p}^*))^2}{2}, \end{aligned} \quad (35)$$

and

$$B = \lim_{\alpha \rightarrow 1^+} \frac{\tilde{p}_i(\alpha) (\log p_i^* - \sum_j \tilde{p}_j(\alpha) \log p_j^*) - p_i^* \log p_i^* - \tilde{p}_i(\alpha) \mathbf{H}^S(\mathbf{p}^*)}{2(\alpha - 1)}. \quad (36)$$

When  $\alpha \rightarrow 1^+$ ,  $B$  becomes again a  $\frac{0}{0}$  indetermination, which we can solve by applying again L’Hôpital’s



rule. Differentiating the numerator and denominator in Eq. 36:

$$\begin{aligned}
B &= \frac{1}{2} \lim_{\alpha \rightarrow 1^+} \left\{ \tilde{p}_i(\alpha) \log p_i^* \left( \sum_j \tilde{p}_j(\alpha) \log p_j^* - \log p_i^* \right) \right. \\
&\quad \left. - \tilde{p}_i(\alpha) \left( \sum_j \tilde{p}_j(\alpha) \log p_j^* - \log p_i^* \right) \left( \sum_j \tilde{p}_j(\alpha) \log p_j^* + H^S(\mathbf{p}^*) \right) \right. \\
&\quad \left. - \tilde{p}_i(\alpha) \sum_j \tilde{p}_j(\alpha) \log p_j^* \left( \sum_k \tilde{p}_k(\alpha) \log p_k^* - \log p_j^* \right) \right\} \\
&= \frac{-p_i^* \log p_i^* (H^S(\mathbf{p}^*) + \log p_i^*) + p_i^* \sum_j p_j^* \log p_j^* (H^S(\mathbf{p}^*) + \log p_j^*)}{2} \\
&= \frac{-H^S(\mathbf{p}^*) p_i^* \log p_i^* - p_i^* (H^S(\mathbf{p}^*))^2 - p_i^* \log^2 p_i^* + p_i^* \sum_j p_j^* \log^2 p_j^*}{2}. \tag{37}
\end{aligned}$$

Finally, summing Eq. 35 and Eq. 37, we get

$$\left. \frac{\partial p_i^*}{\partial \alpha} \right|_{\alpha=1} = \frac{-p_i^* \log^2 p_i^* + p_i^* \sum_j p_j^* \log^2 p_j^*}{2}. \tag{38}$$

### C.3 Summary

To sum up, we have the following expression for the Jacobian of  $\alpha$ -entmax with respect to  $\alpha$ :

$$\frac{\partial p_i^*}{\partial \alpha} = \begin{cases} \frac{p_i^* - \tilde{p}_i(\alpha)}{(\alpha-1)^2} - \frac{p_i^* \log p_i^* + \tilde{p}_i(\alpha) H^S(\mathbf{p}^*)}{\alpha-1}, & \text{for } \alpha > 1 \\ \frac{-p_i^* \log^2 p_i^* + p_i^* \sum_j p_j^* \log^2 p_j^*}{2}, & \text{for } \alpha = 1. \end{cases} \tag{39}$$