

Learnable Sparsity and Weak Supervision for Data-efficient, Transparent, and Compact Neural Models

Gonçalo M. Correia

Jury: André Martins, Mário Figueiredo, Ivan Titov, Wilker Aziz, Isabel Trancoso

Deep learning successes

Deep learning successes

- Subset of machine learning that uses **neural networks**

Deep learning successes

- Subset of machine learning that uses **neural networks**
- Powerful tool for learning representations of any data

Deep learning successes

- Subset of machine learning that uses **neural networks**
- Powerful tool for learning representations of any data
- Remarkable results

Deep learning successes

- Subset of machine learning t
- Powerful tool for learning re
- Remarkable results

A robot wrote this entire article. Are you scared yet, human?

GPT-3



Deep learning successes

≡ WIRED

SUBSCRIBE

WILL KNIGHT BUSINESS MAR 29, 2021 7:00 AM

This AI Can Generate Convincing Text—and Anyone Can Use It

The makers of Eleuther hope it will be an open source alternative to GPT-3, the well-known language program from OpenAI.

A robot wrote this entire article. Are you scared yet, human?

GPT-3

The
Guardian
News website of the year

Deep learning successes

≡ WIRED

SUBSCRIBE

WILL KNIGHT

BUSINESS MAR 29, 2021 7:00 AM

This AI Can Generate Convincing Text—Anyone Can Use It

The makers of Eleuther hope it's a source alternative to GPT-3, the language program from OpenAI

A robot wrote this entire article. Are you scared yet, human?

GPT-3

The

SCIENCE

Danny's workmate is called GPT-3. You've probably read its work without realising it's an AI

ABC Science / By technology reporter James Purtill

Posted Sat 28 May 2022 at 7:30pm

- Subs
- Pow
- Rem

Deep learning successes

≡ WIRED

SUBSCRIBE

WILL KNIGHT

BUSINESS MAR 29, 2021 7:00 AM

This AI Can Generate Convincing Text—[SCIENCE](#)

Forbes

INNOVATION

Are AI Systems About To Outperform Humans?

A robot wrote this entire article. Are you scared yet, human?

CDT-2

The

arkmate is called
e probably read its
ut realising it's an AI

hnology reporter [James Purtill](#)

Posted Sat 28 May 2022 at 7:30pm

Deep learning successes Artificial intelligence beats eight world champions at bridge

Victory marks milestone for AI as bridge requires more human skills than other strategy games

INNOVATION

Are AI Systems About To Outperform Humans?

Deep learning successes

robot wrote this entire article. Are you scared yet, human?

DT-2

The

arkmate is called

'e probably read its
ut realising it's an AI

hnology reporter [James Purtill](#)

Posted Sat 28 May 2022 at 7:30pm

Deep learning successes

Artificial intelligence beats eight world champions at bridge

Victory marks milestone for AI
bridge requires more human skill than other strategy games

INNOVATION

Are AI Systems About To Outperform Humans?

Posted Sat

robot wrote this

AI 'outperforms' doctors diagnosing breast cancer



Fergus Walsh
Medical correspondent
@BBCFergusWalsh

Deep learning limitations and drawbacks

Deep learning limitations and drawbacks

- Requires a lot of data

Deep learning limitations and drawbacks

- Requires a lot of data
- Hard to understand and interpret reasons behind decisions

Deep learning limitations and drawbacks

- Requires a lot of data
- Hard to understand and interpret reasons behind decisions
- Requires a lot of computation

Deep learning limitations and drawbacks

- Requires a lot of data
- Hard to understand and interpret
- Requires a lot of computation

The screenshot shows a white rectangular box with a dark border. At the top left is the WIRED logo. To its right is a blue 'SUBSCRIBE' button. Below the logo, the title 'AI Can Do Great Things—if It Doesn't Burn the Planet' is displayed in large, bold, black capital letters. Underneath the title is a paragraph of text: 'The computing power required for AI landmarks, such as recognizing images and defeating humans at Go, increased 300,000-fold from 2012 to 2018.' To the right of the text, the word 'sions' is partially visible, suggesting it's part of a larger sentence cut off by the image's edge.

= WIRED

SUBSCRIBE

AI Can Do Great Things—if It Doesn't Burn the Planet

The computing power required for AI landmarks, such as recognizing images and defeating humans at Go, increased 300,000-fold from 2012 to 2018.

sions

Deep learning limitations and drawbacks

Forbes

AI

- Req
 - Hard
 - Req
- ## Overcoming AI's Transparency Paradox

≡ WIRED

SUBSCRIBE

AI Can Do Great Things—if It Doesn't Burn the Planet

The computing power required for AI landmarks, such as recognizing images and defeating humans at Go, increased 300,000-fold from 2012 to 2018.

sions

Deep learning limitations and drawbacks

- Req
- Hard
- Req

Forbes

AI

Overcoming Transparency Paradox



≡ WIRED

SUBSCRIBE

AI Can Do Great Things—if It Doesn't Burn the Planet

The computing power required for AI

Harvard
Business
Review

AI Can Outperform Doctors. So Why Don't Patients Trust It?

