

Adaptively Sparse Transformers

Gonçalo Correia Instituto de Telecomunicações, Lisbon

Vlad Niculae IT

André Martins IT & Unbabel

TL;DL (Too Long; Didn't Listen)

We replace the softmax function in Transformers with α -entmax, creating a model that is able to learn its own sparsity in each attention head.

A bit of context... On Seq2Seq

United Nations elections end today

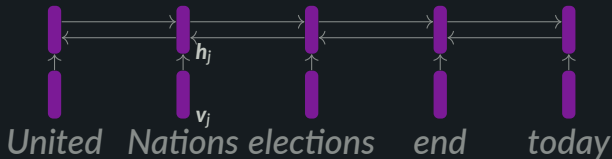
A bit of context... On Seq2Seq

Encoder

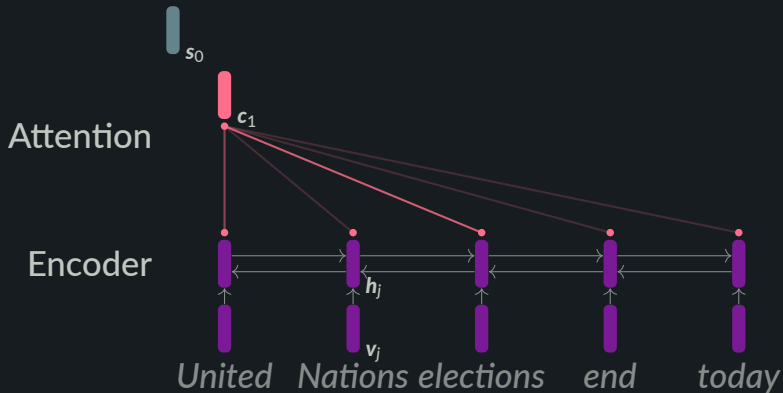

United Nations elections end today

A bit of context... On Seq2Seq

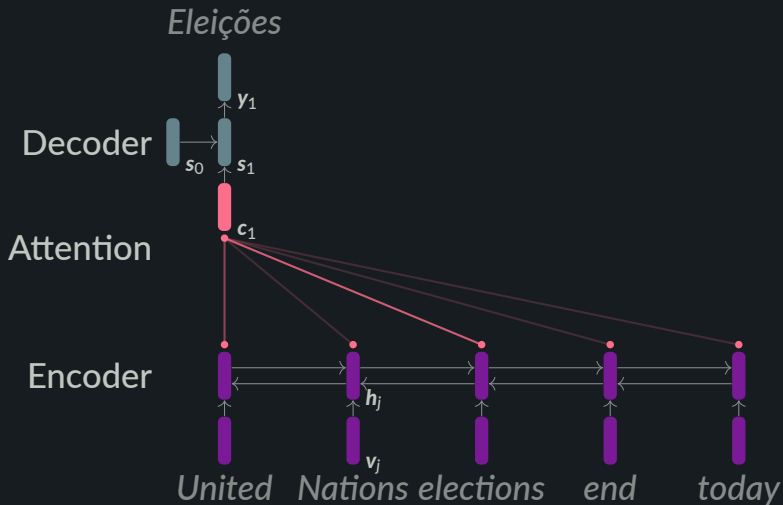
Encoder



