

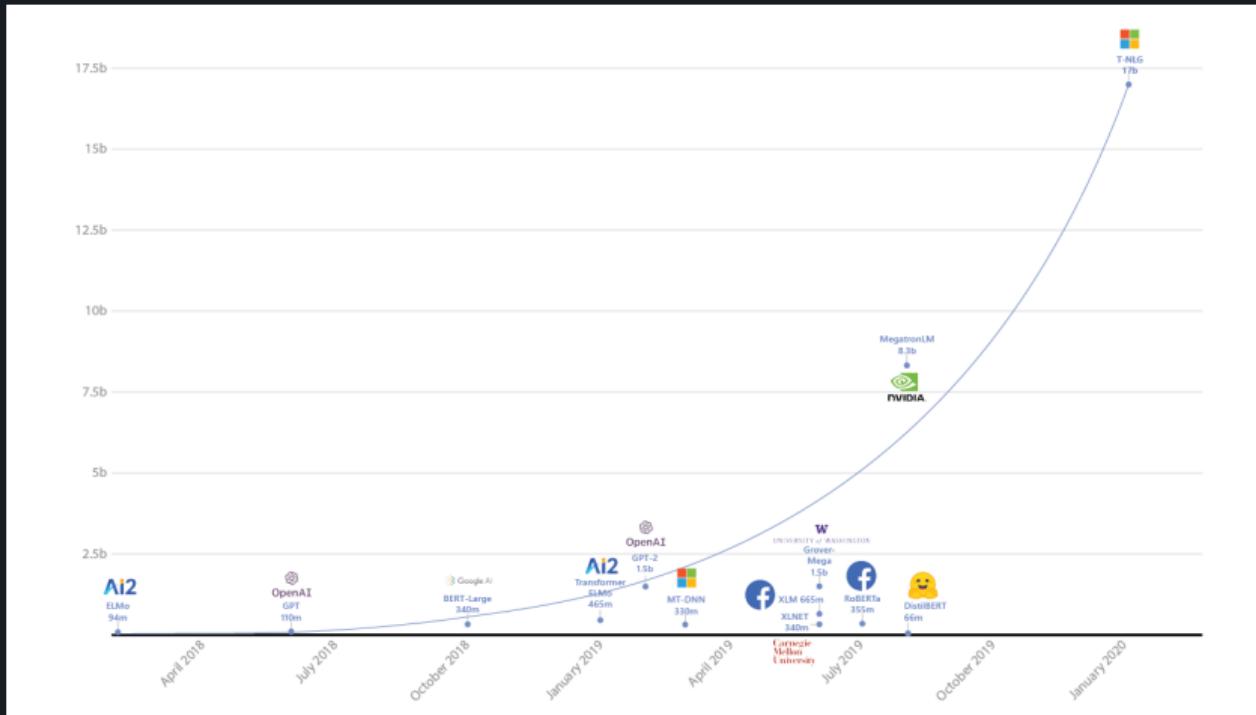
# **Thesis Proposal:**

# **Neural Machine Translation with Latent Context**

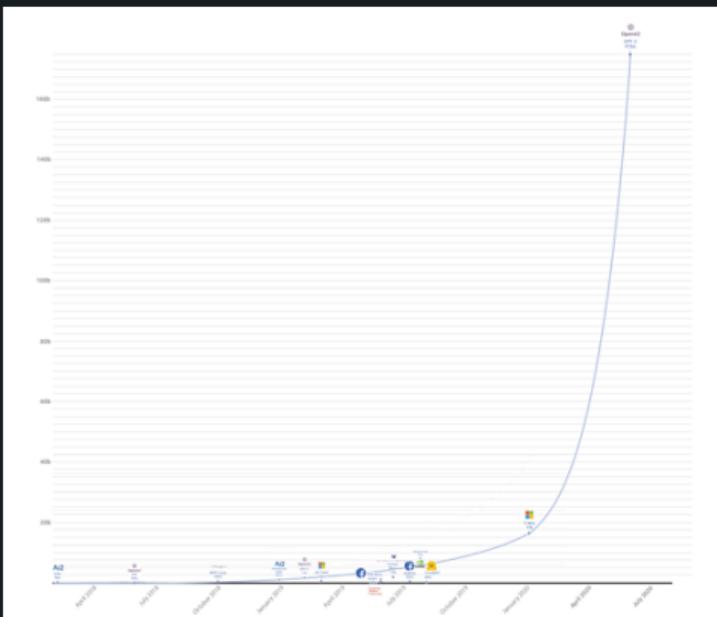
**Gonçalo Correia**

Committee: André Martins, Mário Figueiredo, Vlad Niculae, Wilker Aziz

# Neural networks: big, uninterpretable, and powerful black-boxes?



# Neural networks: big, uninterpretable, and powerful black-boxes?



# Key concepts of this thesis

- Transparency

# Key concepts of this thesis

- Transparency
- Sparsity

# Key concepts of this thesis

- Transparency
- Sparsity
- Compactness

# Key concepts of this thesis

- Transparency
- Sparsity
- Compactness
- Data efficiency

# Key concepts of this thesis

- Transparency
- Sparsity
- Compactness
- Data efficiency

**Les artistes sont démunis**  
  
**The artists don't have any money**

# Neural Machine Translation (NMT)

- Advances in optimization → ditch structure in favour of direct decomposition

# Neural Machine Translation (NMT)

- Advances in optimization → ditch structure in favour of direct decomposition
- Good for high-resource language, fails in low-resource settings

# A bit of context... On Seq2Seq

*United Nations elections end today*

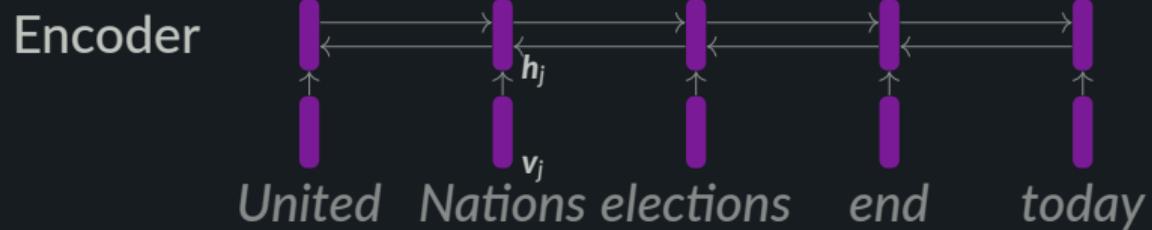
Slide Credit: Vlad Niculae

# A bit of context... On Seq2Seq

Encoder

The diagram illustrates an encoder processing the sentence "United Nations elections end today". The words are represented by vertical purple bars, and the word "elections" is annotated with a subscript  $v_j$ . Below each bar, the corresponding word is written in a light gray font: "United", " $v_j$  Nations", "elections", "end", and "today".

# A bit of context... On Seq2Seq



# A bit of context... On Seq2Seq

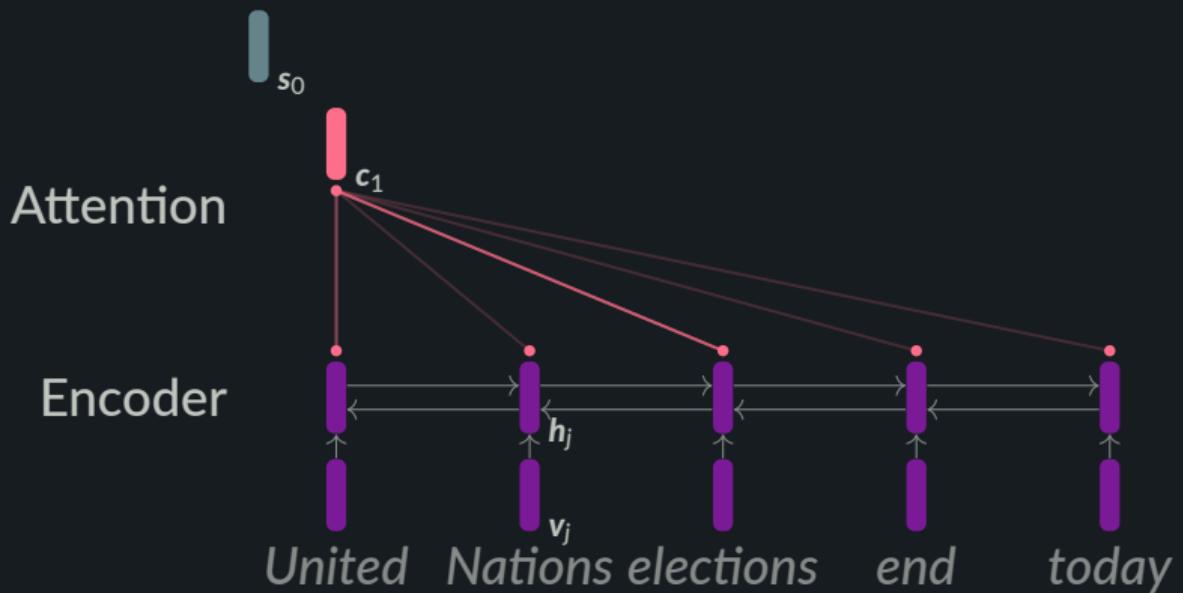
**attention weights**  
computed with  
*softmax*:

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

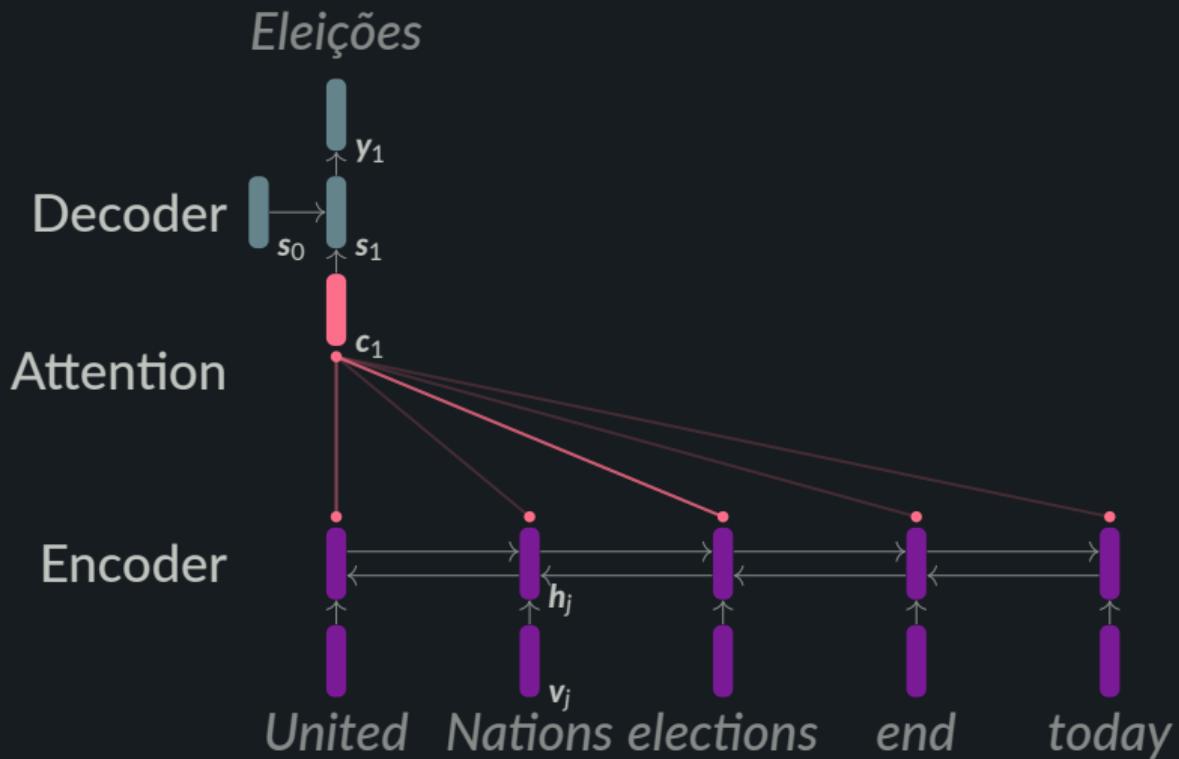
$$z_j = s_t^\top W^{(a)} h_j$$

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

$$c_t = \sum_j \pi_j h_j$$



# A bit of context... On Seq2Seq



**attention weights**  
computed with  
*softmax*:

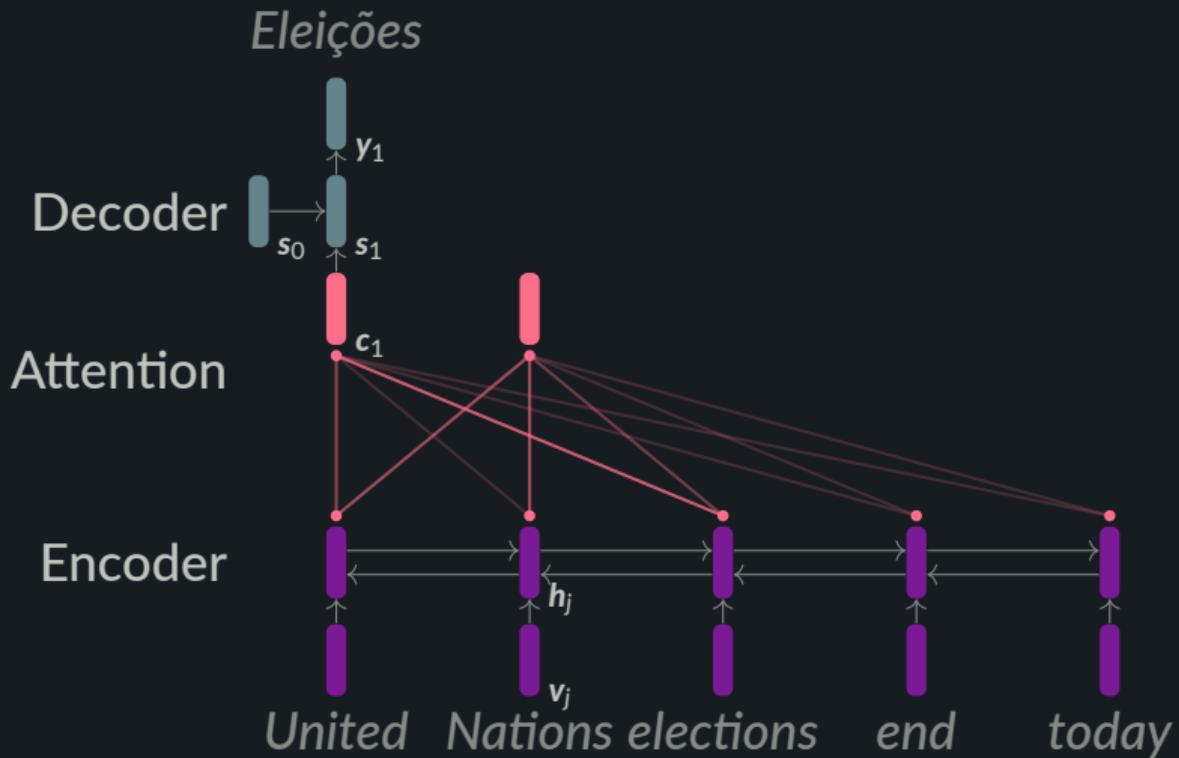
for some decoder state  $s_t$ ,  
compute contextually  
weighted average of input  $c_t$ :