by Chiara Longoni and Carey K. Morewedge

Deep learning limitations and drawbacks

≡ WIRED

SUBSCRIBE

Forbes

AI Can Do Great Things—if It Doesn't Burn the Planet

Computing power required for AI

rd
ess
N

'Dangerous' AI offers to write fake news

By Jane Wakefield
Technology reporter

Outperform
So Why Don't
Trust It?

oni and Carey K. Morewedge

Key concepts of this thesis

**Learnable Sparsity and Weak Supervision
for Data-efficient, Transparent, and Compact
Neural Models**

Key concepts of this thesis

Learnable Sparsity and Weak Supervision
for **Data-efficient**, Transparent, and Compact
Neural Models

Key concepts of this thesis

Learnable Sparsity and Weak Supervision
for **Data-efficient**, **Transparent**, and **Compact**
Neural Models

Key concepts of this thesis

Learnable Sparsity and Weak Supervision
for Data-efficient, Transparent, and Compact
Neural Models

Key concepts of this thesis

**Learnable Sparsity and Weak Supervision
for Data-efficient, Transparent, and Compact
Neural Models**

Key concepts of this thesis

**Learnable Sparsity and Weak Supervision
for Data-efficient, Transparent, and Compact
Neural Models**

Published work of this thesis

Published work of this thesis

- Automatic Post-Editing using **weak supervision** for **data-efficiency** (ACL)

Published work of this thesis

- Automatic Post-Editing using **weak supervision** for **data-efficiency** (ACL)
- Letting transformer **learn sparsity** of its attentions for **transparency** (EMNLP)

Published work of this thesis

- Automatic Post-Editing using **weak supervision** for **data-efficiency** (ACL)
- Letting transformer **learn sparsity** of its attentions for **transparency** (EMNLP)
- General strategy for efficiently training discrete latent variable models, to have **compactness** (NeurIPS)

Table of Contents

A Simple and Effective Approach to APE with Transfer Learning

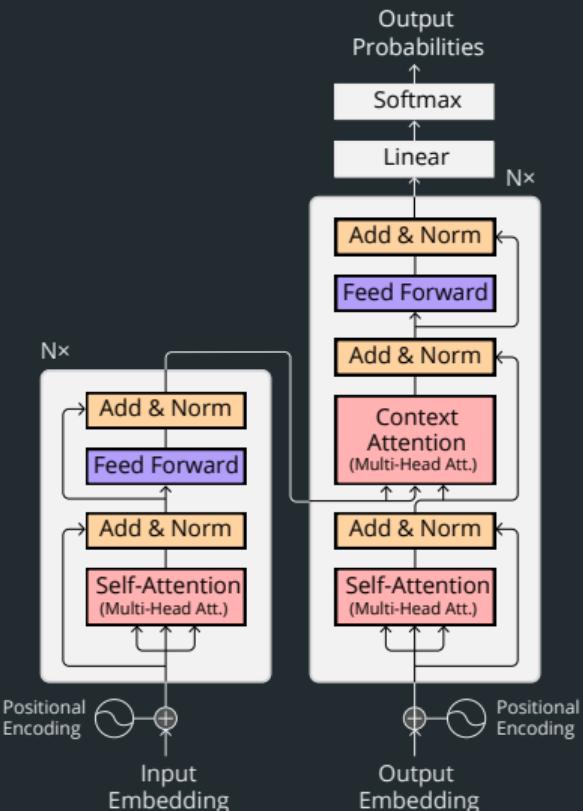
Adaptively Sparse Transformers

Efficient Marg. of Discrete Latent Variables via Sparsity

Conclusions

A bit of context on transformers

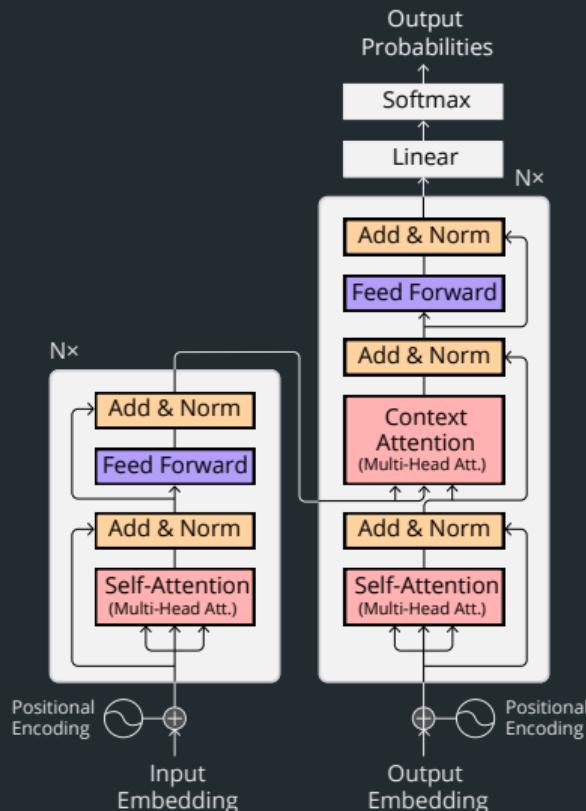
What if... Attention is all you need?