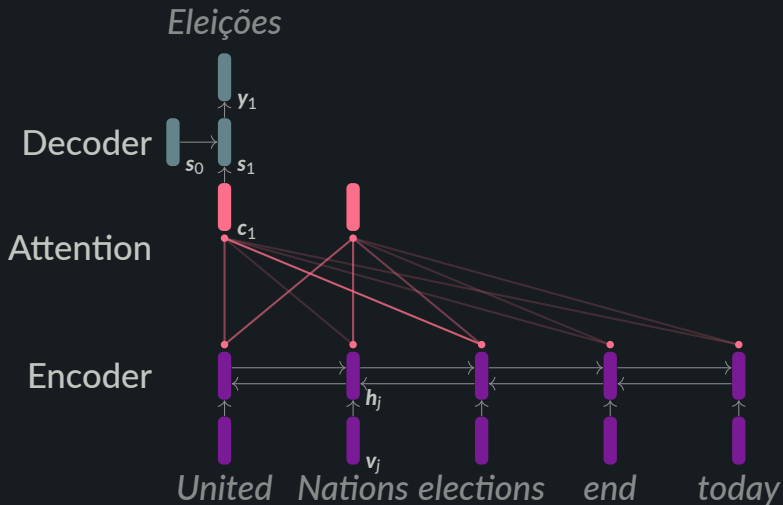
A bit of context... On Seq2Seq



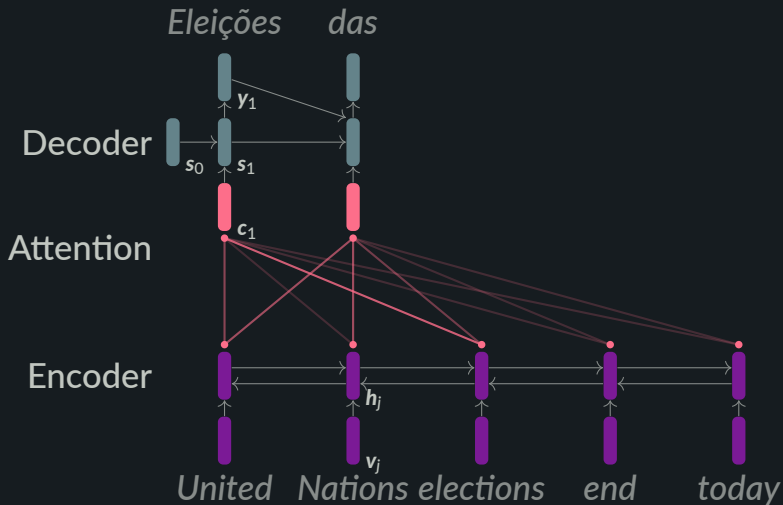
A bit of context... On Seq2Seq



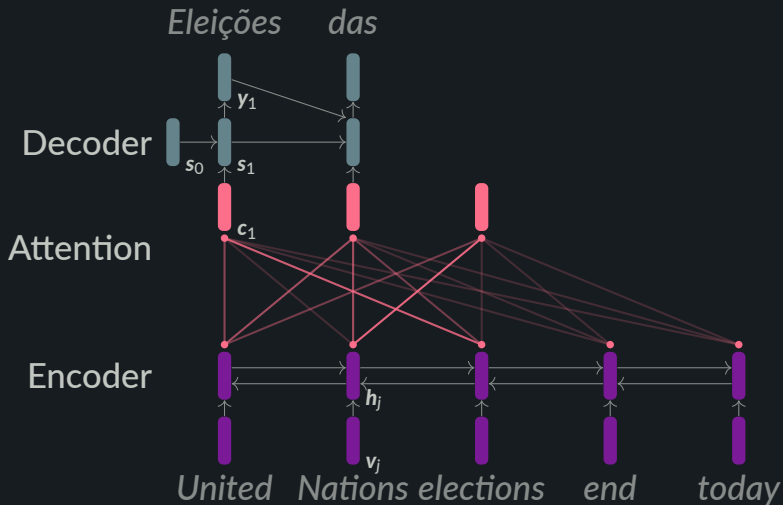
A bit of context... On Seq2Seq



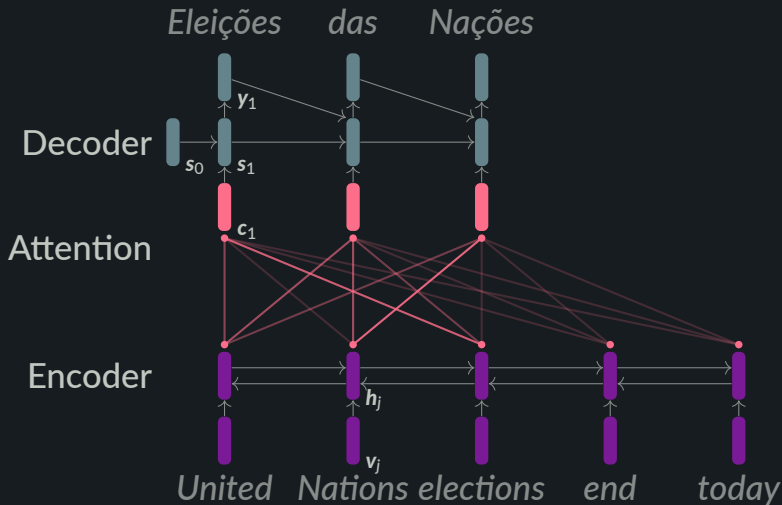
A bit of context... On Seq2Seq



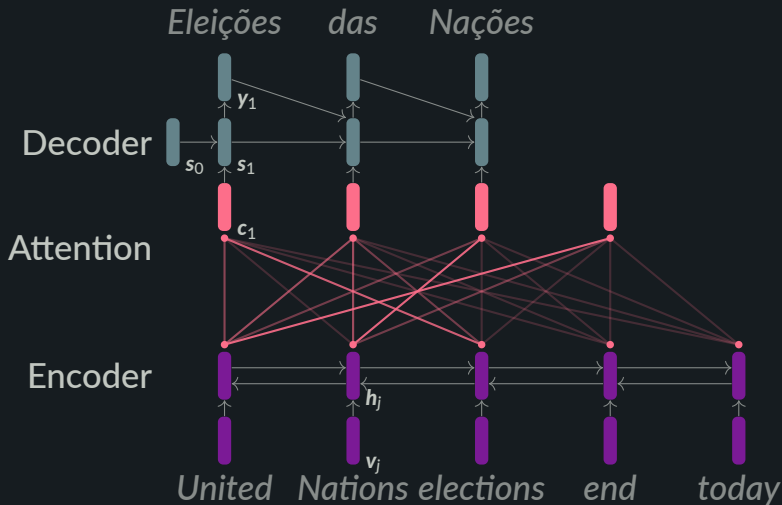
A bit of context... On Seq2Seq



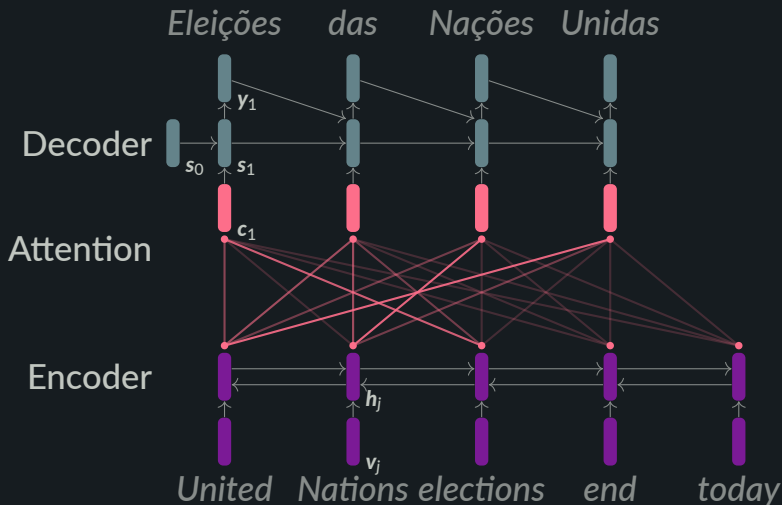
A bit of context... On Seq2Seq



A bit of context... On Seq2Seq



A bit of context... On Seq2Seq



attention weights
computed with
softmax:

for some decoder state s_t ,
compute contextually
weighted average of input c_t :