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

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

$$c_t = \sum_j \pi_j h_j$$

# A bit of context... On Seq2Seq



**attention weights**  
computed with  
*softmax*:

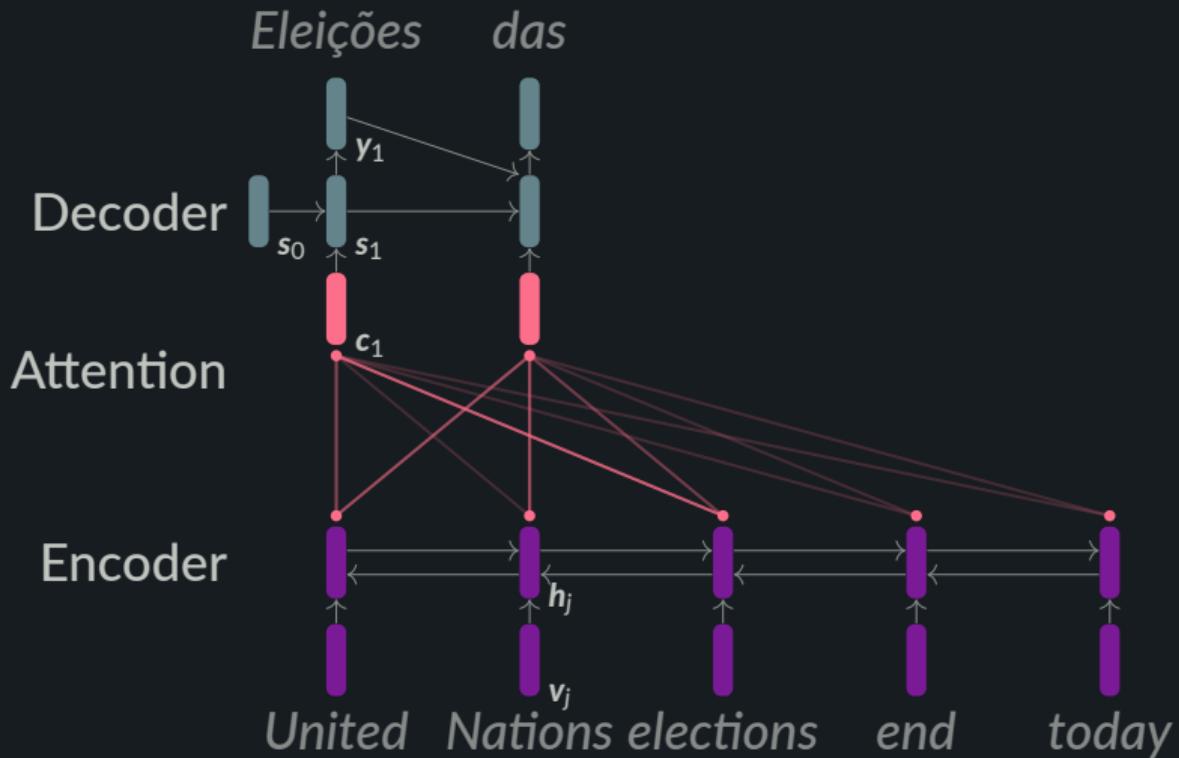
for some decoder state  $s_t$ ,  
compute contextually  
weighted average of input  $c_t$ :

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

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

$$c_t = \sum_j \pi_j h_j$$

# A bit of context... On Seq2Seq



**attention weights**  
computed with  
*softmax*:

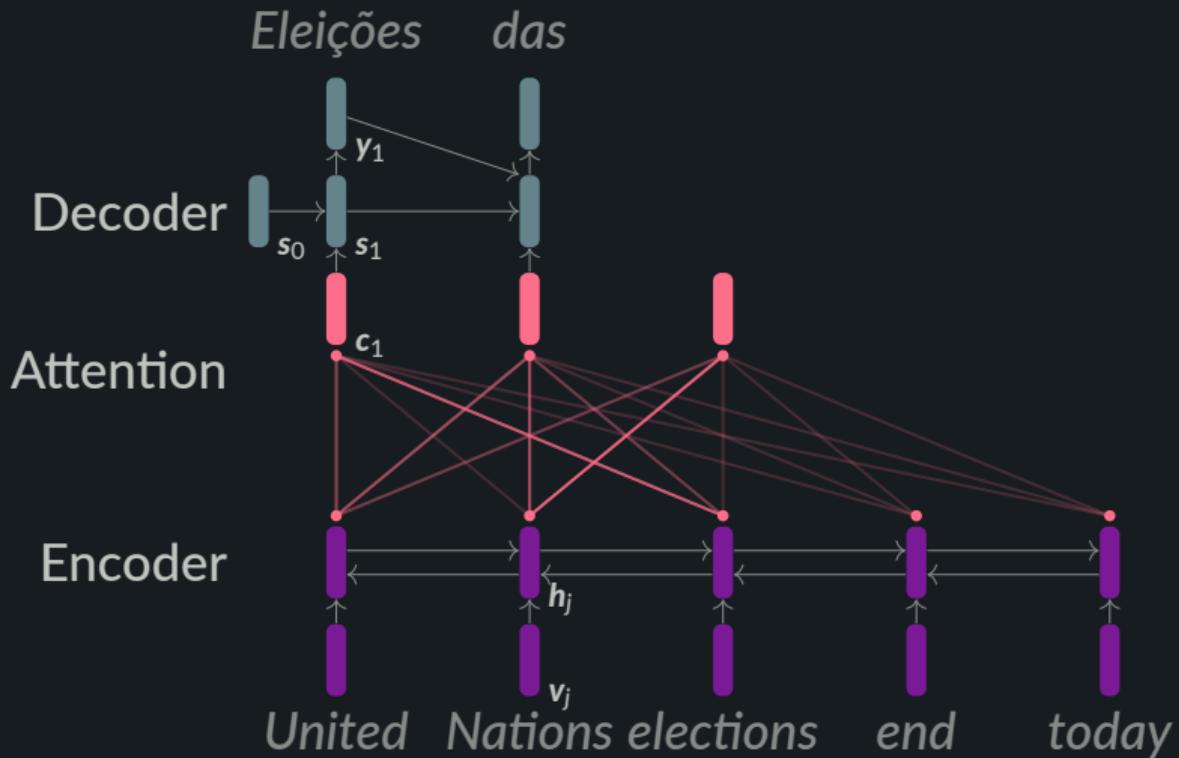
for some decoder state  $s_t$ ,  
compute contextually  
weighted average of input  $c_t$ :

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

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

$$c_t = \sum_j \pi_j h_j$$

# A bit of context... On Seq2Seq



**attention weights**  
computed with  
*softmax*:

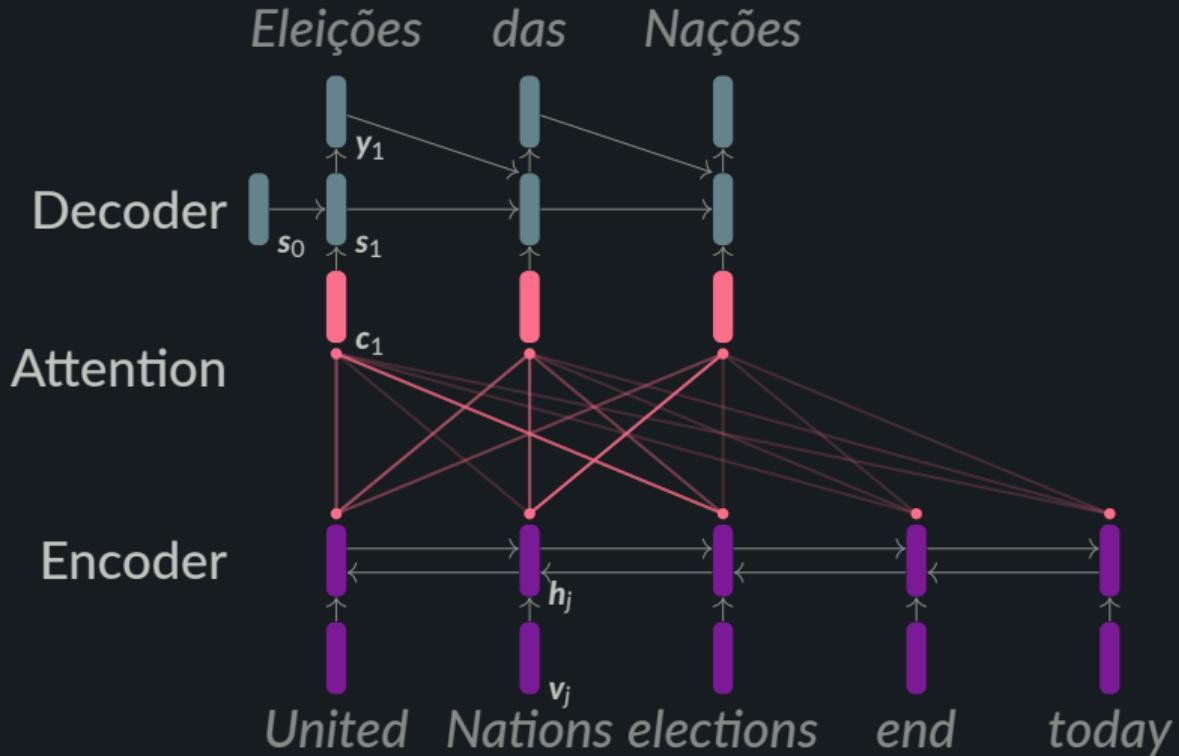
for some decoder state  $s_t$ ,  
compute contextually  
weighted average of input  $c_t$ :

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

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

$$c_t = \sum_j \pi_j h_j$$

# A bit of context... On Seq2Seq



**attention weights**  
computed with  
*softmax*:

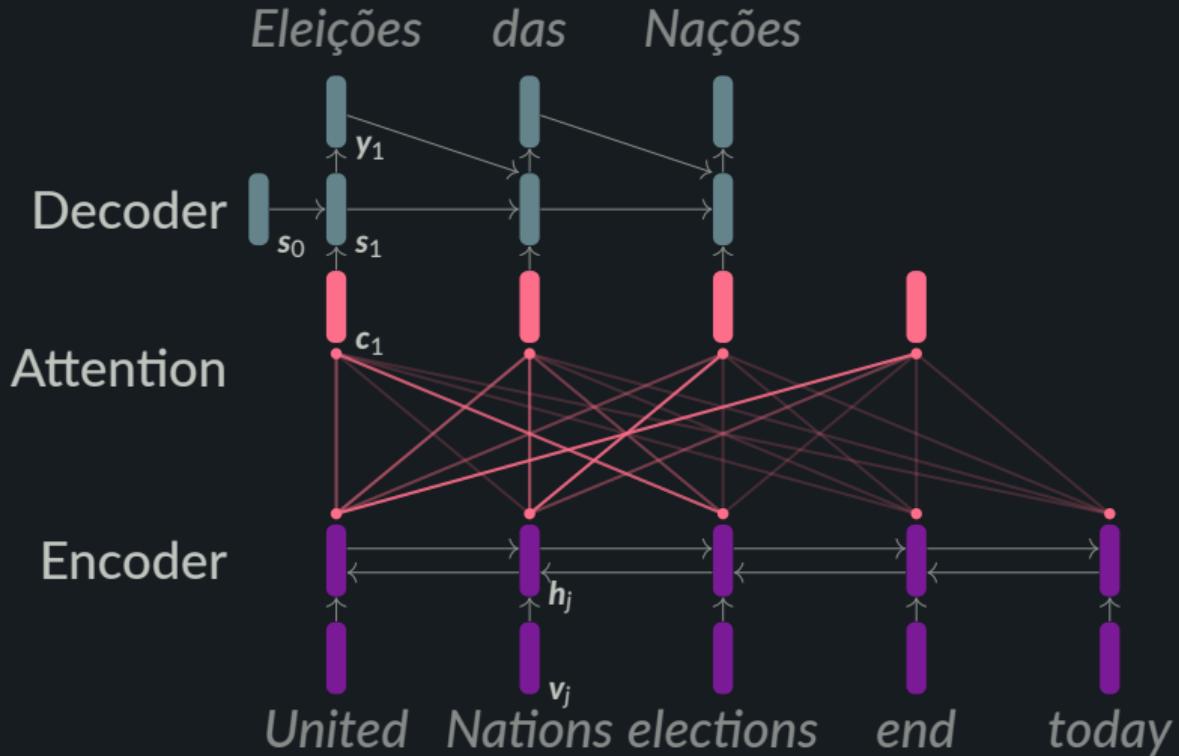
for some decoder state  $s_t$ ,  
compute contextually  
weighted average of input  $c_t$ :

