

# Adaptively Sparse Transformers

**Gonçalo Correia** Instituto de Telecomunicações, Lisbon

Vlad Niculae IT

André Martins IT & Unbabel

# Introduction

Transformers are amazing at achieving top of the leaderboard results!

# Introduction

Transformers are amazing at achieving top of the leaderboard results!

But they seem overparameterized...

# Introduction

Transformers are amazing at achieving top of the leaderboard results!

But they seem overparameterized...

Attention heads aid visualization but they are completely **dense**.

# Introduction

Transformers are amazing at achieving top of the leaderboard results!

But they seem overparameterized...

Attention heads aid visualization but they are completely **dense**.

Our solution is to bet on **sparsity**:

- for interpretability
- for discovering linguistic structure
- for efficiency

# Introduction

Transformers are amazing at achieving top of the leaderboard results!

But they seem overparameterized...

Attention heads aid visualization but they are completely **dense**.

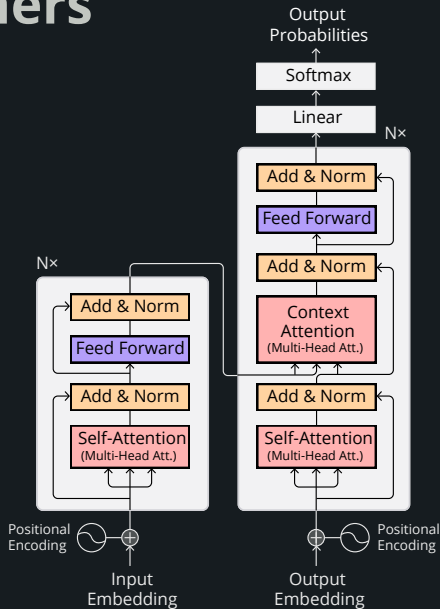
Our solution is to bet on **sparsity**:

- **for interpretability**
- **for discovering linguistic structure**
- for efficiency

# Transformers

In each attention head:

$$\bar{V} = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V.$$



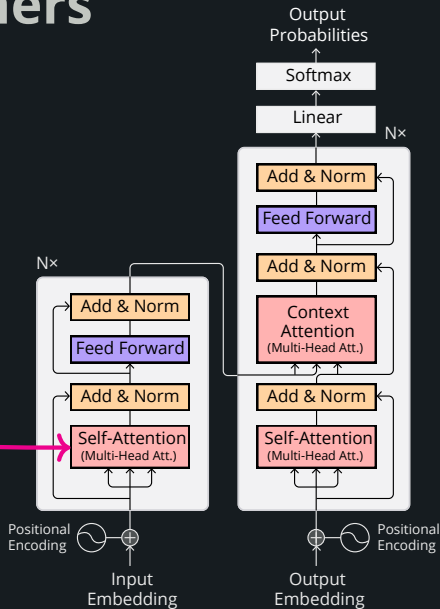
# Transformers

In each attention head:

$$\bar{V} = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V.$$

Attention in three places:

- Self-attention in the encoder





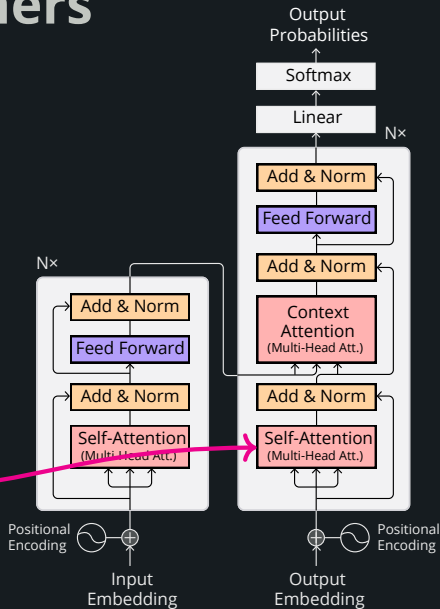
# Transformers

In each attention head:

$$\bar{\mathbf{V}} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}.$$

Attention in three places:

- Self-attention in the encoder
- Self-attention in the decoder



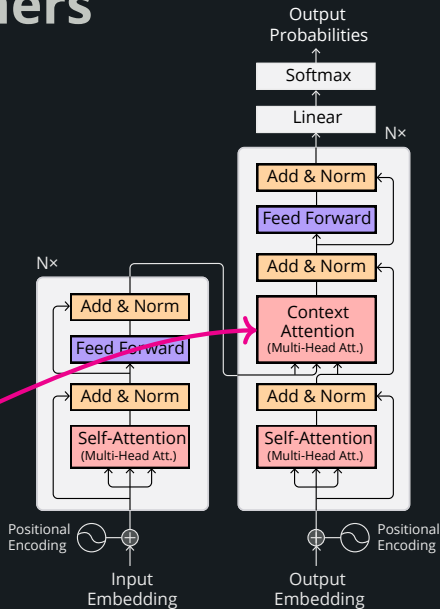
# Transformers

In each attention head:

$$\bar{\mathbf{V}} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}.$$

Attention in three places:

- Self-attention in the encoder
- Self-attention in the decoder
- Contextual attention



# Sparse Transformers

# Sparse Transformers

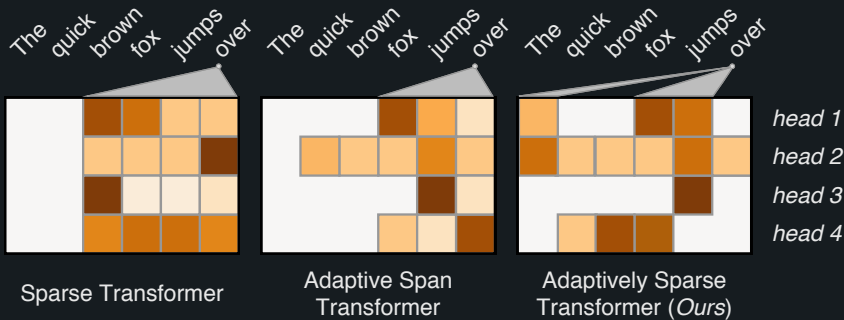
**Key idea:** replace softmax in attention heads by a sparse normalizing function! 🙌

# Adaptively Sparse Transformers

**Key idea:** replace softmax in attention heads by a sparse normalizing function! 🙌

**Another key idea:** use a normalizing function that is adaptively sparse via a learnable  $\alpha$ ! 🙌 🙌 🙌

# Related Work: Other Sparse Transformers



Our model allows **non-contiguous** attention for each head.

# What is softmax?

Softmax exponentiates and normalizes:  $\text{softmax}(z_i) := \frac{\exp(z_i)}{\sum_j \exp(z_j)}$

# What is softmax?

Softmax exponentiates and normalizes:  $\text{softmax}(z_i) := \frac{\exp(z_i)}{\sum_j \exp(z_j)}$

It's fully dense:  $\text{softmax}(z) > 0$



# What is softmax?

Softmax exponentiates and normalizes:  $\text{softmax}(z_i) := \frac{\exp(z_i)}{\sum_j \exp(z_j)}$

It's fully dense:  $\text{softmax}(\mathbf{z}) > \mathbf{0}$

Argmax can be written as:

$$\text{argmax}(\mathbf{z}) := \arg \max_{\mathbf{p} \in \Delta} \mathbf{z}^T \mathbf{p}$$

# What is softmax?

Softmax exponentiates and normalizes:  $\text{softmax}(z_i) := \frac{\exp(z_i)}{\sum_j \exp(z_j)}$

It's fully dense:  $\text{softmax}(\mathbf{z}) > \mathbf{0}$

Argmax can be written as:

$$\text{argmax}(\mathbf{z}) := \arg \max_{\mathbf{p} \in \Delta} \mathbf{z}^T \mathbf{p}$$

- Retrieves a **one-hot vector** for the highest scored index.

# What is softmax?

Softmax exponentiates and normalizes:  $\text{softmax}(z_i) := \frac{\exp(z_i)}{\sum_j \exp(z_j)}$

It's fully dense:  $\text{softmax}(\mathbf{z}) > \mathbf{0}$

Argmax can be written as:

$$\text{argmax}(\mathbf{z}) := \arg \max_{\mathbf{p} \in \Delta} \mathbf{z}^T \mathbf{p}$$

- Retrieves a **one-hot vector** for the highest scored index.
- Sometimes used as hard attention, but not differentiable!

# $\Omega$ -Regularized Argmax

For convex  $\Omega$ , define the  $\Omega$ -regularized argmax transformation:

$$\mathbf{argmax}_{\Omega}(\mathbf{z}) := \arg \max_{\mathbf{p} \in \Delta} \mathbf{z}^{\top} \mathbf{p} - \Omega(\mathbf{p})$$

# $\Omega$ -Regularized Argmax

For convex  $\Omega$ , define the  $\Omega$ -regularized argmax transformation:

$$\mathbf{argmax}_{\Omega}(\mathbf{z}) := \arg \max_{\mathbf{p} \in \Delta} \mathbf{z}^{\top} \mathbf{p} - \Omega(\mathbf{p})$$

- **Argmax** corresponds to no regularization,  $\Omega \equiv 0$

# $\Omega$ -Regularized Argmax

For convex  $\Omega$ , define the  $\Omega$ -regularized argmax transformation:

$$\mathbf{argmax}_{\Omega}(\mathbf{z}) := \arg \max_{\mathbf{p} \in \Delta} \mathbf{z}^{\top} \mathbf{p} - \Omega(\mathbf{p})$$