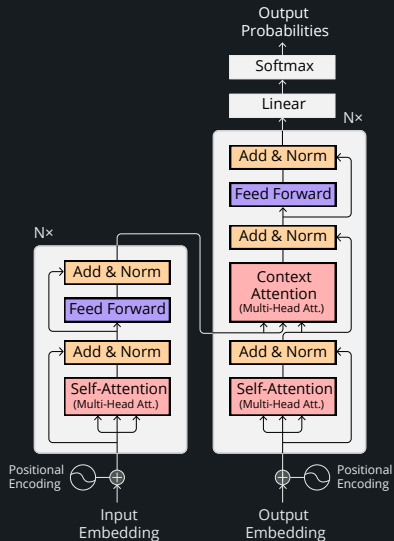
$$z_j = s_t^T \mathbf{W}^{(a)} h_j$$

$$\pi_j = \text{softmax}_j(\mathbf{z})$$

$$c_t = \sum_j \pi_j h_j$$

A bit of context... on Transformers

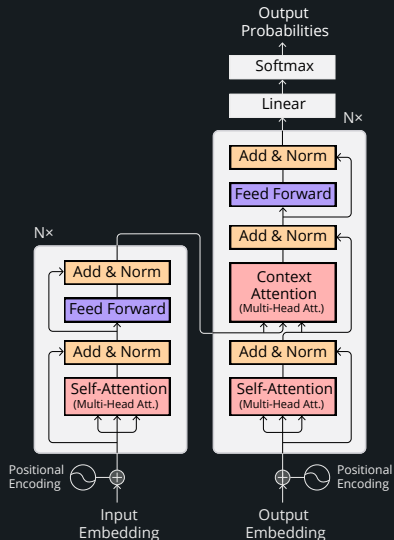
What if... Attention is all you need?



A bit of context... on Transformers

What if... Attention is all you need?

Key idea: Instead of Recurrent Neural Networks (RNNs), let's use attention mechanisms!

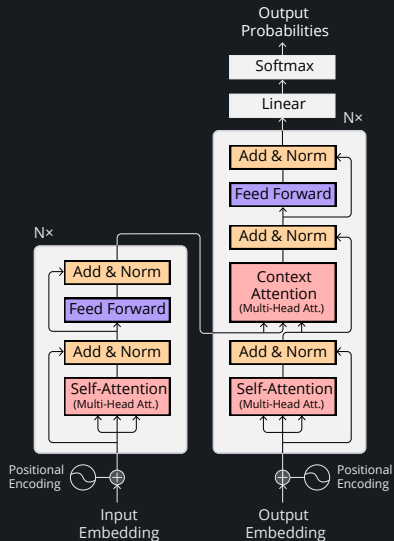


A bit of context... on Transformers

What if... Attention is all you need?

Key idea: Instead of Recurrent Neural Networks (RNNs), let's use attention mechanisms!

- In place of the RNNs, use self-attention

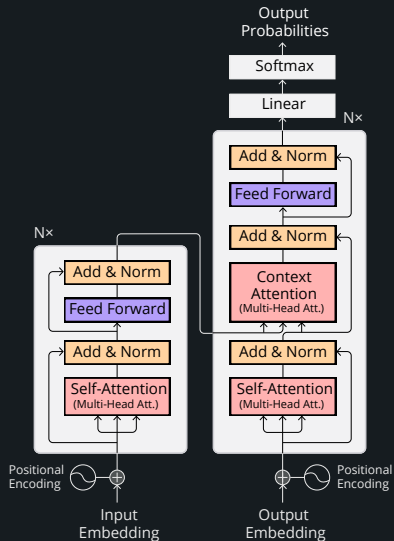


A bit of context... on Transformers

What if... Attention is all you need?

Key idea: Instead of Recurrent Neural Networks (RNNs), let's use attention mechanisms!

- In place of the RNNs, use self-attention
- Do this with multiple heads (i.e. attention mechanisms in parallel)

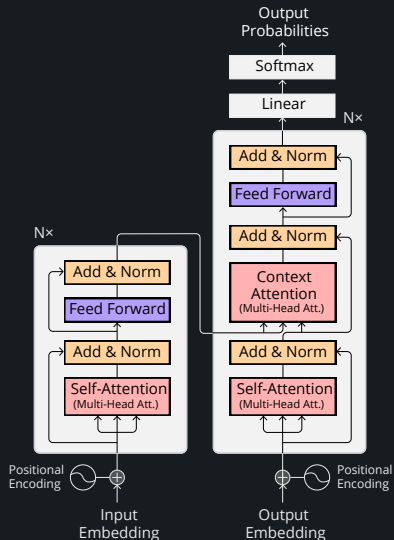


A bit of context... on Transformers

What if... Attention is all you need?

Key idea: Instead of Recurrent Neural Networks (RNNs), let's use attention mechanisms!

- In place of the RNNs, use self-attention
- Do this with multiple heads (i.e. attention mechanisms in parallel)
- ... and do it through several layers



Introduction

Transformers have been achieving incredible SOTA results in the last couple of years!

Introduction

Transformers have been achieving incredible SOTA results in the last couple of years!

But they seem overparameterized...

Introduction

Transformers have been achieving incredible SOTA results in the last couple of years!

But they seem overparameterized...

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

Introduction

Transformers have been achieving incredible SOTA results in the last couple of years!