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

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

$$c_t = \sum_j \pi_j h_j$$

# A bit of context... On Seq2Seq



**attention weights**  
computed with  
*softmax*:

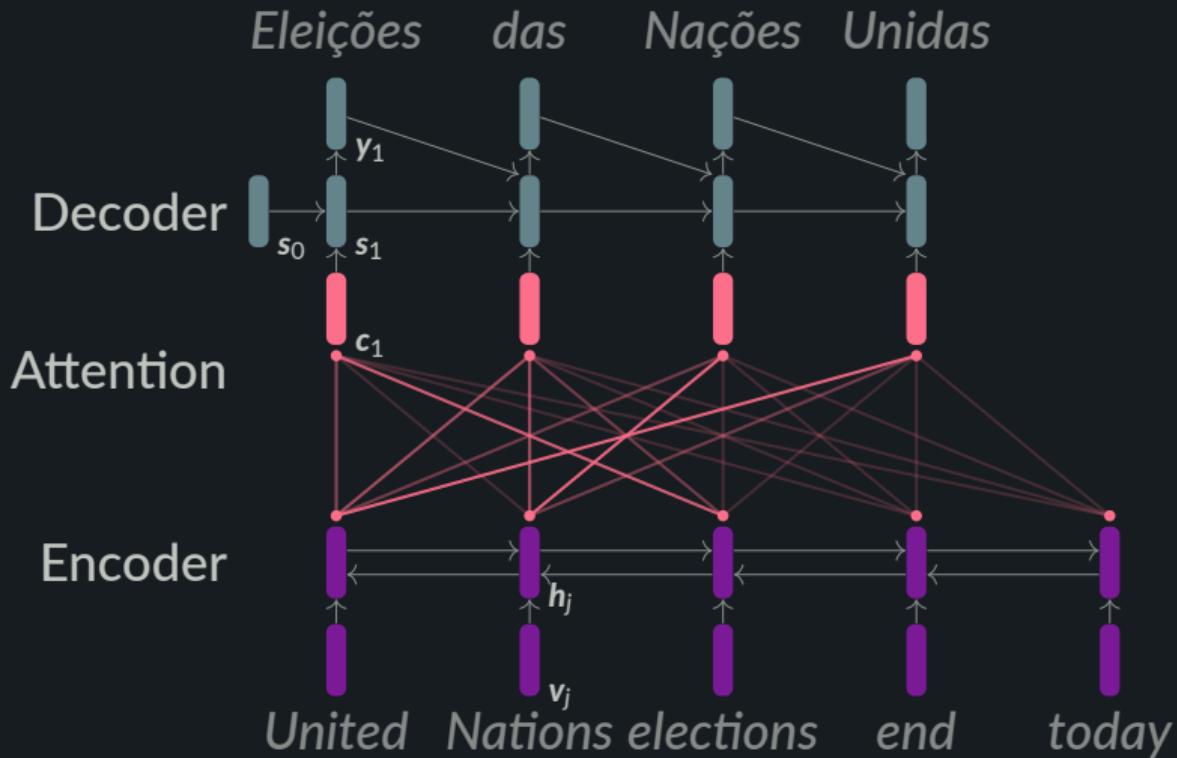
for some decoder state  $s_t$ ,  
compute contextually  
weighted average of input  $c_t$ :

$$z_j = s_t^\top \mathbf{W}^{(a)} \mathbf{h}_j$$

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

$$\mathbf{c}_t = \sum_j \pi_j \mathbf{h}_j$$

# A bit of context... On Seq2Seq



**attention weights**  
computed with  
*softmax*:

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

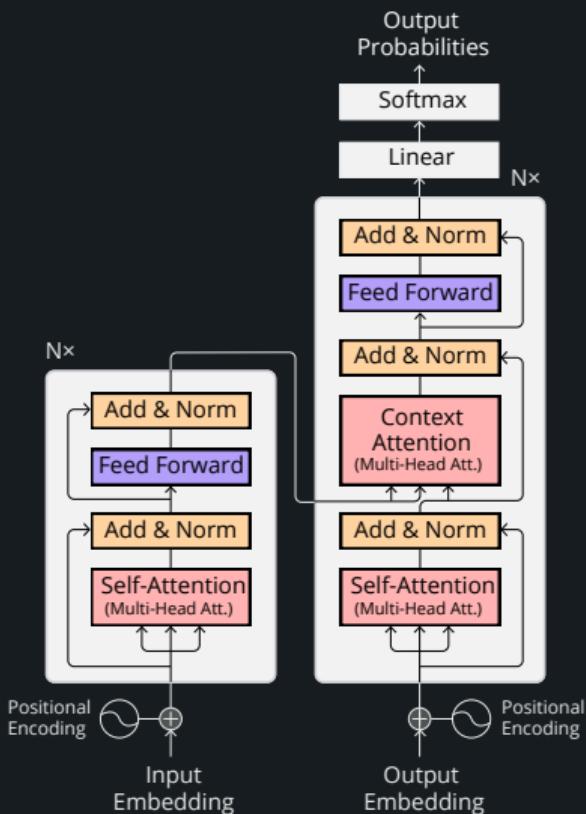
$$z_j = s_t^\top W^{(a)} h_j$$

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

$$c_t = \sum_j \pi_j h_j$$

# A bit of context... on Transformers

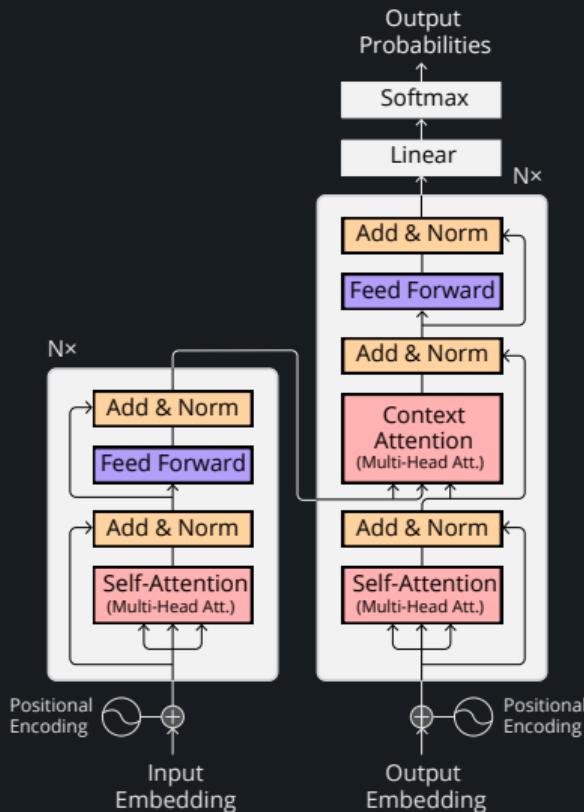
What if... Attention is all you need?



# 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!

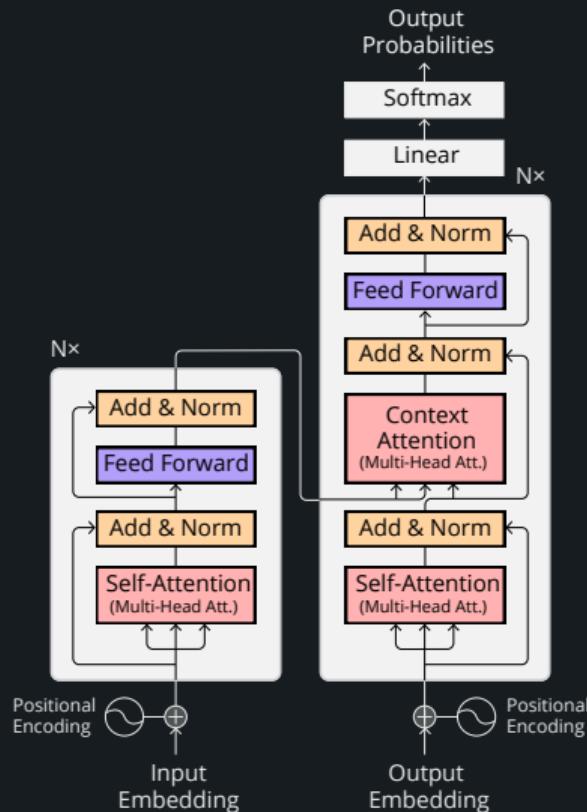


# A bit of context... on Transformers

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

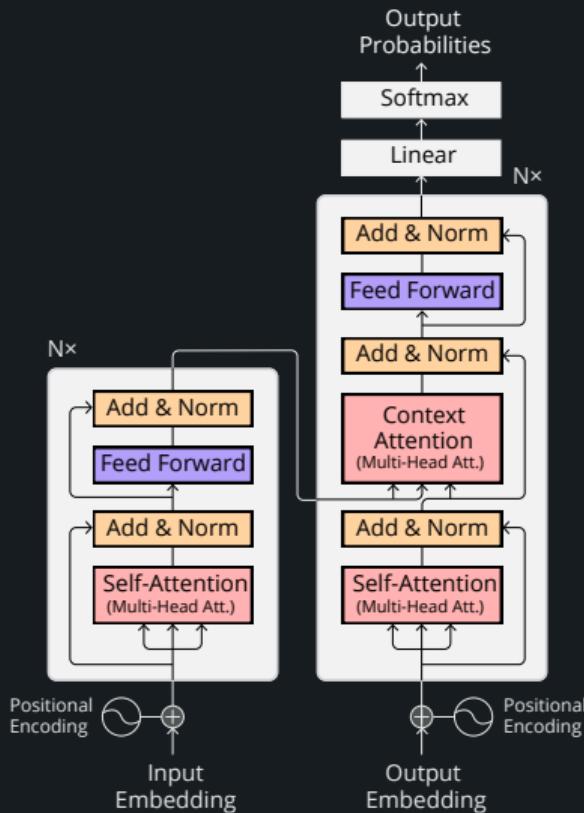


# A bit of context... on Transformers

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)

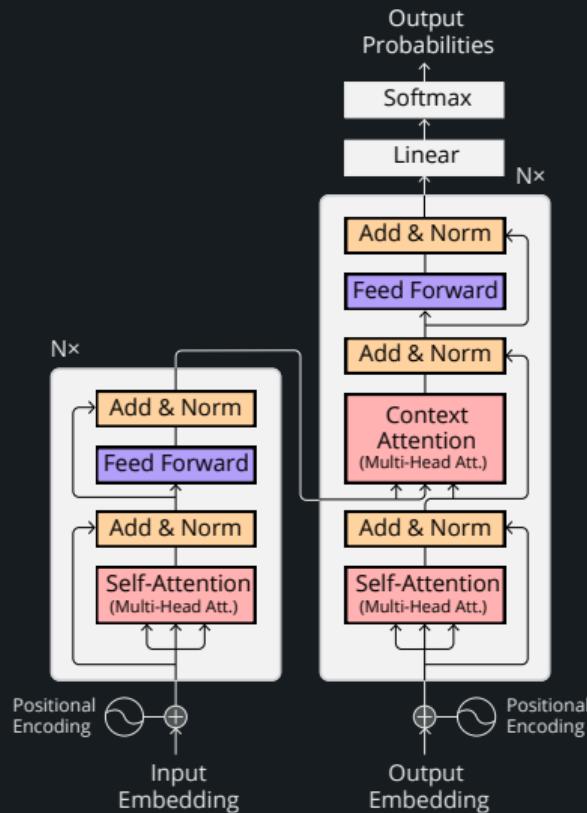


# A bit of context... on Transformers

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

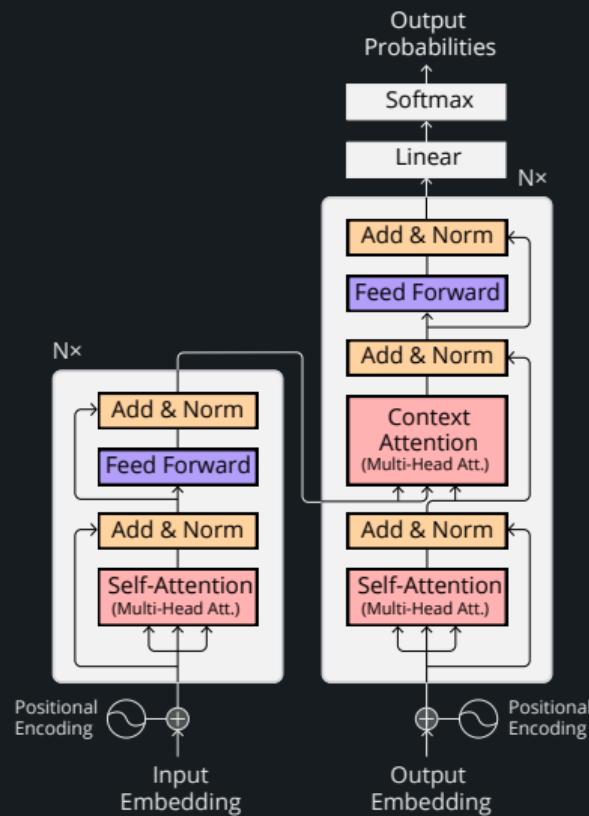


# A bit of context... on Transformers

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
- Inspiration for big general-purpose models like BERT!