- **Argmax** corresponds to **no regularization**,  $\Omega \equiv 0$
- **Softmax** amounts to **entropic regularization**,  $\Omega(\mathbf{p}) = \sum_{i=1}^K p_i \log p_i$

# $\Omega$ -Regularized Argmax

For convex  $\Omega$ , define the  $\Omega$ -regularized argmax transformation:

$$\mathbf{argmax}_{\Omega}(\mathbf{z}) := \arg \max_{\mathbf{p} \in \Delta} \mathbf{z}^{\top} \mathbf{p} - \Omega(\mathbf{p})$$

- **Argmax** corresponds to **no regularization**,  $\Omega \equiv 0$
- **Softmax** amounts to **entropic regularization**,  $\Omega(\mathbf{p}) = \sum_{i=1}^K p_i \log p_i$
- **Sparsemax** amounts to  $\ell_2$ -regularization,  $\Omega(\mathbf{p}) = \frac{1}{2} \|\mathbf{p}\|^2$ .

# $\Omega$ -Regularized Argmax

For convex  $\Omega$ , define the  $\Omega$ -regularized argmax transformation:

$$\mathbf{argmax}_{\Omega}(\mathbf{z}) := \arg \max_{\mathbf{p} \in \Delta} \mathbf{z}^{\top} \mathbf{p} - \Omega(\mathbf{p})$$

- **Argmax** corresponds to **no regularization**,  $\Omega \equiv 0$
- **Softmax** amounts to **entropic regularization**,  $\Omega(\mathbf{p}) = \sum_{i=1}^K p_i \log p_i$
- **Sparsemax** amounts to  $\ell_2$ -regularization,  $\Omega(\mathbf{p}) = \frac{1}{2} \|\mathbf{p}\|^2$ .

Is there something in-between?



# $\alpha$ -Entmax

Parametrized by  $\alpha \geq 0$ :

$$\Omega_{\alpha}(\mathbf{p}) := \begin{cases} \frac{1}{\alpha(\alpha-1)} \left( 1 - \sum_{i=1}^K p_i^{\alpha} \right) & \text{if } \alpha \neq 1 \\ \sum_{i=1}^K p_i \log p_i & \text{if } \alpha = 1. \end{cases}$$

# $\alpha$ -Entmax

Parametrized by  $\alpha \geq 0$ :

$$\Omega_{\alpha}(\mathbf{p}) := \begin{cases} \frac{1}{\alpha(\alpha-1)} \left( 1 - \sum_{i=1}^K p_i^{\alpha} \right) & \text{if } \alpha \neq 1 \\ \sum_{i=1}^K p_i \log p_i & \text{if } \alpha = 1. \end{cases}$$

- **Argmax** corresponds to  $\alpha \rightarrow \infty$

# $\alpha$ -Entmax

Parametrized by  $\alpha \geq 0$ :

$$\Omega_{\alpha}(\mathbf{p}) := \begin{cases} \frac{1}{\alpha(\alpha-1)} \left( 1 - \sum_{i=1}^K p_i^{\alpha} \right) & \text{if } \alpha \neq 1 \\ \sum_{i=1}^K p_i \log p_i & \text{if } \alpha = 1. \end{cases}$$

- **Argmax** corresponds to  $\alpha \rightarrow \infty$
- **Softmax** amounts to  $\alpha \rightarrow 1$

# $\alpha$ -Entmax

Parametrized by  $\alpha \geq 0$ :

$$\Omega_{\alpha}(\mathbf{p}) := \begin{cases} \frac{1}{\alpha(\alpha-1)} \left( 1 - \sum_{i=1}^K p_i^{\alpha} \right) & \text{if } \alpha \neq 1 \\ \sum_{i=1}^K p_i \log p_i & \text{if } \alpha = 1. \end{cases}$$

- **Argmax** corresponds to  $\alpha \rightarrow \infty$
- **Softmax** amounts to  $\alpha \rightarrow 1$
- **Sparsemax** amounts to  $\alpha = 2$ .