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 have been achieving incredible SOTA results in the last couple of years!

But they seem overparameterized...

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

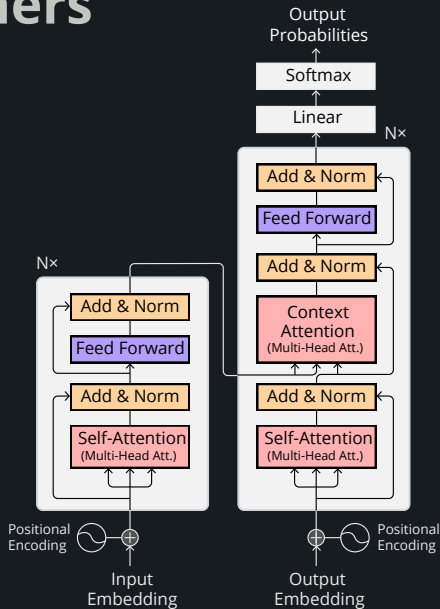
Our solution is to bet on **sparsity**:

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

Transformers

In each attention head:

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



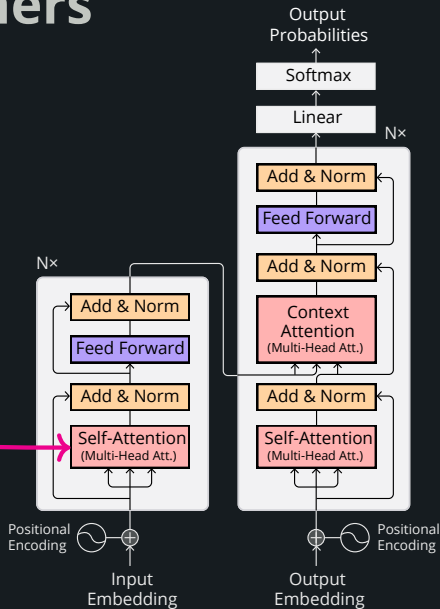
Transformers

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



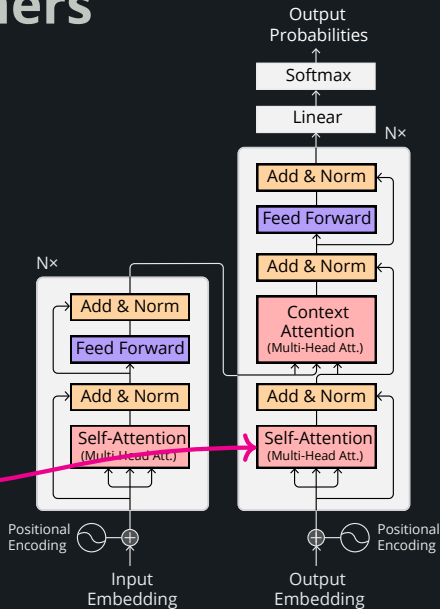
Transformers

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



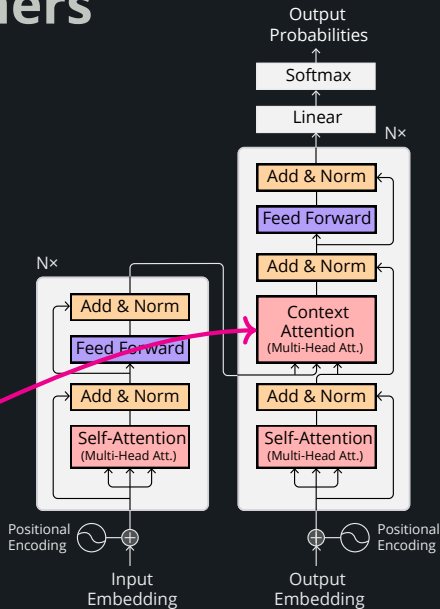
Transformers

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

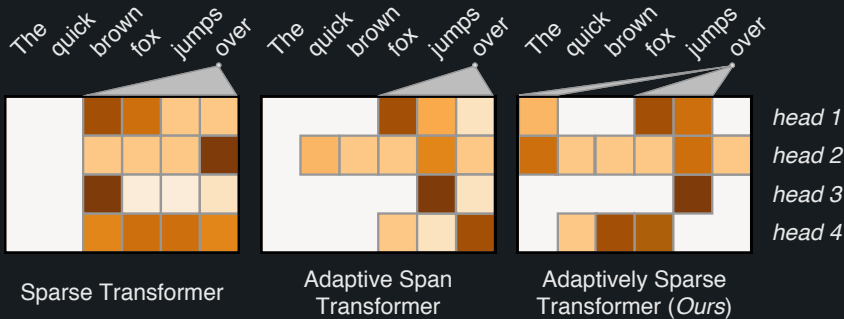
Key idea: replace softmax in attention heads by a sparse normalizing function! 🙌

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 α ! 🙌 🙌 🙌

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:

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

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: $\text{softmax}(\mathbf{z}) > 0$

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!

Ω -Regularized Argmax

For convex Ω , define the Ω -regularized argmax transformation:

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

Ω -Regularized Argmax

For convex Ω , define the Ω -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, $\Omega \equiv 0$

Ω -Regularized Argmax

For convex Ω , define the Ω -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**, $\Omega(\mathbf{p}) = \sum_{i=1}^K p_i \log p_i$

Ω -Regularized Argmax

For convex Ω , define the Ω -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 ℓ_2 -regularization, $\Omega(\mathbf{p}) = \frac{1}{2} \|\mathbf{p}\|^2$.

Ω -Regularized Argmax

For convex Ω , define the Ω -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 ℓ_2 -regularization, $\Omega(\mathbf{p}) = \frac{1}{2} \|\mathbf{p}\|^2$.