# Outline

# Outline

- Automatic Post-Editing (APE) using BERT (ACL 2019)

# Outline

- Automatic Post-Editing (APE) using BERT (ACL 2019)
- Letting Transformer learn how sparse it wants its attention (EMNLP 2019)

# Outline

- Automatic Post-Editing (APE) using BERT (ACL 2019)
- Letting Transformer learn how sparse it wants its attention (EMNLP 2019)
- General strategy for efficient training discrete latent variables (NeurIPS 2020)

# Outline

- Automatic Post-Editing (APE) using BERT (ACL 2019)
- Letting Transformer learn how sparse it wants its attention (EMNLP 2019)
- General strategy for efficient training discrete latent variables (NeurIPS 2020)
- Future work: powerful structured latent variable models for NMT

# Table of Contents

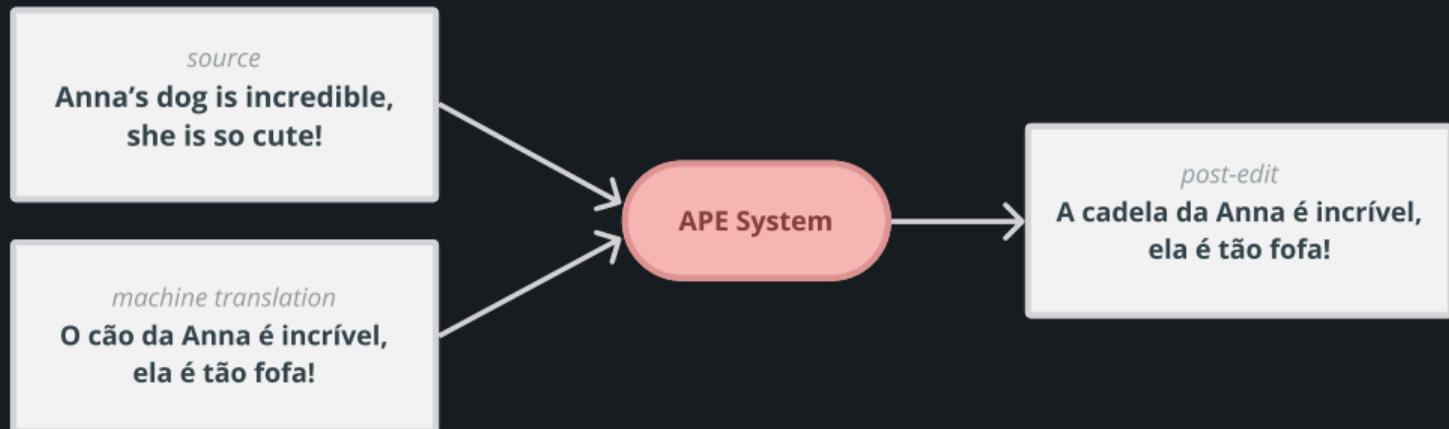
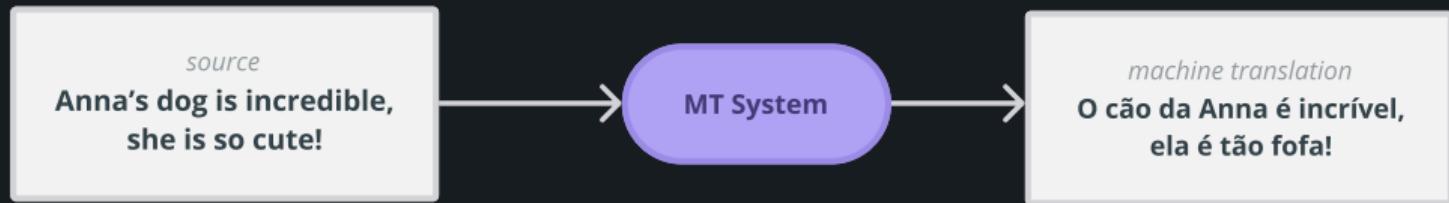
A Simple and Effective Approach to Automatic Post-Editing with Transfer Learning

Adaptively Sparse Transformers

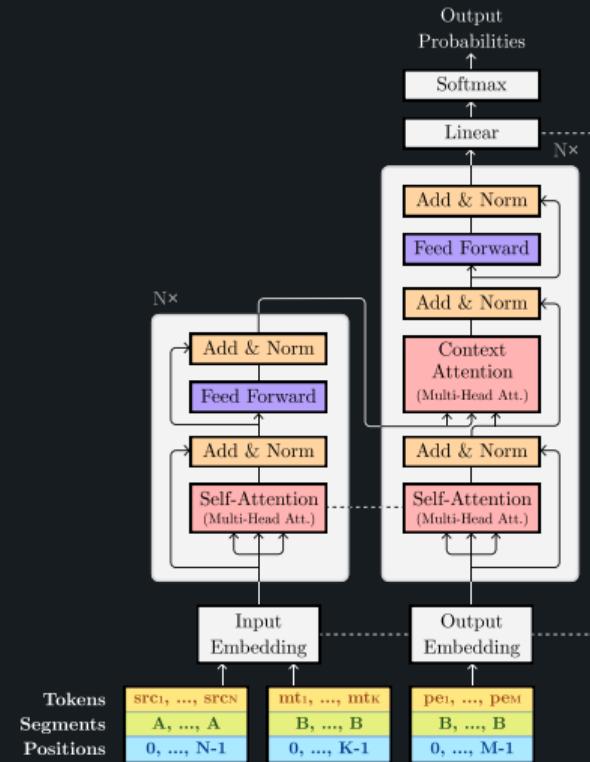
Efficient Marginalization of Discrete and Structured Latent Variables via Sparsity

Proposed Work

# What is APE?

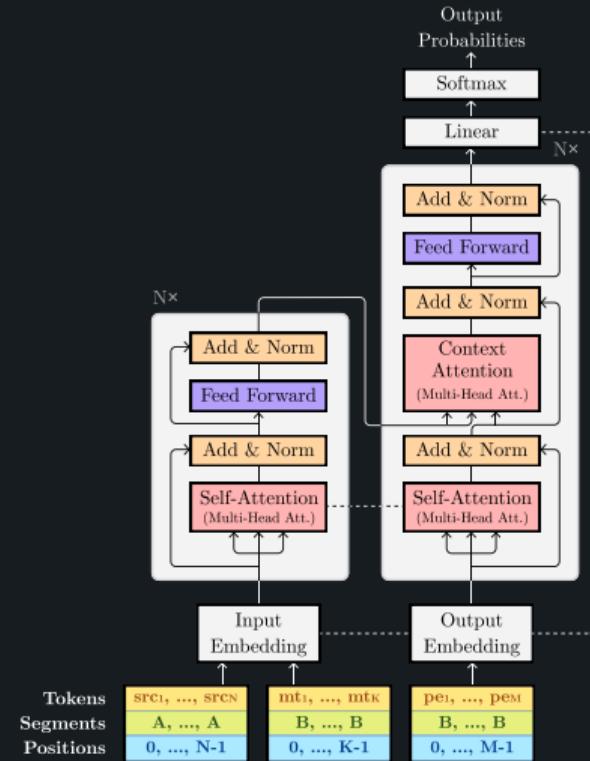


# BERT for APE



# BERT for APE

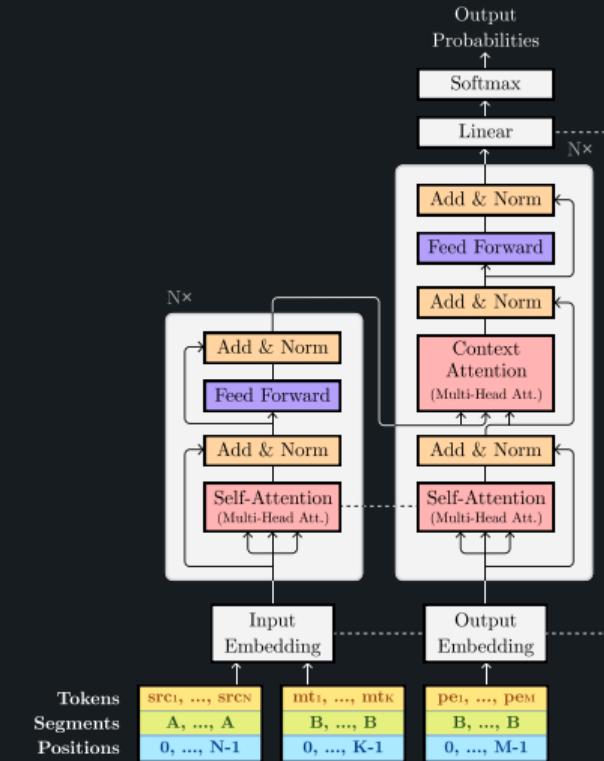
Key idea: Use BERT to do APE



# BERT for APE

Key idea: Use BERT to do APE

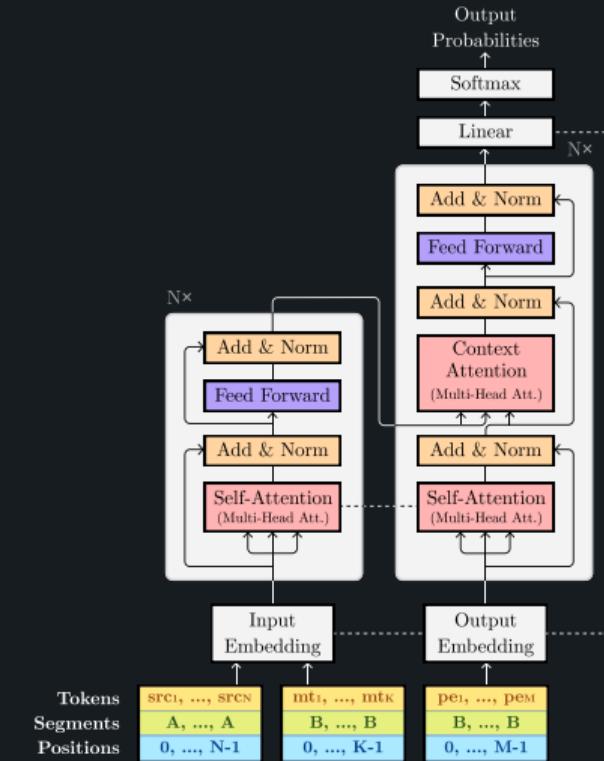
- Prior to this work, BERT was mainly used for simple classification tasks



# BERT for APE

Key idea: Use BERT to do APE

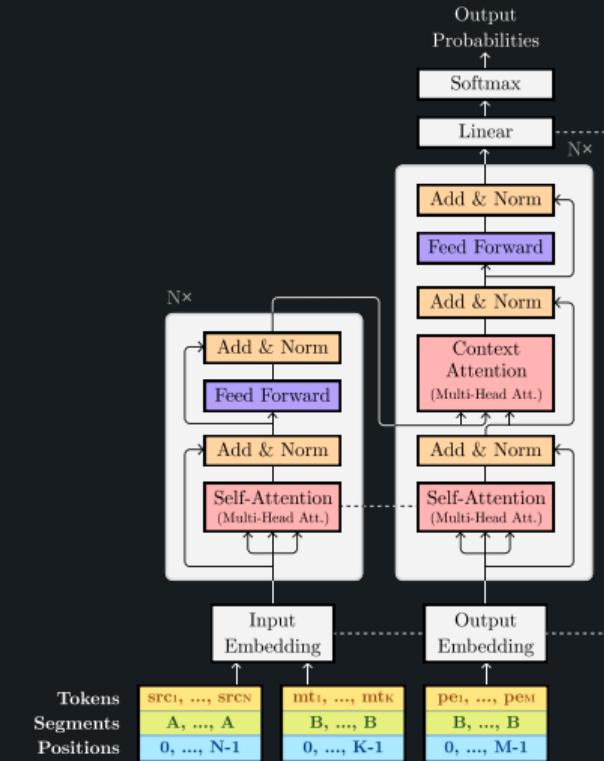
- Prior to this work, BERT was mainly used for simple classification tasks
- We introduced an effective method to use BERT in a generation task (APE)



# BERT for APE

Key idea: Use BERT to do APE

- Prior to this work, BERT was mainly used for simple classification tasks
- We introduced an effective method to use BERT in a generation task (APE)
- Smart parameter sharing between encoder and decoder



# Table of Contents

A Simple and Effective Approach to Automatic Post-Editing with Transfer Learning

Adaptively Sparse Transformers

Efficient Marginalization of Discrete and Structured Latent Variables via Sparsity

Proposed Work

# Getting to know attention heads better

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

# Getting to know attention heads better

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

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

- for interpretability
- for discovering linguistic structure
- for efficiency

# Getting to know attention heads better

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

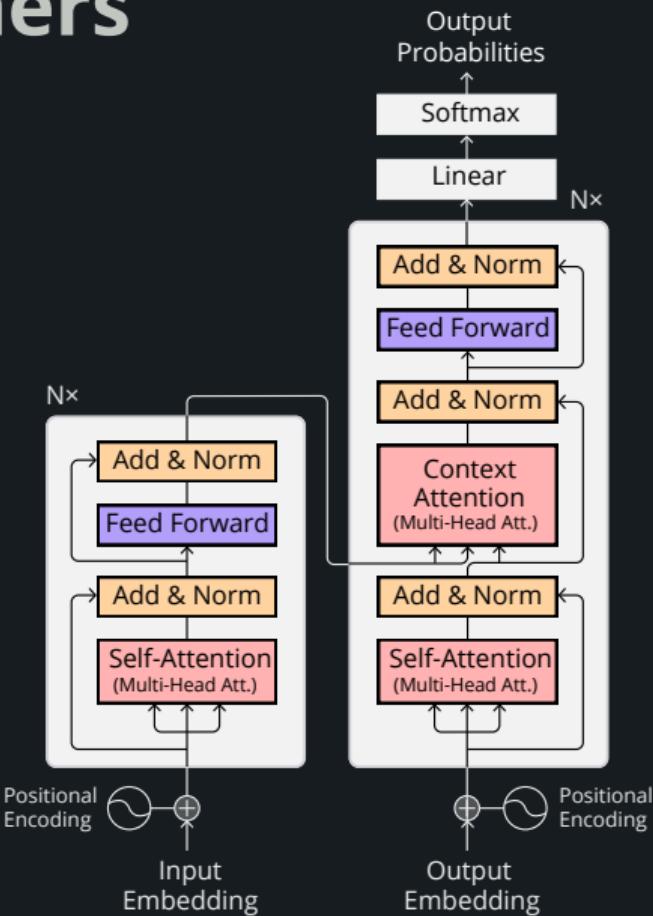
Our solution is to bet on **sparsity**:

- for interpretability
- for discovering linguistic structure
- for efficiency

# Transformers

In each attention head:

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



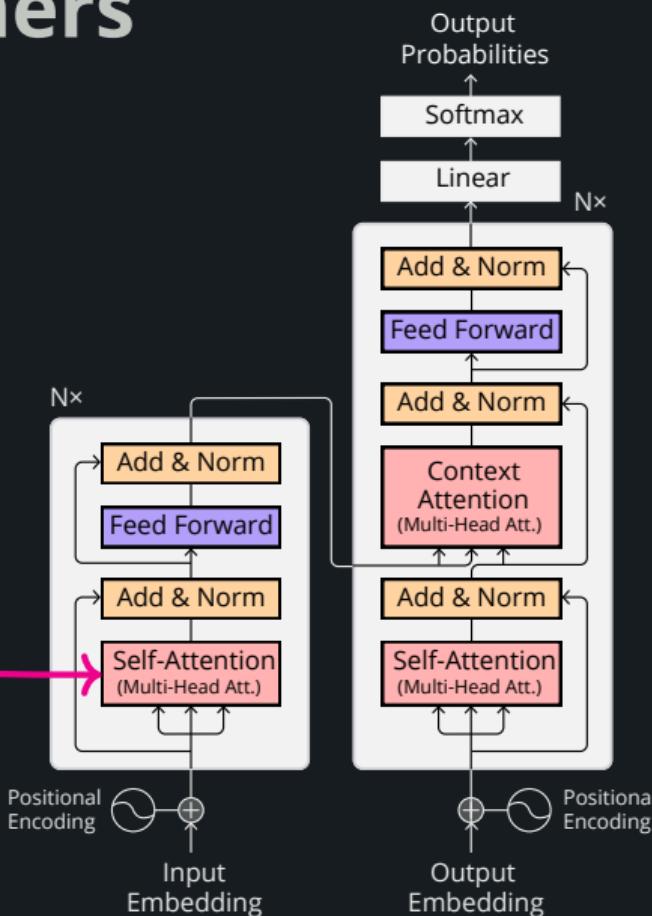
# Transformers

In each attention head:

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

Attention in three places:

- Self-attention in the encoder



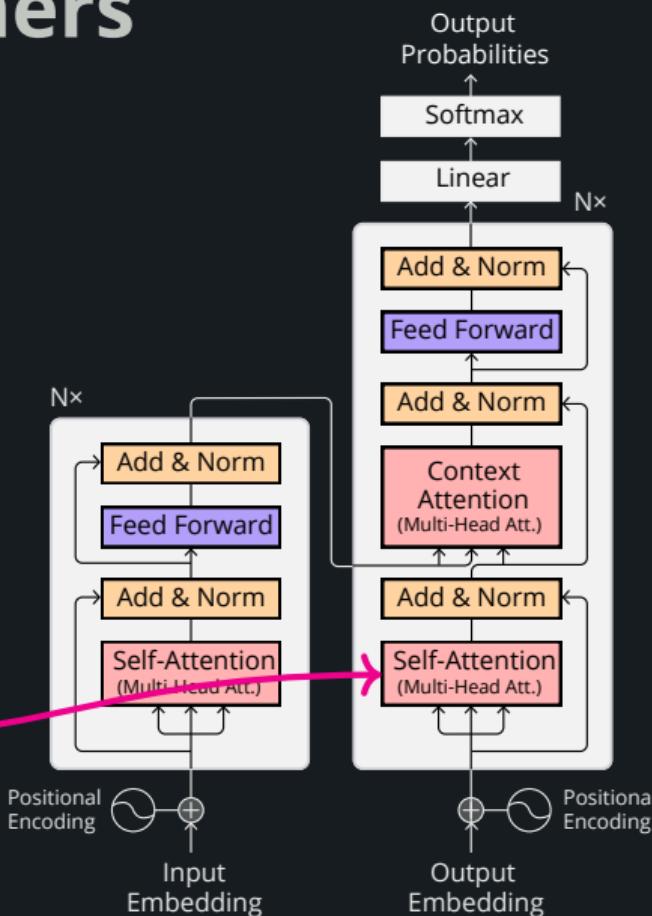
# Transformers

In each attention head:

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

Attention in three places:

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



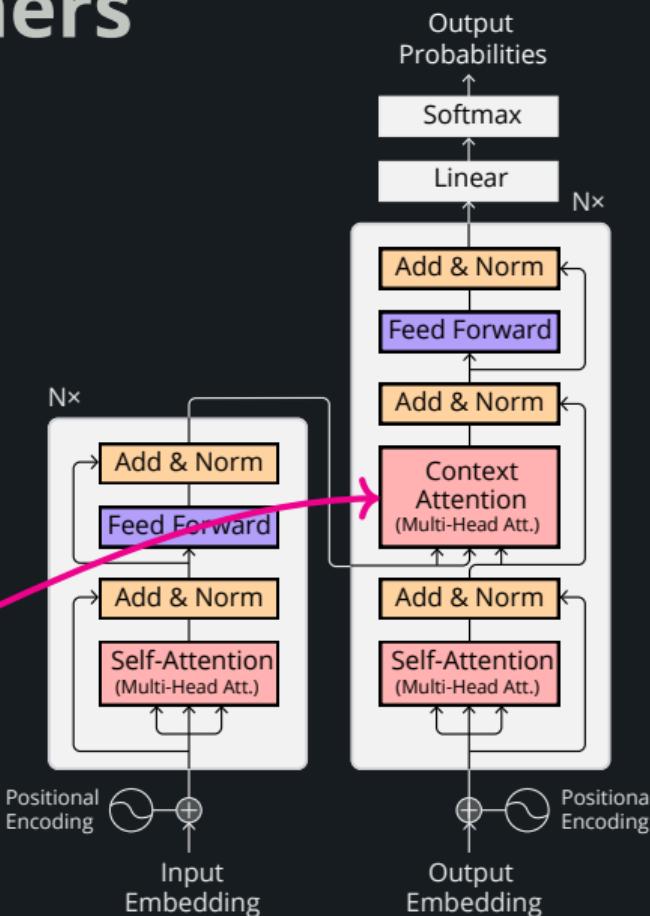
# Transformers

In each attention head:

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

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  $\alpha$ ! 

# 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: **softmax(z) > 0**

# $\alpha$ -Entmax

Parametrized by  $\alpha \geq 0$ :

# $\alpha$ -Entmax

Parametrized by  $\alpha \geq 0$ :

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

# $\alpha$ -Entmax

Parametrized by  $\alpha \geq 0$ :

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

# $\alpha$ -Entmax

Parametrized by  $\alpha \geq 0$ :

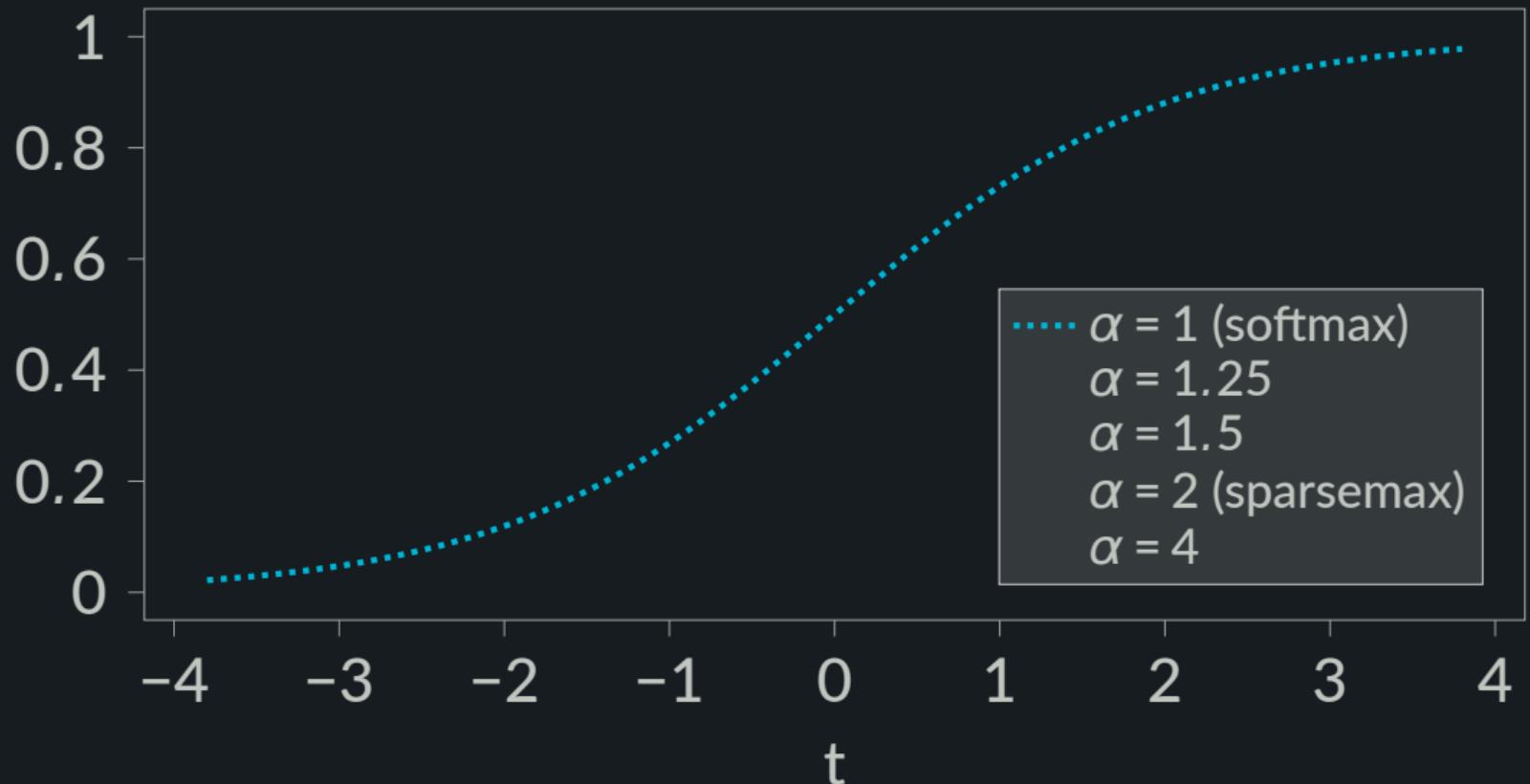
- **Argmax** corresponds to  $\alpha \rightarrow \infty$
- **Softmax** amounts to  $\alpha \rightarrow 1$
- **Sparsemax** amounts to  $\alpha = 2$ .