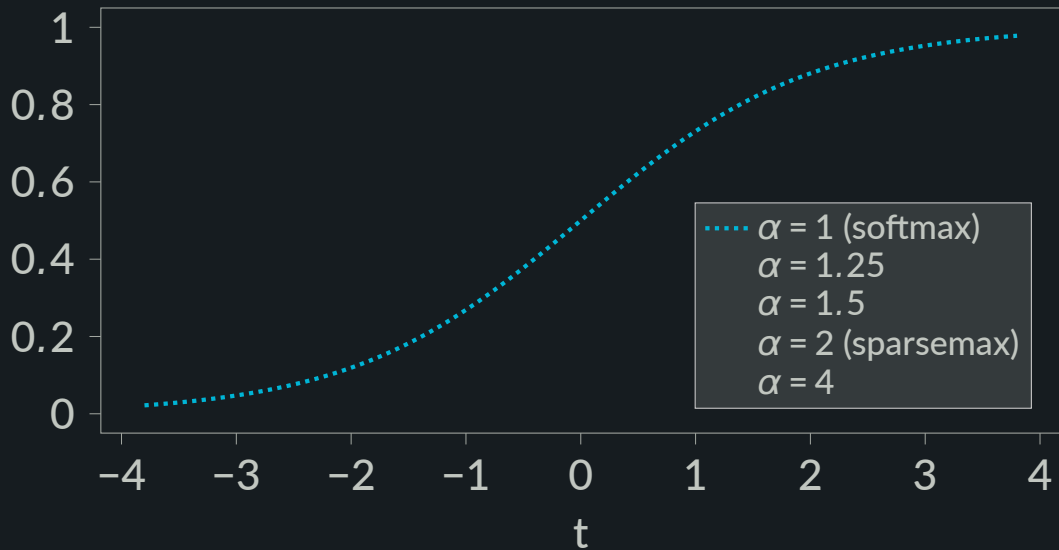
# $\alpha$ -Entmax

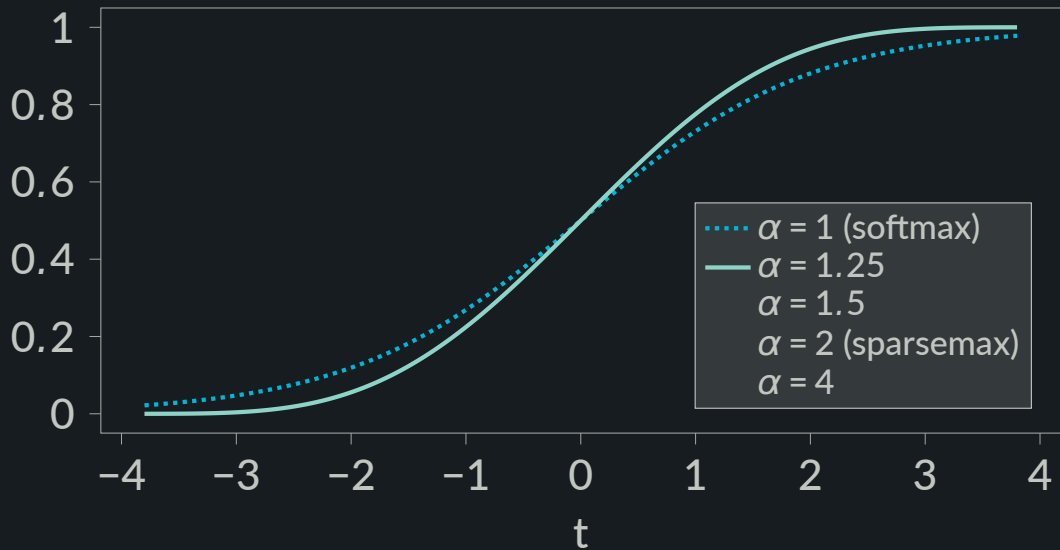
Parametrized by  $\alpha \geq 0$ :

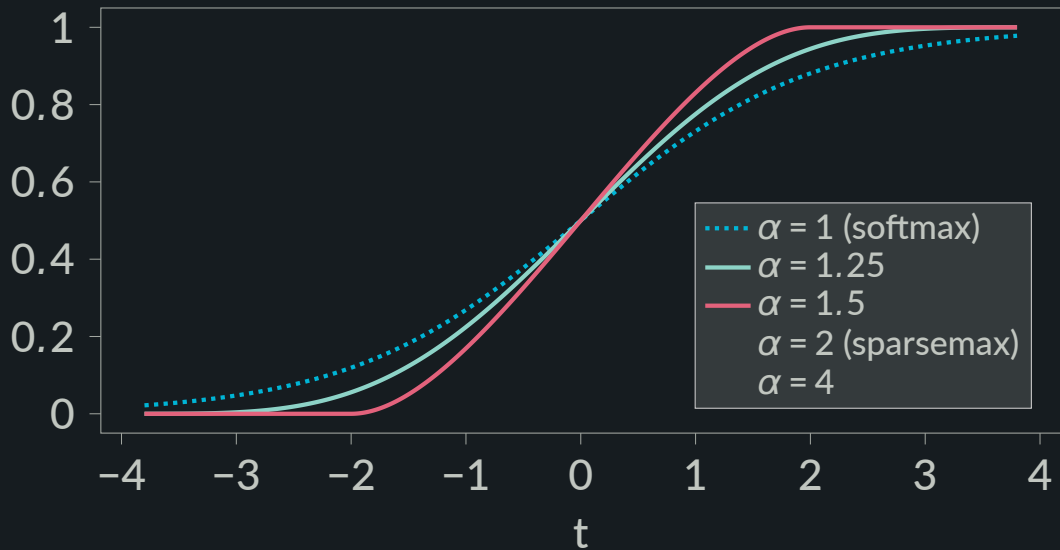
$$\Omega_{\alpha}(\mathbf{p}) := \begin{cases} \frac{1}{\alpha(\alpha-1)} \left( 1 - \sum_{i=1}^K p_i^{\alpha} \right) & \text{if } \alpha \neq 1 \\ \sum_{i=1}^K p_i \log p_i & \text{if } \alpha = 1. \end{cases}$$

- **Argmax** corresponds to  $\alpha \rightarrow \infty$
- **Softmax** amounts to  $\alpha \rightarrow 1$
- **Sparsemax** amounts to  $\alpha = 2$ .