Is there something in-between?

α -Entmax

Parametrized by $\alpha \geq 0$:

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

α -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 $\alpha \rightarrow \infty$

α -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 $\alpha \rightarrow 1$

α -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 $\alpha = 2$.

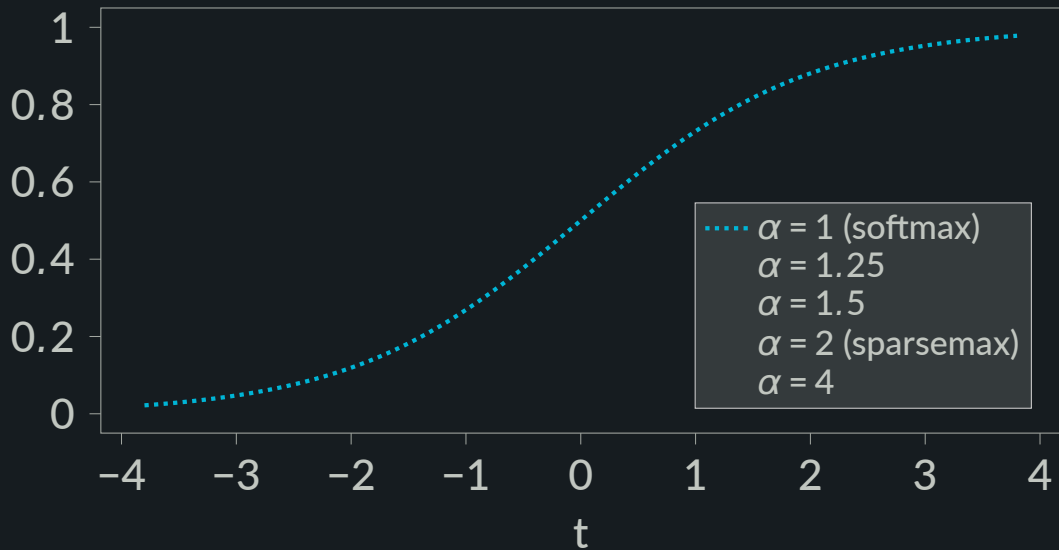
α -Entmax

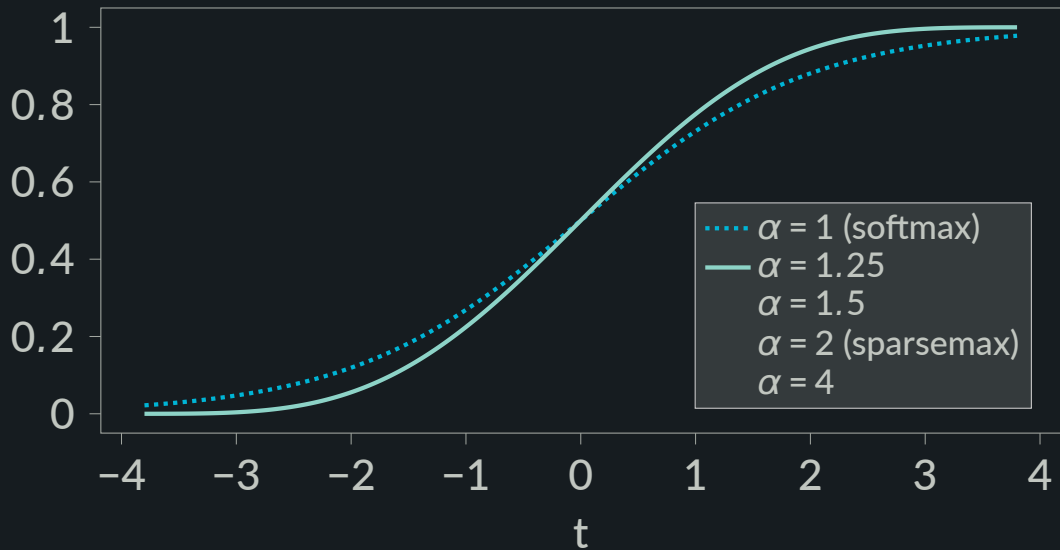
Parametrized by $\alpha \geq 0$:

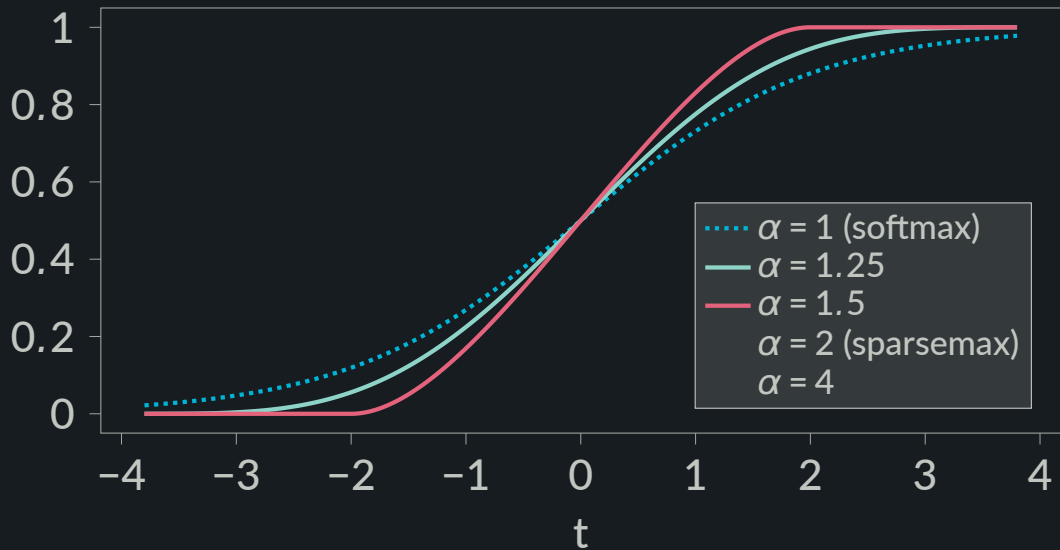
$$\Omega_{\alpha}(\mathbf{p}) := \begin{cases} \frac{1}{\alpha(\alpha-1)} \left(1 - \sum_{i=1}^K p_i^{\alpha} \right) & \text{if } \alpha \neq 1 \\ \sum_{i=1}^K p_i \log p_i & \text{if } \alpha = 1. \end{cases}$$