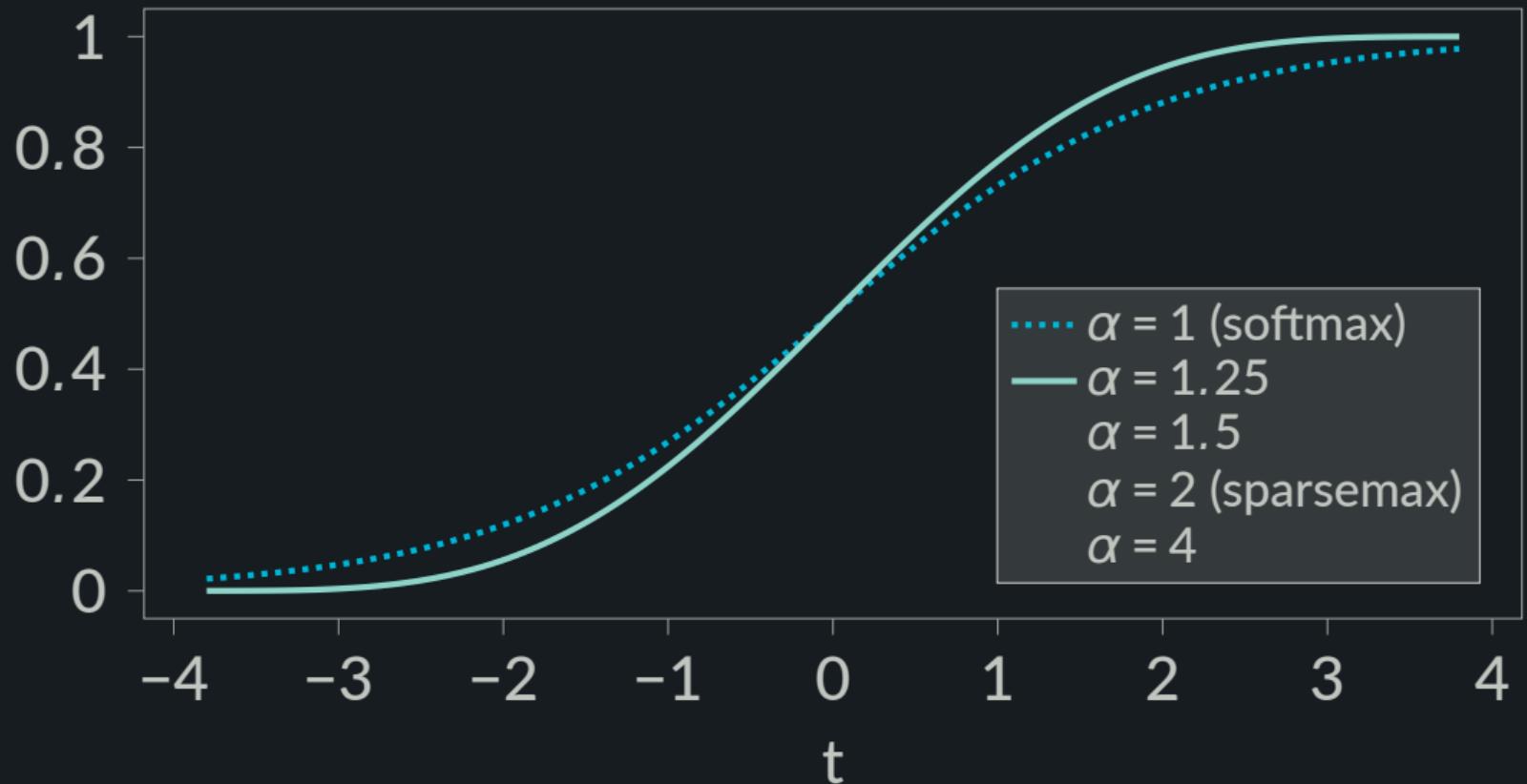
# $\alpha$ -Entmax

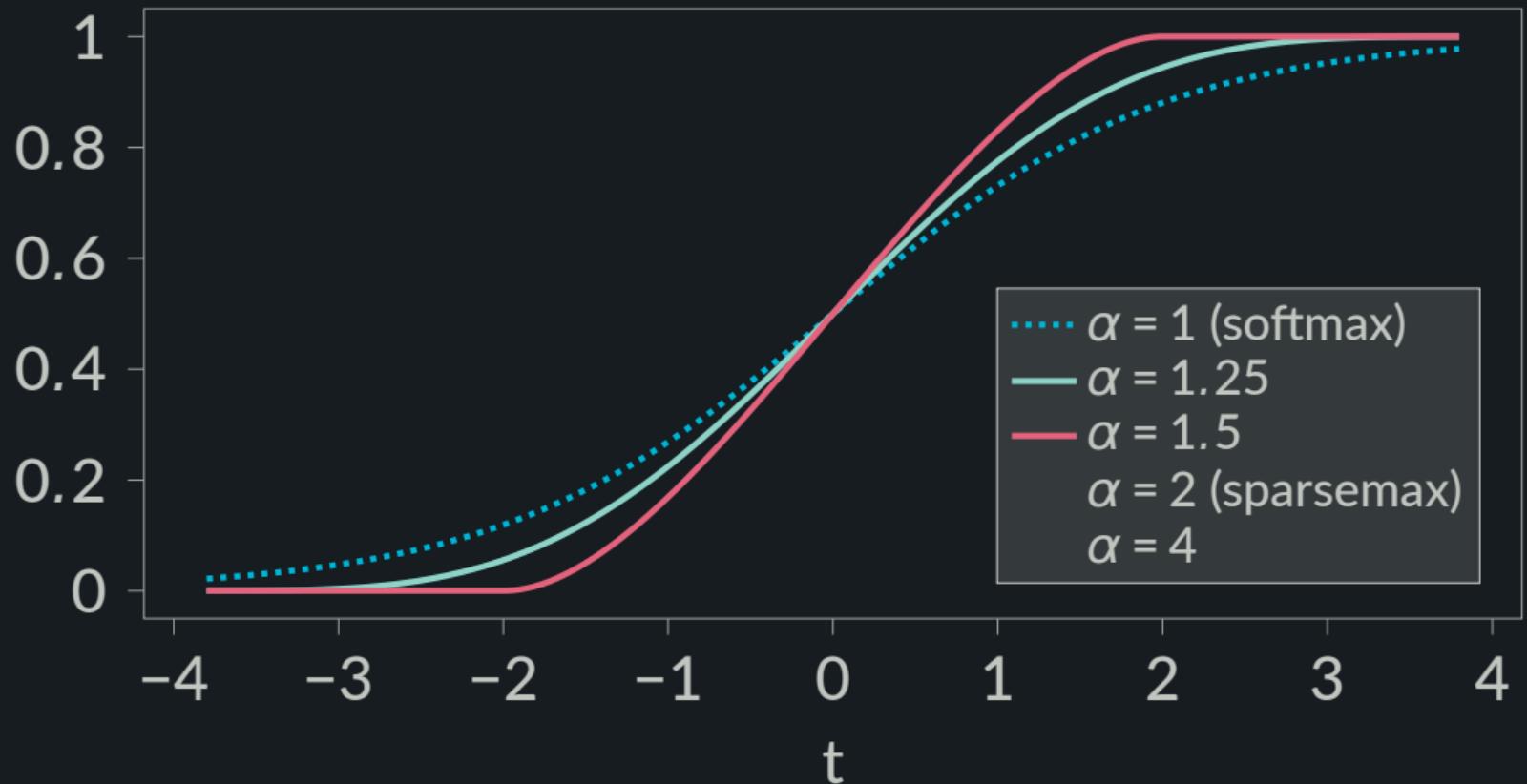
Parametrized by  $\alpha \geq 0$ :

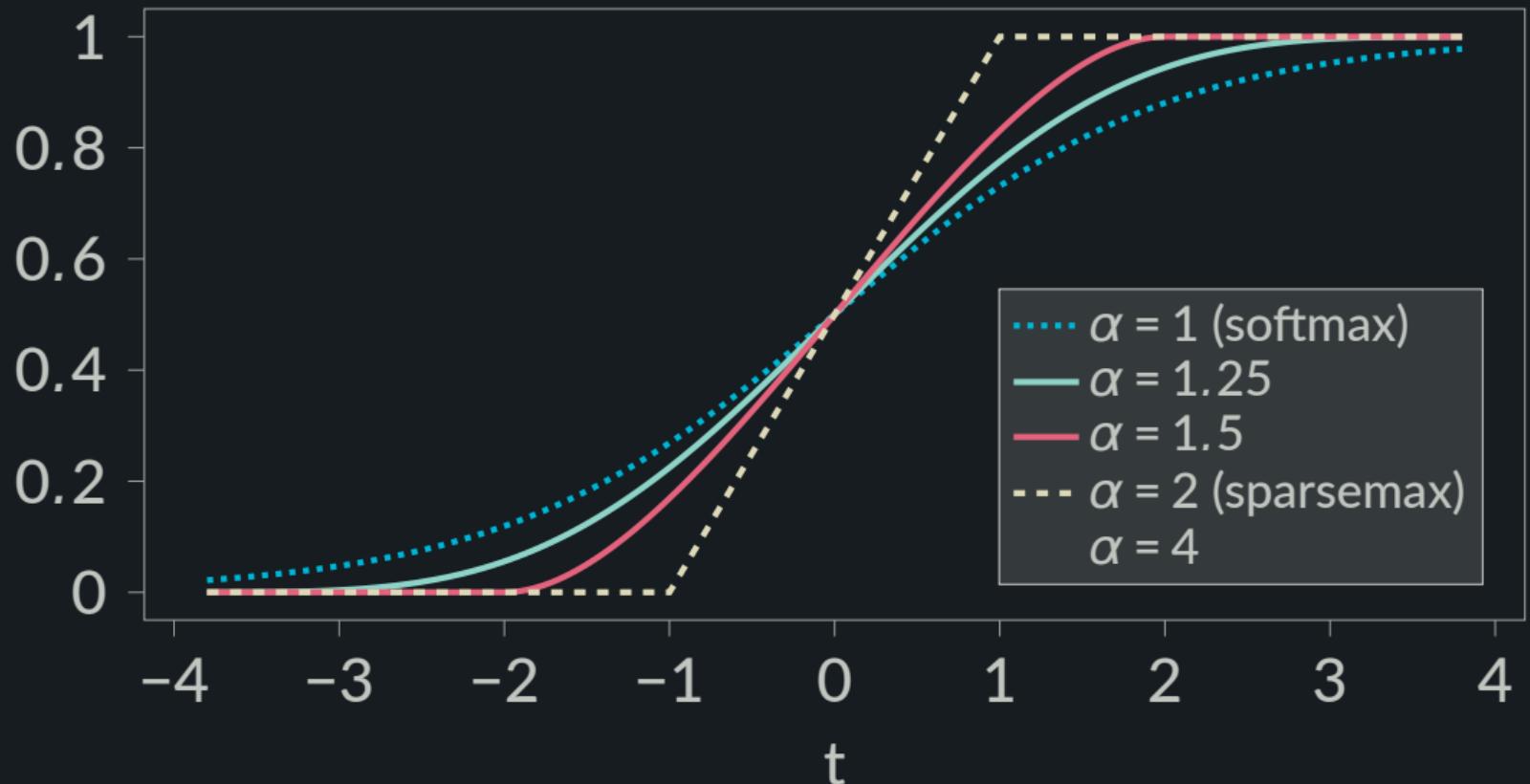
- **Argmax** corresponds to  $\alpha \rightarrow \infty$
- **Softmax** amounts to  $\alpha \rightarrow 1$
- **Sparsemax** amounts to  $\alpha = 2$ .

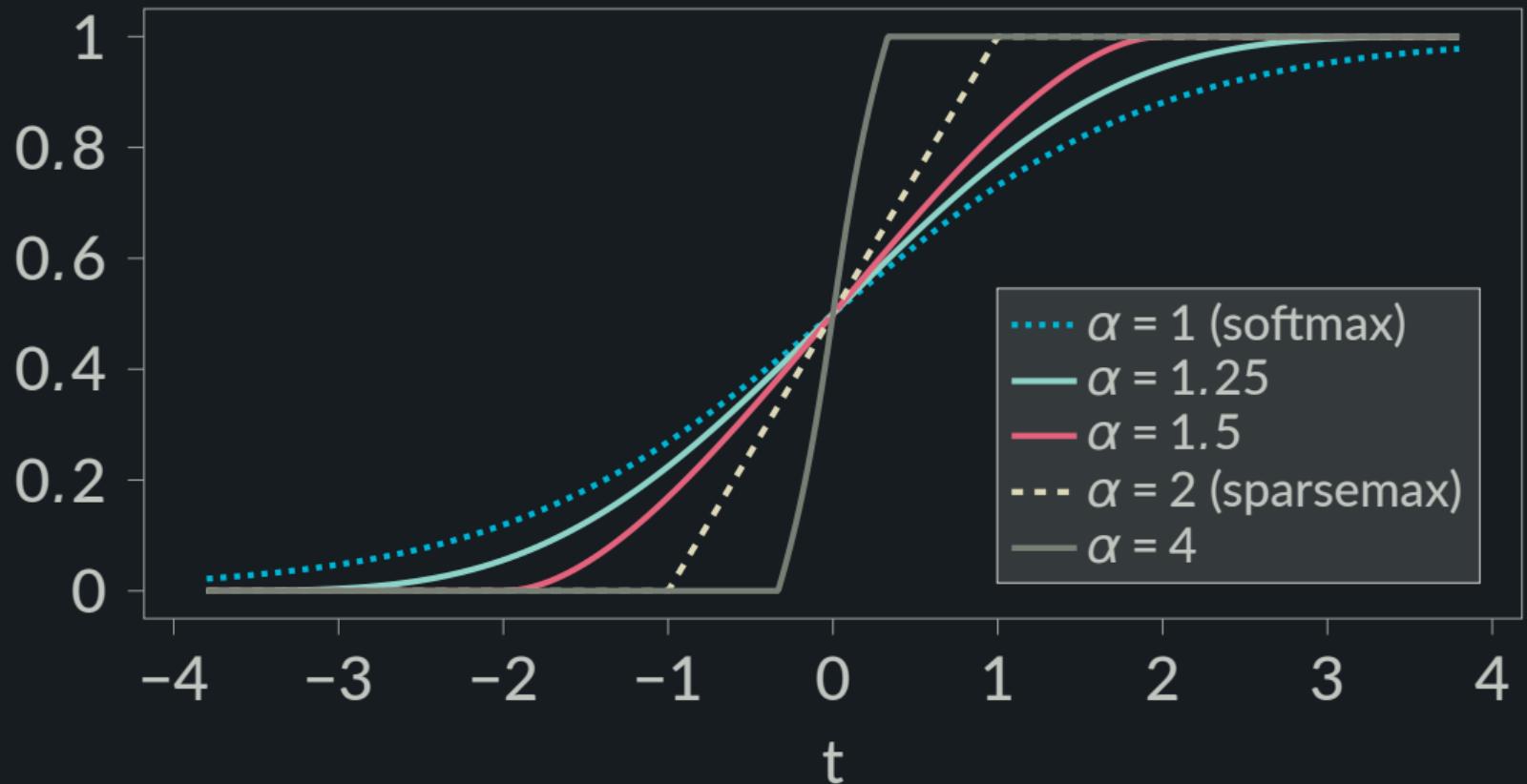
Key result: can be sparse for  $\alpha > 1$ , propensity for sparsity increases with  $\alpha$ .











# Learning $\alpha$

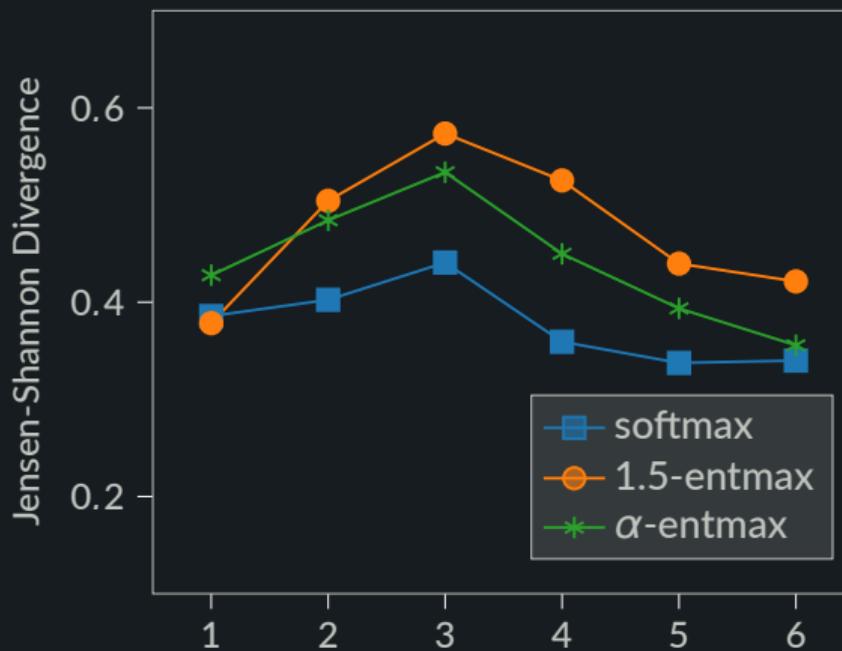
# Learning $\alpha$

Key contribution:

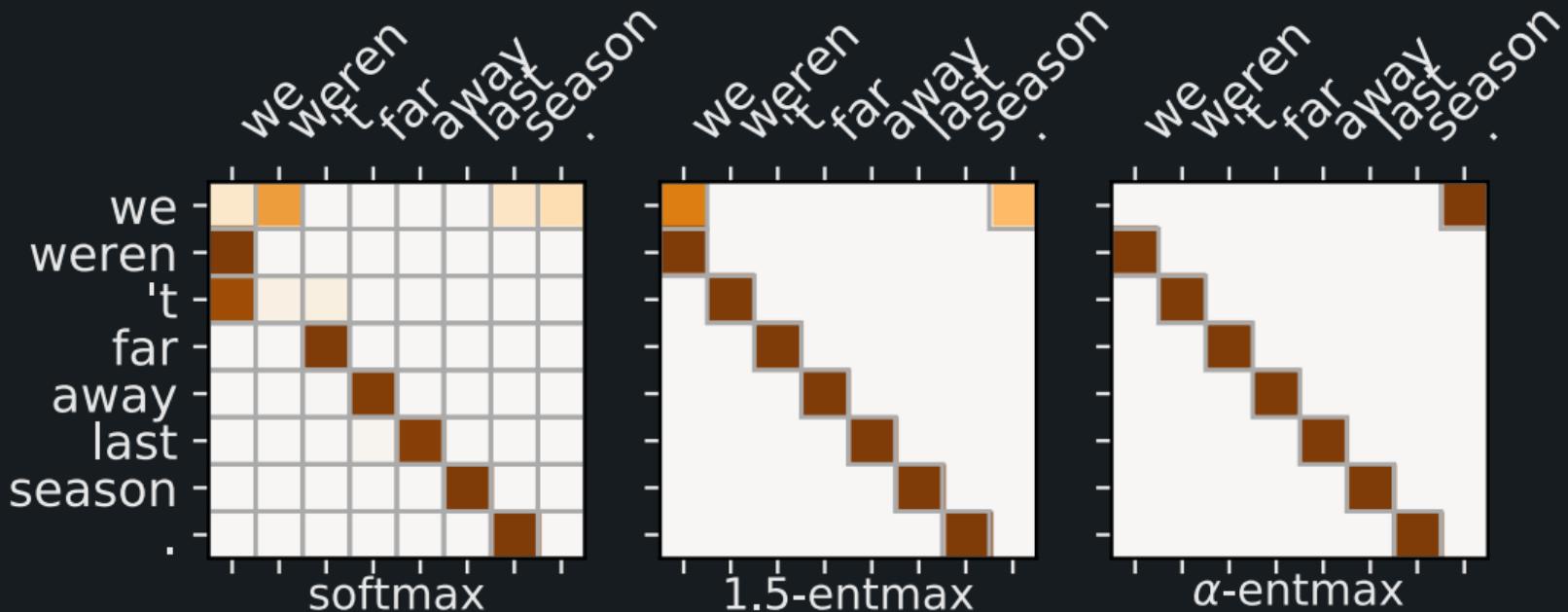
a closed-form expression for  $\frac{\partial \alpha\text{-entmax}(\mathbf{z})}{\partial \alpha}$



# Head Diversity per Layer

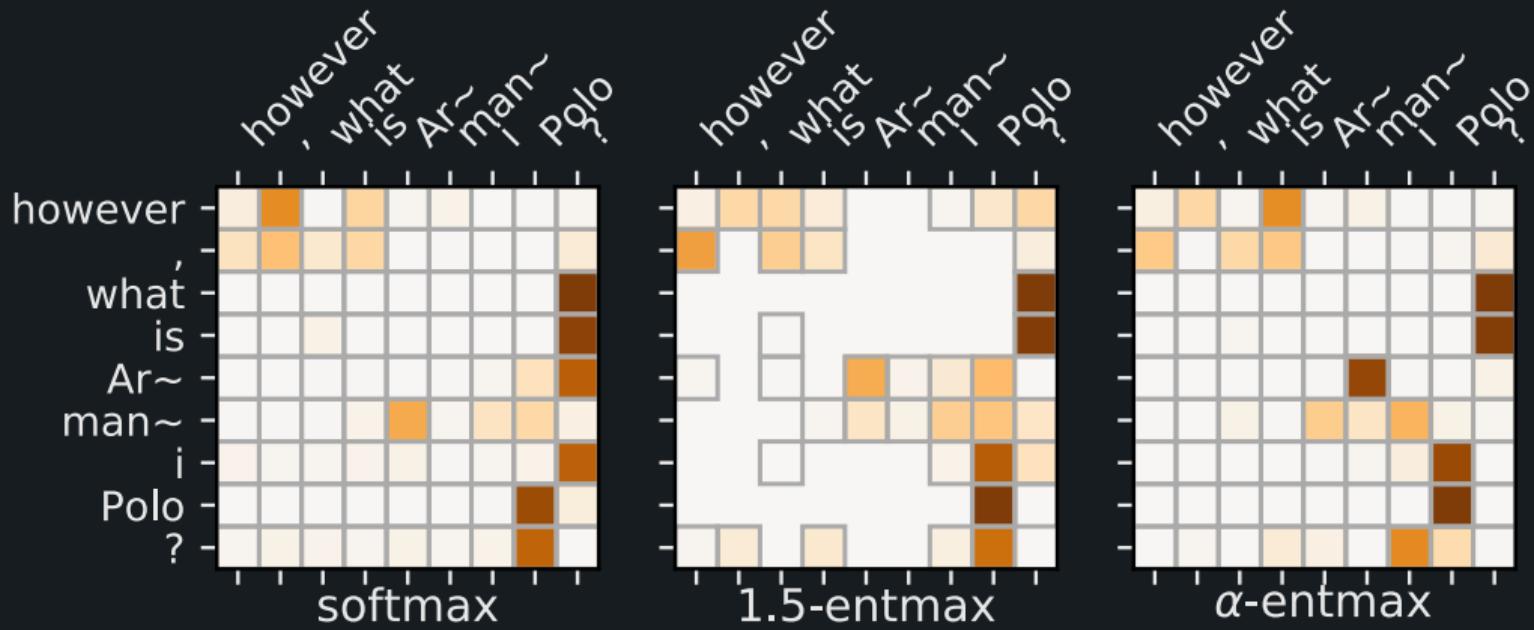


# Previous Position Head



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

# Interrogation-Detecting Head



Learned  $\alpha = 1.05$ .

# Table of Contents

A Simple and Effective Approach to Automatic Post-Editing with Transfer Learning

Adaptively Sparse Transformers

Efficient Marginalization of Discrete and Structured Latent Variables via Sparsity

Proposed Work

# Latent Variable Models

Latent variable  $z$  can be

# Latent Variable Models

Latent variable  $z$  can be continuous



Source: Bouges et al., 2013

# Latent Variable Models

Latent variable  $z$  can be **continuous**, **discrete**



# Latent Variable Models

Latent variable  $z$  can be continuous, discrete, or structured



Source: Liu et al., 2015

# Training Discrete or Structured Latent Variable Models

Latent variable  $z$  can be

# Training Discrete or Structured Latent Variable Models

Latent variable  $z$  can be discrete



# Training Discrete or Structured Latent Variable Models

Latent variable  $z$  can be **discrete** or **structured**



# Training Discrete or Structured Latent Variable Models

Latent variable  $z$  can be **discrete** or **structured**

$\pi(z|x, \theta)$ : distribution over possible  $z$



# Training Discrete or Structured Latent Variable Models

Latent variable  $z$  can be **discrete** or **structured**

$\pi(z|x, \theta)$ : distribution over possible  $z$



# Training Discrete or Structured Latent Variable Models

Latent variable  $z$  can be **discrete** or **structured**

$\pi(z|x, \theta)$ : distribution over possible  $z$

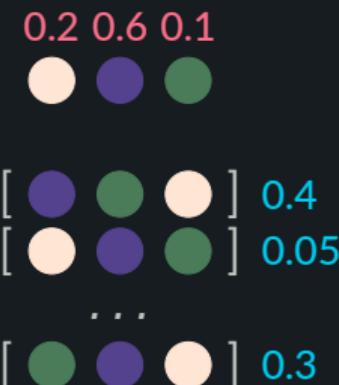


# Training Discrete or Structured Latent Variable Models

Latent variable  $z$  can be **discrete** or **structured**

$\pi(z|x, \theta)$ : distribution over possible  $z$

$\ell(x, z; \theta)$ : downstream loss: ELBO, Log-Likelihood, (...)



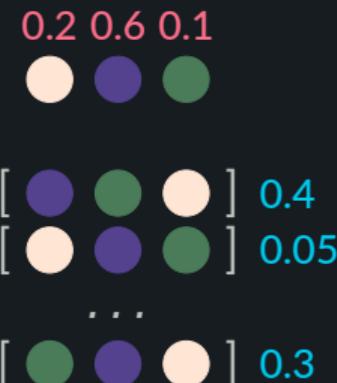
# Training Discrete or Structured Latent Variable Models

Latent variable  $z$  can be **discrete** or **structured**

$\pi(z|x, \theta)$ : distribution over possible  $z$

$\ell(x, z; \theta)$ : downstream loss: ELBO, Log-Likelihood, (...)

To train, we need to compute the following expectation:



# Training Discrete or Structured Latent Variable Models

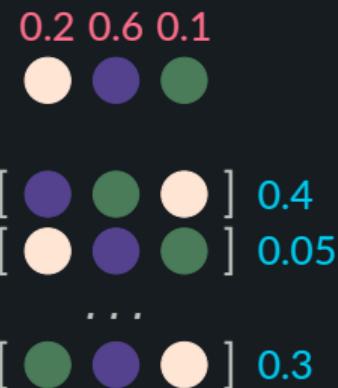
Latent variable  $z$  can be **discrete** or **structured**

$\pi(z|x, \theta)$ : distribution over possible  $z$

$\ell(x, z; \theta)$ : downstream loss: ELBO, Log-Likelihood, (...)

To train, we need to compute the following expectation:

$$\mathcal{L}_x(\theta) = \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta)$$

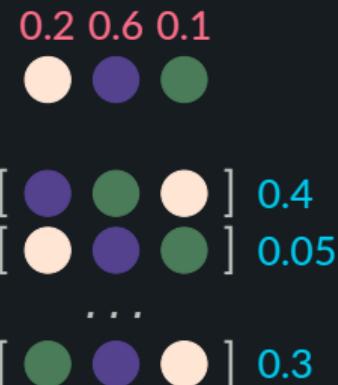


# Training Discrete or Structured Latent Variable Models

Latent variable  $z$  can be **discrete** or **structured**

$\pi(z|x, \theta)$ : distribution over possible  $z$

$\ell(x, z; \theta)$ : downstream loss: ELBO, Log-Likelihood, (...)



To train, we need to compute the following expectation:

$$\mathcal{L}_x(\theta) = \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta)$$

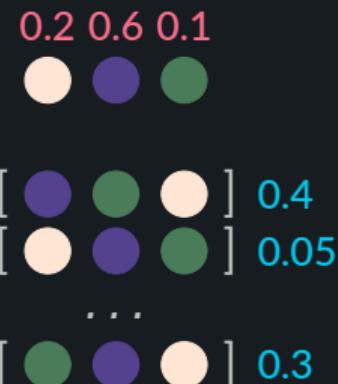
If  $\mathcal{Z}$  is **large**, this sum can get very expensive due to  $\ell(x, z; \theta)$ ! 🍷

# Training Discrete or Structured Latent Variable Models

Latent variable  $z$  can be **discrete** or **structured**

$\pi(z|x, \theta)$ : distribution over possible  $z$

$\ell(x, z; \theta)$ : downstream loss: ELBO, Log-Likelihood, (...)



To train, we need to compute the following expectation:

$$\mathcal{L}_x(\theta) = \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta)$$

If  $\mathcal{Z}$  is **combinatorial**, this can be intractable to compute!



# Current Solutions

If  $\mathcal{Z}$  is large, exact gradient computation is prohibitive

# Current Solutions

If  $\mathcal{Z}$  is large, exact gradient computation is prohibitive

One option: SFE (aka REINFORCE) → unbiased but high variance

# Current Solutions

If  $\mathcal{Z}$  is large, exact gradient computation is prohibitive

One option: SFE (aka REINFORCE) → unbiased but high variance

Another option: Gumbel-Softmax → continuous relaxation, biased estimation

# Current Solutions

If  $\mathcal{Z}$  is large, exact gradient computation is prohibitive

One option: SFE (aka REINFORCE) → unbiased but high variance

Another option: Gumbel-Softmax → continuous relaxation, biased estimation

New option: use sparsity! 

# Current Solutions

If  $\mathcal{Z}$  is large, exact gradient computation is prohibitive

One option: SFE (aka REINFORCE) → unbiased but high variance

Another option: Gumbel-Softmax → continuous relaxation, biased estimation

New option: use sparsity! 

no need for sampling → no variance

# Current Solutions

If  $\mathcal{Z}$  is large, exact gradient computation is prohibitive

One option: SFE (aka REINFORCE) → unbiased but high variance

Another option: Gumbel-Softmax → continuous relaxation, biased estimation

New option: use sparsity! 

no need for sampling → no variance

no relaxation into the continuous space

# Taking a step back...

Does the expectation over possible  $z$  need to be expensive?

# Taking a step back...

Does the expectation over possible  $z$  need to be expensive?

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta) \\ &= \pi(z_1|x, \theta) \ell(x, z_1; \theta) + \pi(z_2|x, \theta) \ell(x, z_2; \theta) + \dots \\ &\quad + \pi(z_i|x, \theta) \ell(x, z_i; \theta) + \dots + \pi(z_N|x, \theta) \ell(x, z_N; \theta)\end{aligned}$$

# Taking a step back...

Does the expectation over possible  $z$  need to be expensive?