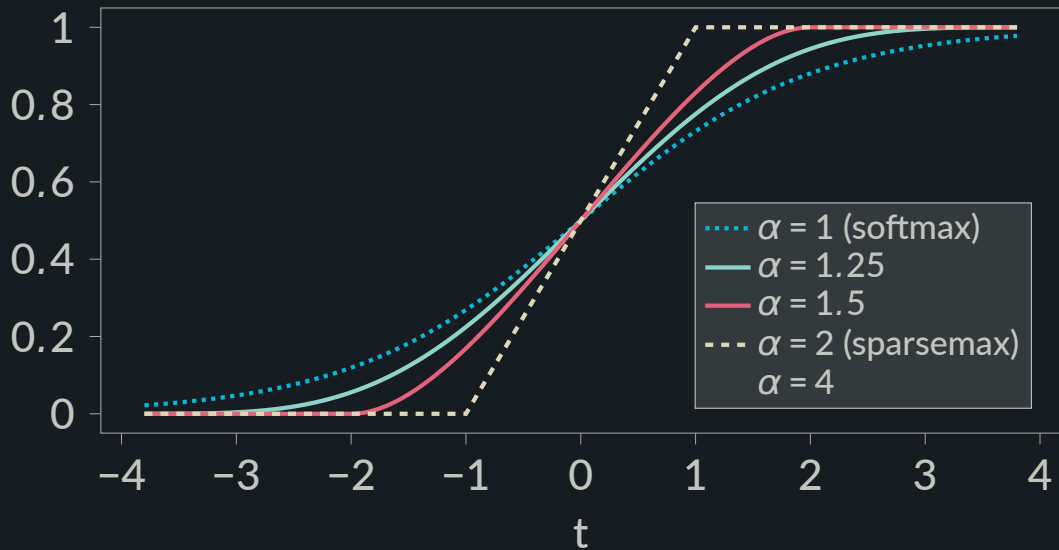
**Key result:** can be sparse for  $\alpha > 1$ , propensity for sparsity increases with  $\alpha$ .

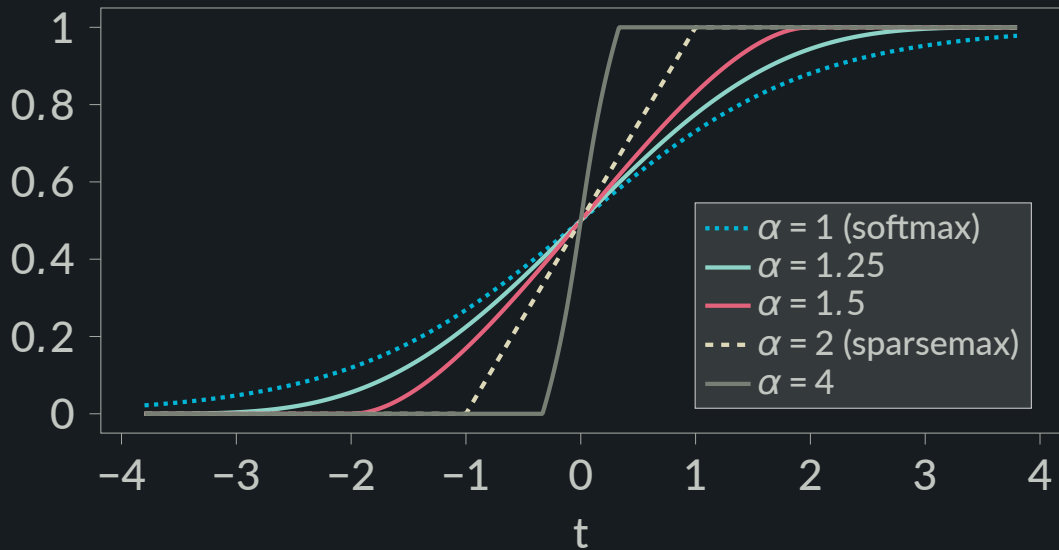












# Learning $\alpha$

# Learning $\alpha$

Key contribution:

a closed-form expression for  $\frac{\partial \alpha\text{-entmax}(\mathbf{z})}{\partial \alpha}$  🙌

# Learning $\alpha$

Key contribution:

a closed-form expression for  $\frac{\partial \alpha\text{-entmax}(\mathbf{z})}{\partial \alpha}$  🙌

Requires argmin differentiation → see paper for details!