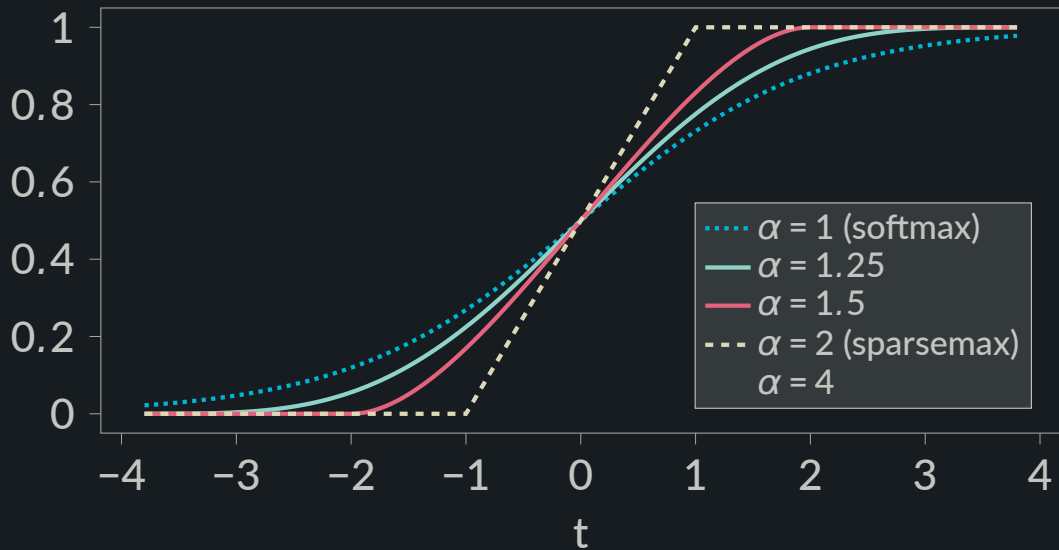
- **Argmax** corresponds to $\alpha \rightarrow \infty$
- **Softmax** amounts to $\alpha \rightarrow 1$
- **Sparsemax** amounts to $\alpha = 2$.

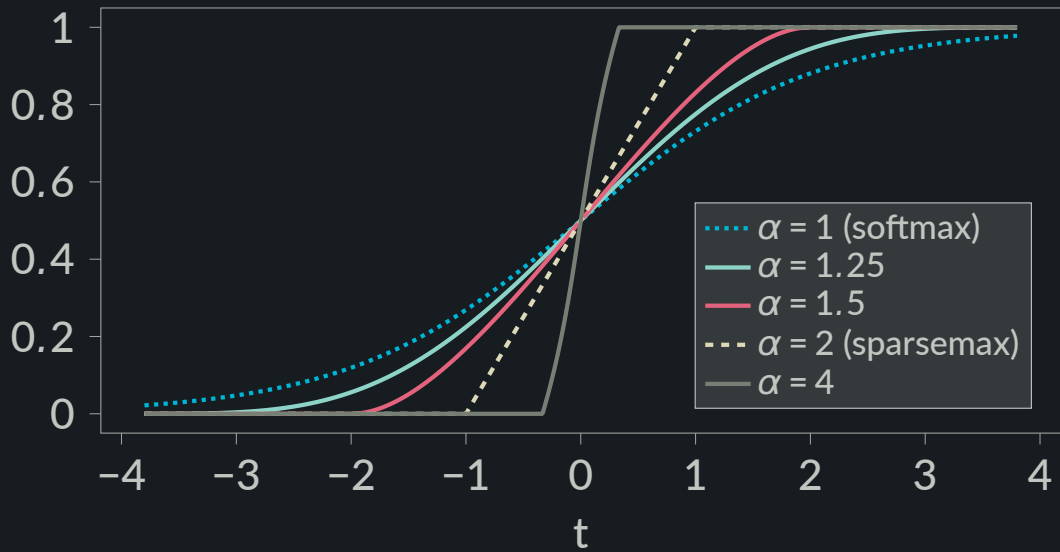
Key result: can be sparse for $\alpha > 1$, propensity for sparsity increases with α .











Learning α

Learning α

Key contribution:

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



Learning α

Key contribution:

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

Requires argmin differentiation → see paper for details!