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta) \\ &= \pi(z_1|x, \theta) \ell(x, z_1; \theta) + \pi(z_2|x, \theta) \ell(x, z_2; \theta) + \dots \\ &\quad + \pi(z_i|x, \theta) \ell(x, z_i; \theta) + \dots + \pi(z_N|x, \theta) \ell(x, z_N; \theta)\end{aligned}$$

Usually we normalize  $\pi$  with softmax  $\propto \exp(s) \Rightarrow \pi(z_i|x, \theta) > 0$

# Sparse normalizers

We use **sparsemax**, **top- $k$  sparsemax** and **SparseMAP** to allow efficient marginalization

# Sparse normalizers

We use **sparsemax**, **top- $k$  sparsemax** and **SparseMAP** to allow efficient marginalization

These functions are able to assign **probabilities of exactly zero!**

# Sparse normalizers

We use **sparsemax**, **top- $k$  sparsemax** and **SparseMAP** to allow efficient marginalization

These functions are able to assign **probabilities of exactly zero!**

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta) \\ &= \pi(z_1|x, \theta) \ell(x, z_1; \theta) + \underbrace{\pi(z_2|x, \theta)}_{=0} \ell(x, z_2; \theta) + \dots \\ &\quad + \pi(z_i|x, \theta) \ell(x, z_i; \theta) + \dots + \underbrace{\pi(z_N|x, \theta)}_{=0} \ell(x, z_N; \theta)\end{aligned}$$

# Sparse normalizers

We use **sparsemax**, **top- $k$  sparsemax** and **SparseMAP** to allow efficient marginalization

These functions are able to assign **probabilities of exactly zero!**

$$\begin{aligned}
 \mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x, \theta) \ell(x, z; \theta) \\
 &= \pi(z_1|x, \theta) \ell(x, z_1; \theta) + \underbrace{\pi(z_2|x, \theta) \ell(x, z_2; \theta)}_{=0} + \dots \\
 &\quad + \pi(z_i|x, \theta) \ell(x, z_i; \theta) + \dots + \underbrace{\pi(z_N|x, \theta) \ell(x, z_N; \theta)}_{=0}
 \end{aligned}$$

No need for computing  $\ell(x, z; \theta)$  for all  $z \in \mathcal{Z}$ !

# Results

We test our methods for models with discrete latent variables,

# Results

We test our methods for models with discrete latent variables,

- Semi-Supervised VAE

# Results

We test our methods for models with discrete latent variables,

- Semi-Supervised VAE
- Emergent communication

# Results

We test our methods for models with discrete latent variables,

- Semi-Supervised VAE
- Emergent communication

but also in models with an exponentially large set of  $\mathcal{Z}$ ,

# Results

We test our methods for models with discrete latent variables,

- Semi-Supervised VAE
- Emergent communication

but also in models with an exponentially large set of  $\mathcal{Z}$ ,

- Bit-vector VAE

# Results

We test our methods for models with discrete latent variables,

- Semi-Supervised VAE
- Emergent communication

but also in models with an exponentially large set of  $\mathcal{Z}$ ,

- Bit-vector VAE

Our methods are top-performers and efficient!

# Table of Contents

A Simple and Effective Approach to Automatic Post-Editing with Transfer Learning

Adaptively Sparse Transformers

Efficient Marginalization of Discrete and Structured Latent Variables via Sparsity

Proposed Work

# Proposed work: key concepts of this thesis

- Transparency
- Sparsity
- Compactness
- Data efficiency

**Les artistes sont demunis**  
  
**The artists don't have any money**

# Proposed work: key concepts of this thesis

- Transparency
- Sparsity
- Compactness
- Data efficiency
- Inspired by, and data efficient like APE

**Les artistes sont demunis**  
  
**The artists don't have any money**

# Proposed work: key concepts of this thesis

- Transparency
- Sparsity
- Compactness
- Data efficiency
- Inspired by, and data efficient like APE
- Transparent like sparse transformers

**Les artistes sont demunis**  
  
**The artists don't have any money**

# Proposed work: key concepts of this thesis

- Transparency
- Sparsity
- Compactness
- Data efficiency
- Inspired by, and data efficient like APE
- Transparent like sparse transformers
- Compact thanks to latent variables



**Les artistes sont demunis**

The artists don't have any money

# Proposed work: key concepts of this thesis

- Transparency
- Sparsity
- Compactness
- Data efficiency
- Inspired by, and data efficient like APE
- Transparent like sparse transformers
- Compact thanks to latent variables
- Using sparsity as a training method

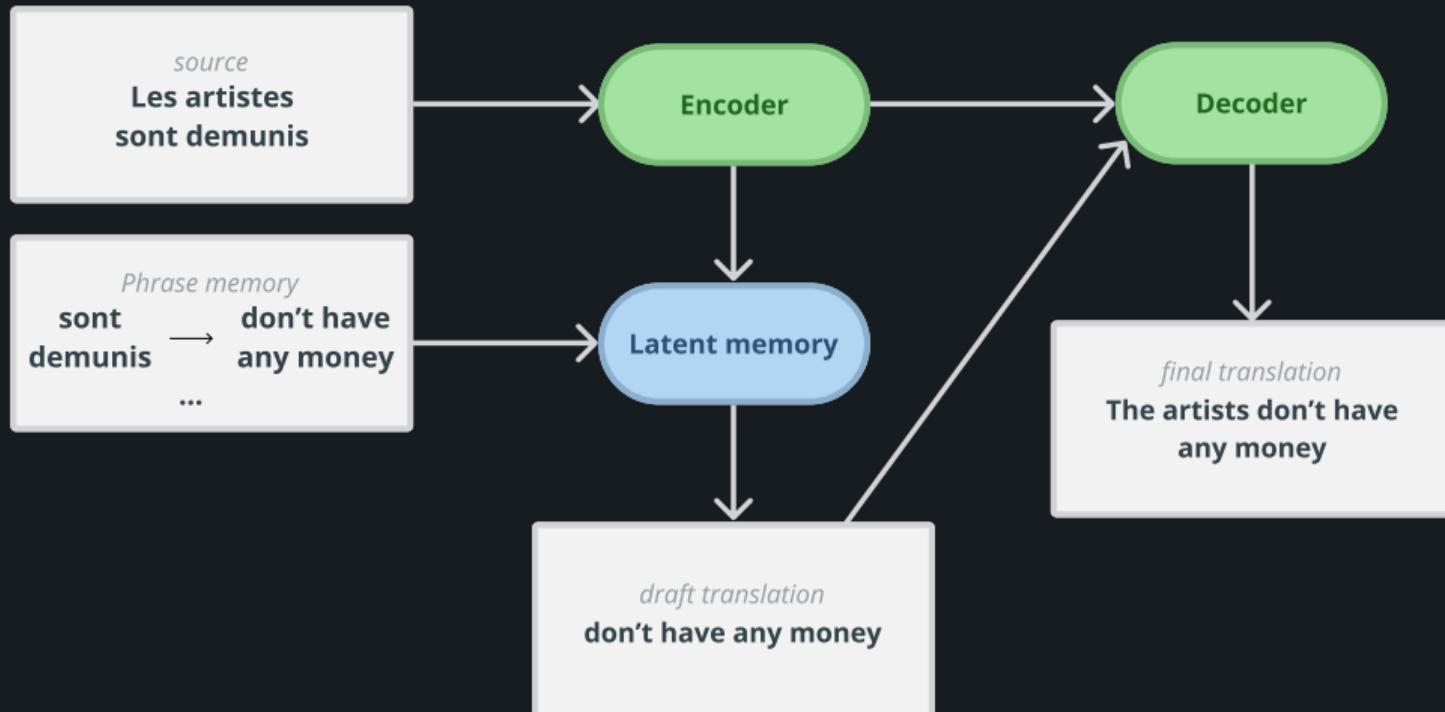


**Les artistes sont demunis**

**The artists don't have any money**

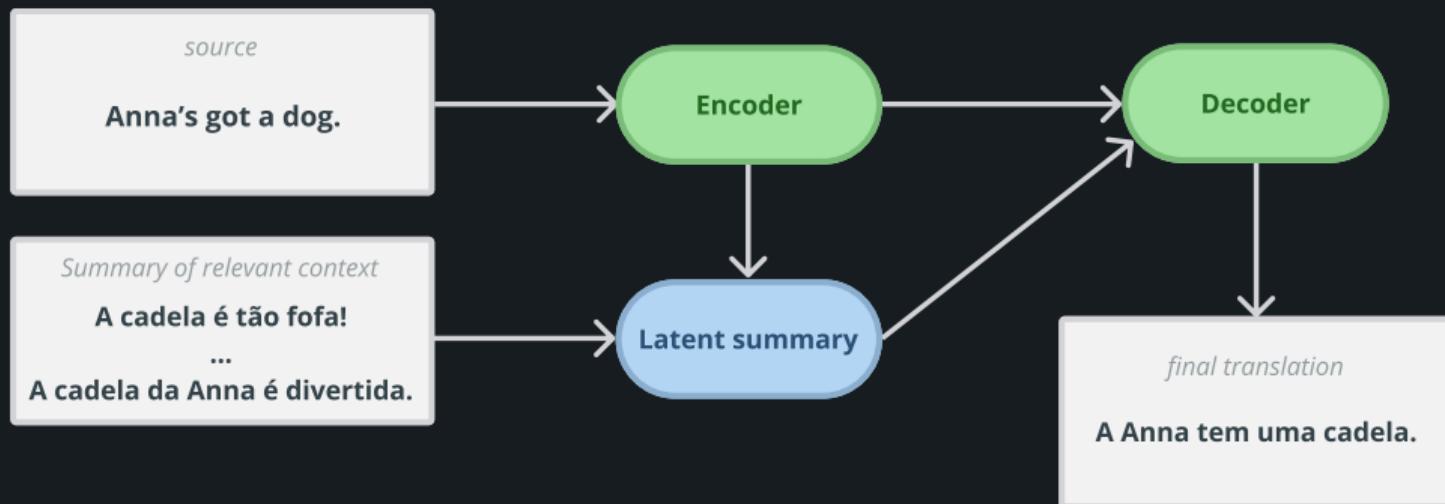
**Idea: create a latent, draft translation**

# Idea: create a latent, draft translation



# Another Idea: using latent summaries in NMT

# Another Idea: using latent summaries in NMT



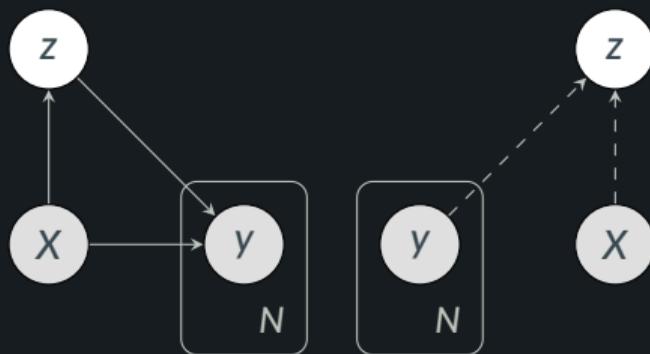
# References I

-  Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio (2015). "Neural machine translation by jointly learning to align and translate". In: *Proc. of ICLR*.
-  Bouges, Pierre, Thierry Chateau, Christophe Blanc, and Gaëlle Loosli (Dec. 2013). "Handling missing weak classifiers in boosted cascade: application to multiview and occluded face detection". In: *EURASIP Journal on Image and Video Processing* 2013, p. 55. DOI: 10.1186/1687-5281-2013-55.
-  Brown, Peter F, Stephen A Della Pietra, Vincent J Della Pietra, and Robert L Mercer (1993). "The mathematics of statistical machine translation: Parameter estimation". In: *Computational linguistics* 19.2, pp. 263–311.
-  Correia, Gonçalo M, Vlad Niculae, and André FT Martins (2019). "Adaptively sparse transformers". In: *Proc. EMNLP*.
-  Correia, Gonçalo M. and André F. T. Martins (July 2019). "A Simple and Effective Approach to Automatic Post-Editing with Transfer Learning". In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, pp. 3050–3056. DOI: 10.18653/v1/P19-1292. URL: <https://www.aclweb.org/anthology/P19-1292>.
-  Correia, Gonçalo M., Vlad Niculae, Wilker Aziz, and André F. T. Martins (2020). "Efficient Marginalization of Discrete and Structured Latent Variables via Sparsity". In: *Proc. NeurIPS*. URL: <https://arxiv.org/abs/2007.01919>.
-  Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2019). "BERT: Pre-training of deep bidirectional transformers for language understanding". In: *Proc. NAACL-HLT*.

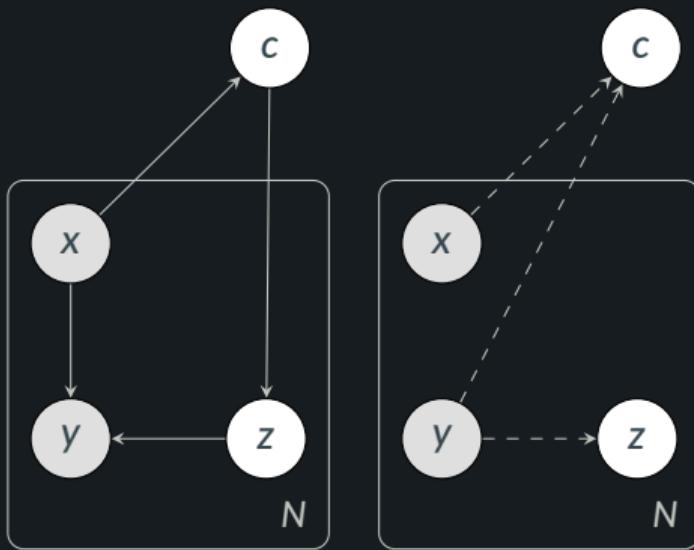
# References II

-  Liu, Ziwei, Ping Luo, Xiaogang Wang, and Xiaoou Tang (Dec. 2015). "Deep Learning Face Attributes in the Wild". In: *Proceedings of International Conference on Computer Vision (ICCV)*.
-  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, André FT Martins, Mathieu Blondel, and Claire Cardie (2018). "SparseMAP: Differentiable sparse structured inference". In: *Proc. of ICML*.
-  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*.
-  Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz 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*.

# Using latent summaries in NMT



# Multi-domain NMT



# Document-level NMT