# Learning $\alpha$

Key commands:

```
:pip install entmax
```

a check [Check github.com/deep-spin/entmax](https://github.com/deep-spin/entmax)

Requires argmin differentiation → see paper for details!

# BLEU Scores

activation	de→en	ja→en	ro→en	en→de
softmax	29.79	21.57	32.70	26.02
1.5-entmax	29.83	22.13	33.10	25.89
$\alpha$ -entmax	29.90	21.74	32.89	26.93

# BLEU Scores

activation	de→en	ja→en	ro→en	en→de
softmax	29.79	21.57	32.70	26.02
1.5-entmax	29.83	22.13	33.10	25.89
$\alpha$ -entmax	29.90	21.74	32.89	26.93

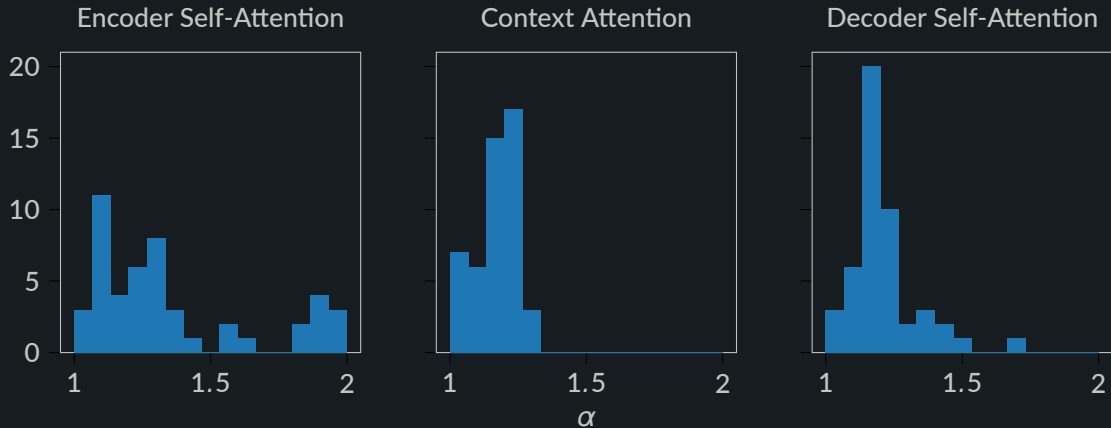


# BLEU Scores

activation	de→en	ja→en	ro→en	en→de
softmax	29.79	21.57	32.70	26.02
1.5-entmax	29.83	22.13	33.10	25.89
$\alpha$ -entmax	29.90	21.74	32.89	26.93

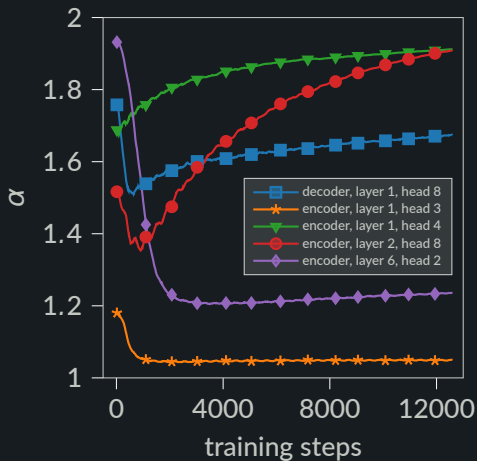
For analysis for other language pairs, see Appendix A.

# Learned $\alpha$

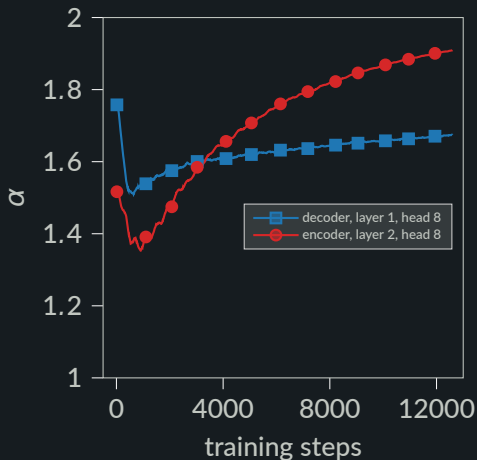


Bimodal for the encoder, mostly unimodal for the decoder.