Learning α

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
α -entmax	29.90	21.74	32.89	26.93

BLEU Scores

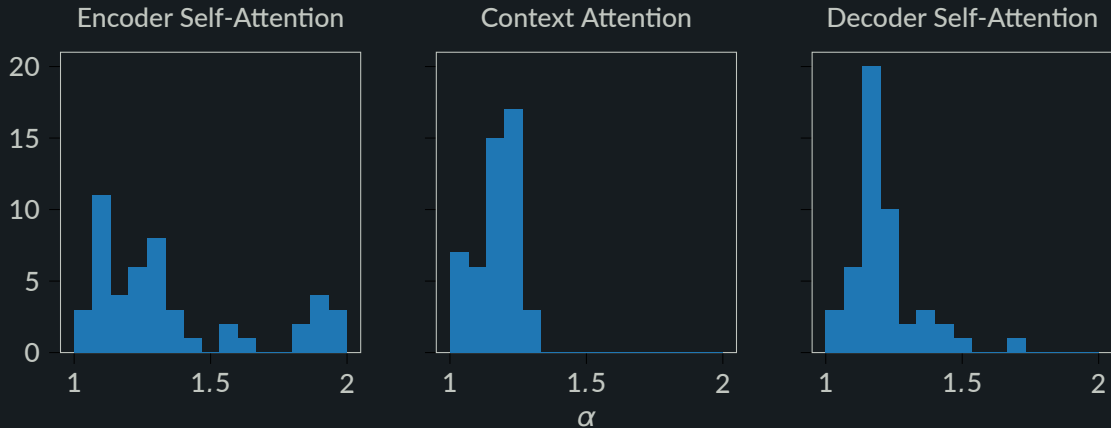
activation	de→en	ja→en	ro→en	en→de
softmax	29.79	21.57	32.70	26.02
1.5-entmax	29.83	22.13	33.10	25.89
α -entmax	29.90	21.74	32.89	26.93

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
α -entmax	29.90	21.74	32.89	26.93

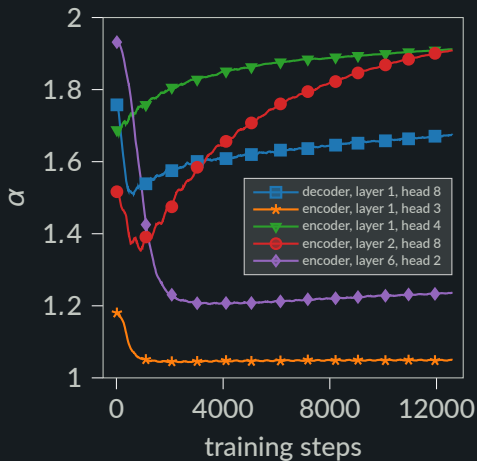
For analysis for other language pairs, see Appendix A.

Learned α

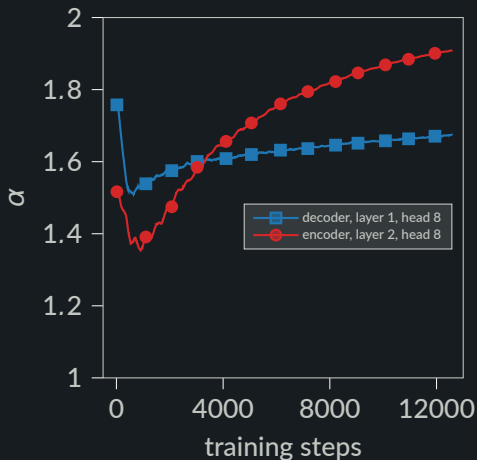


Bimodal for the encoder, mostly unimodal for the decoder.

Trajectories of α During Training

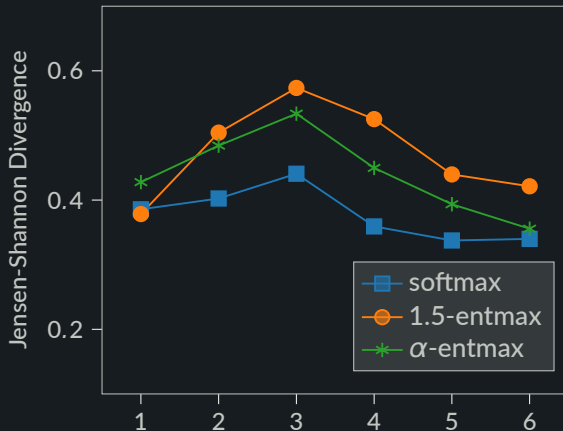


Trajectories of α During Training

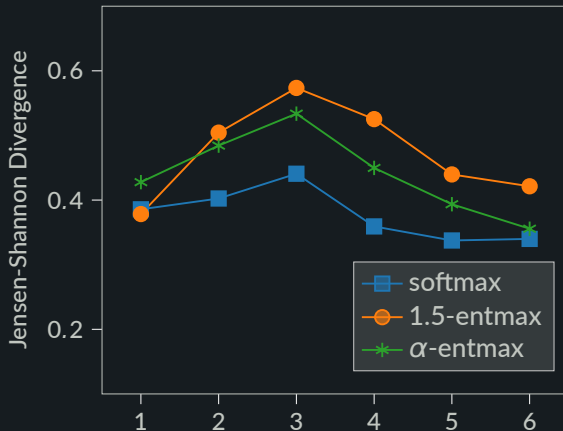


Some heads choose to start dense before becoming sparse.