# Trajectories of $\alpha$ During Training

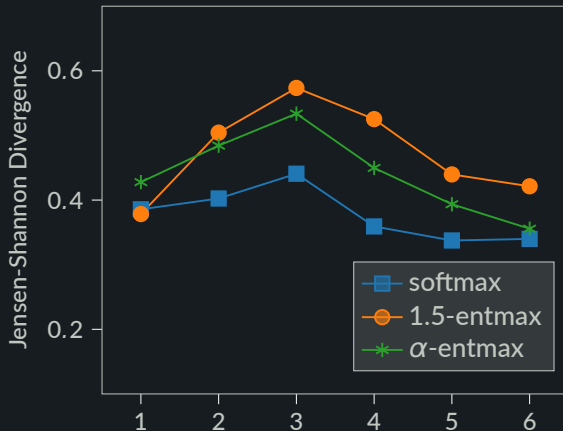


# Trajectories of $\alpha$ During Training

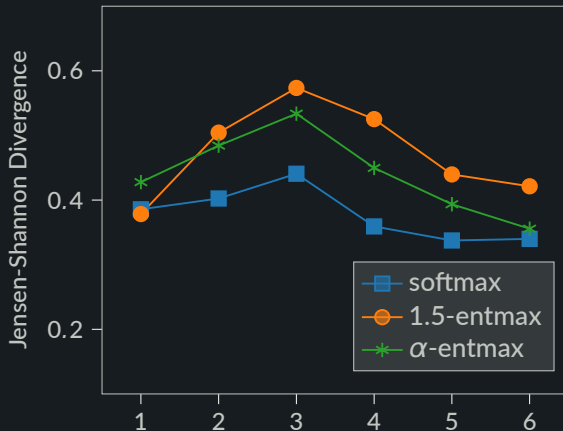


Some heads choose to start dense before becoming sparse.

# Head Diversity per Layer

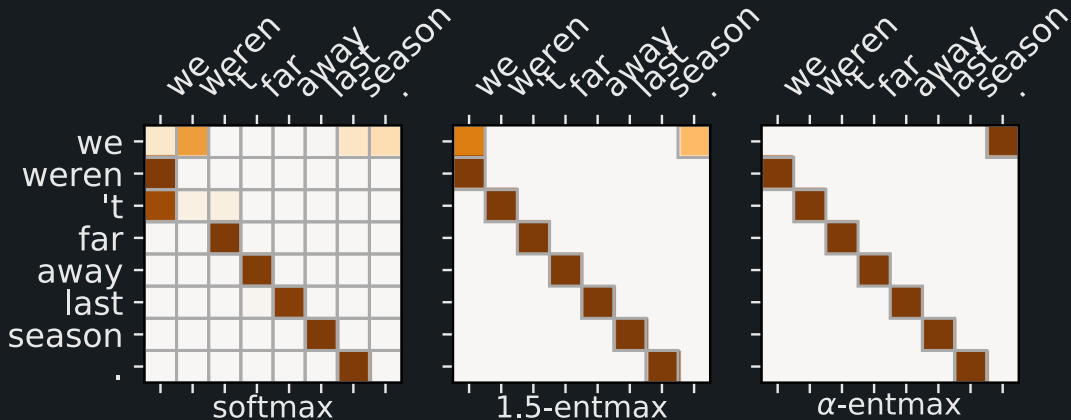


# Head Diversity per Layer



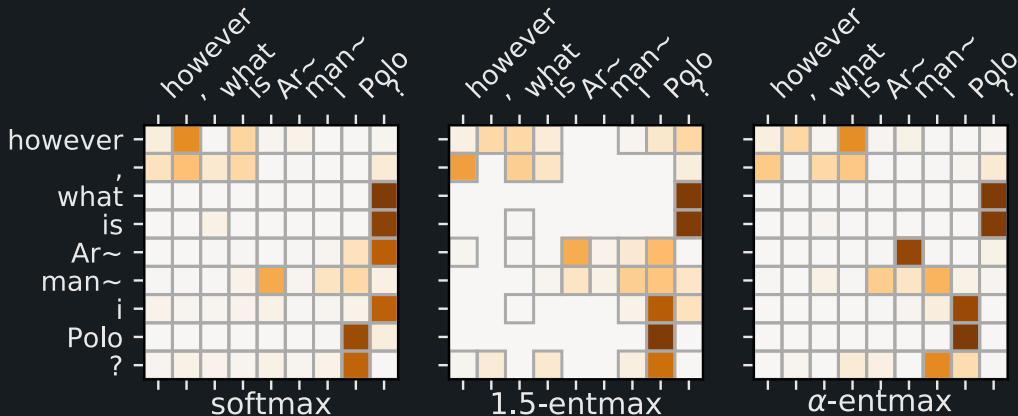
Specialized heads are important as seen in Voita et al. (2019)!

# Previous Position Head



This head role was also found in Voita et al. (2019)! Learned  $\alpha = 1.91$ .

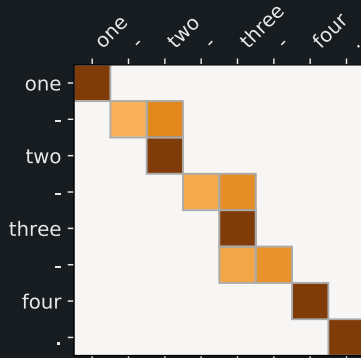
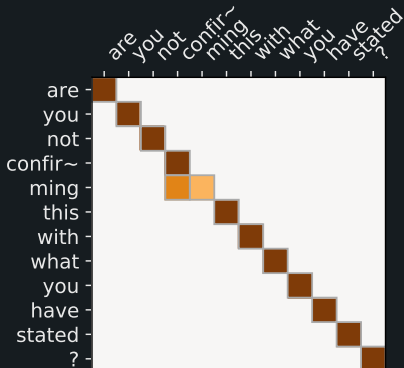
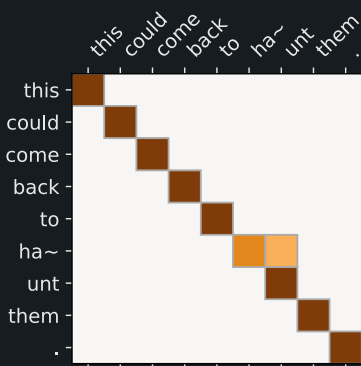
# Interrogation-Detecting Head



Learned  $\alpha = 1.05$ .



# Subword-Merging Head



Learned  $\alpha = 1.91$ .

# Key Takeaways

Introduce **adaptive** sparsity  
for Transformers via  $\alpha$ -entmax with a **gradient learnable  $\alpha$** .

# Key Takeaways

Introduce **adaptive** sparsity  
for Transformers via  $\alpha$ -entmax with a **gradient learnable**  $\alpha$ .

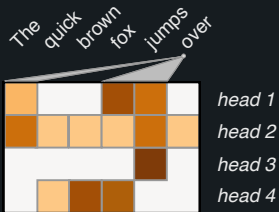
*adaptive sparsity*



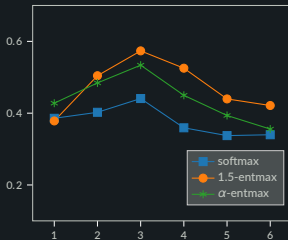
# Key Takeaways

Introduce **adaptive** sparsity  
for Transformers via  $\alpha$ -entmax with a **gradient learnable**  $\alpha$ .

*adaptive sparsity*



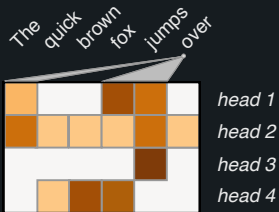
*reduced head redundancy*



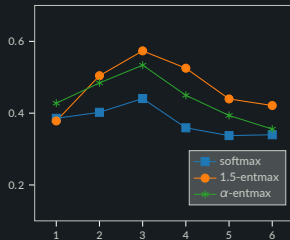
# Key Takeaways

Introduce **adaptive** sparsity  
for Transformers via  $\alpha$ -entmax with a **gradient learnable**  $\alpha$ .

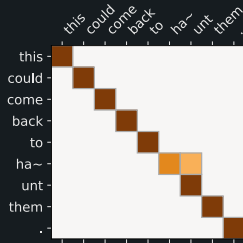
*adaptive sparsity*



*reduced head redundancy*



*clearer head roles*



# Thank you!

## Questions?

```
:pip install entmax
```

Check [github.com/deep-spin/entmax](https://github.com/deep-spin/entmax)

# Acknowledgements



This work was supported by the European Research Council (ERC StG DeepSPIN 758969) and by the Fundação para a Ciência e Tecnologia through contract UID/EEA/50008/2019 and CMUPERI/TIC/0046/2014 (GoLocal).

# References I

- Child, Rewon, Scott Gray, Alec Radford, and Ilya Sutskever (2019). "Generating long sequences with sparse transformers". In: *arXiv preprint arXiv:1904.10509*.
- Martins, André FT and Ramón Fernandez Astudillo (2016). "From softmax to sparsemax: A sparse model of attention and multi-label classification". In: *Proc. of ICML*.
- Niculae, Vlad and Mathieu Blondel (2017). "A Regularized Framework for Sparse and Structured Neural Attention". In: *arXiv preprint arXiv:1705.07704*.
- Peters, Ben, Vlad Niculae, and André F. T. Martins (2019). "Sparse Sequence-to-Sequence Models". In: *Proceedings of the Annual Meeting of the Association for Computational Linguistics*.
- Sukhbaatar, Sainbayar, Edouard Grave, Piotr Bojanowski, and Armand Joulin (2019). "Adaptive Attention Span in Transformers". In: *arXiv preprint arXiv:1905.07799*.
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin (2017). "Attention Is All You Need". In: *Proc. of NeurIPS*.
- Voita, Elena, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov (2019). "Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned". In: *Proc. ACL*.