Head Diversity per Layer

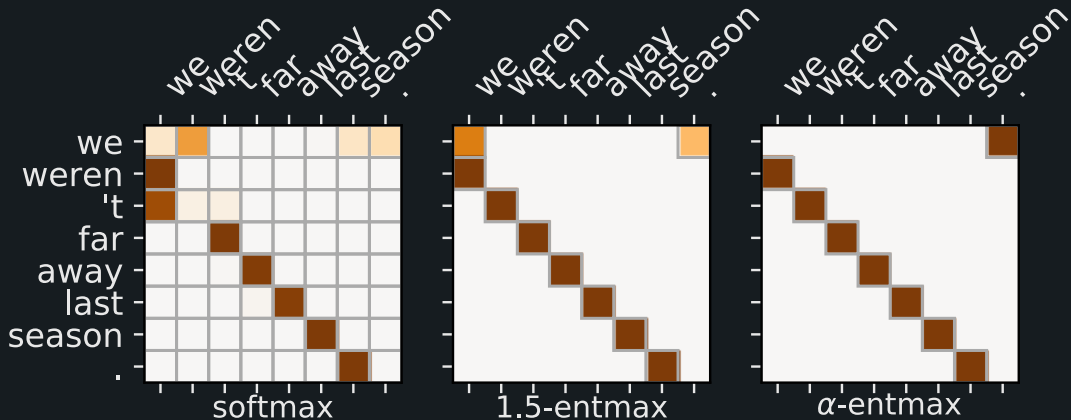


Head Diversity per Layer



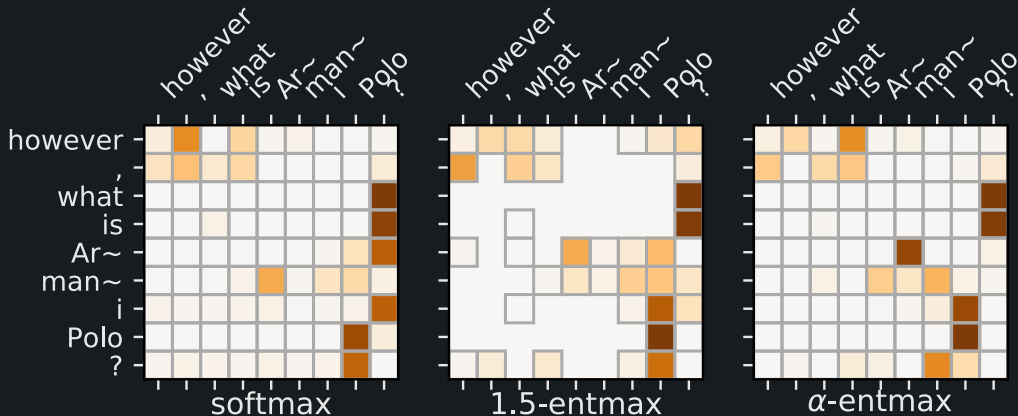
Specialized heads are important as seen in Voita et al. (2019)!

Previous Position Head



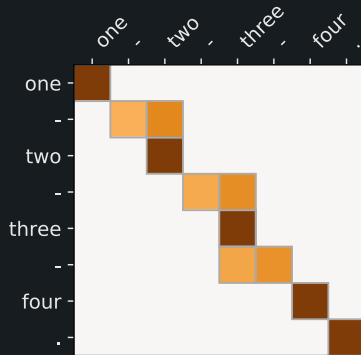
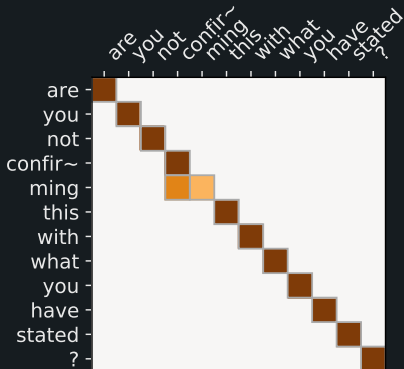
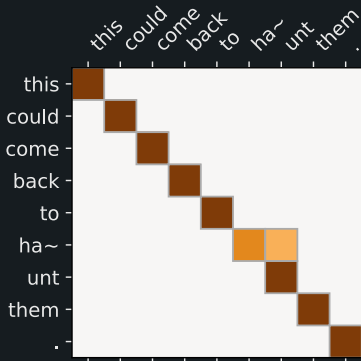
This head role was also found in Voita et al. (2019)! Learned $\alpha = 1.91$.

Interrogation-Detecting Head



Learned $\alpha = 1.05$.

Subword-Merging Head



Learned $\alpha = 1.91$.

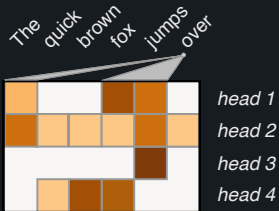
Key Takeaways

Introduce **adaptive** sparsity
for Transformers via α -entmax with a **gradient learnable α** .

Key Takeaways

Introduce **adaptive** sparsity
for Transformers via α -entmax with a **gradient learnable α** .

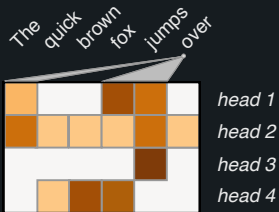
adaptive sparsity



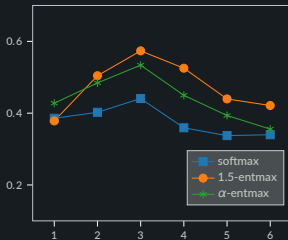
Key Takeaways

Introduce **adaptive** sparsity
for Transformers via α -entmax with a **gradient learnable** α .

adaptive sparsity



reduced head redundancy



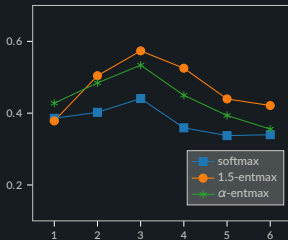
Key Takeaways

Introduce **adaptive** sparsity
for Transformers via α -entmax with a **gradient learnable** α .

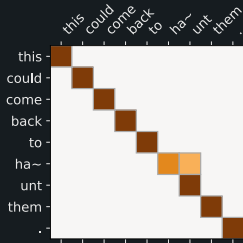
adaptive sparsity



reduced head redundancy



clearer head roles



Thank you!

Questions?

```
:pip install entmax
```

Check 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

- Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio (2015). "Neural machine translation by jointly learning to align and translate". In: *Proc. of ICLR*.
- 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*.