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!

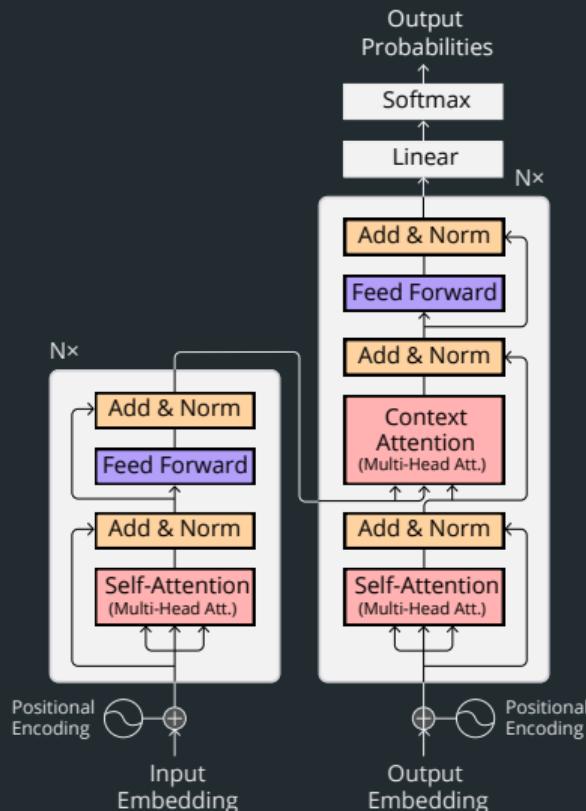


A bit of context on transformers

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

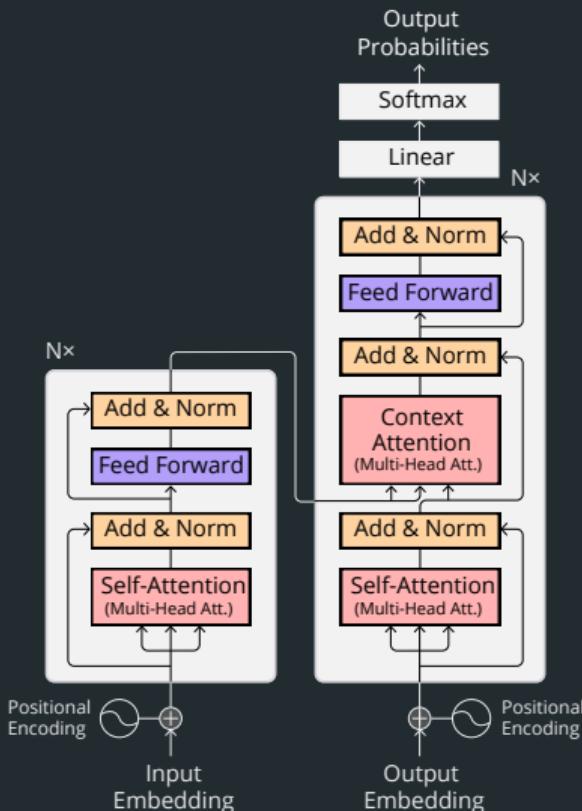


A bit of context on transformers

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)

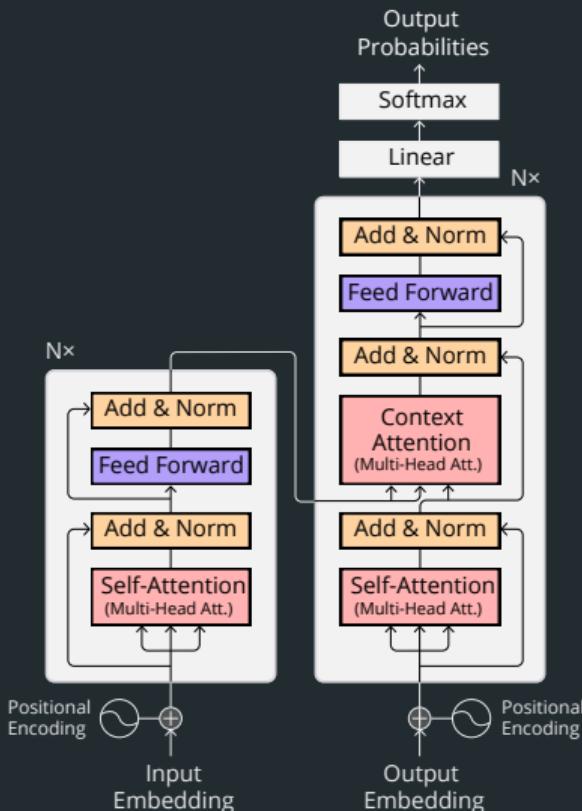


A bit of context on transformers

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

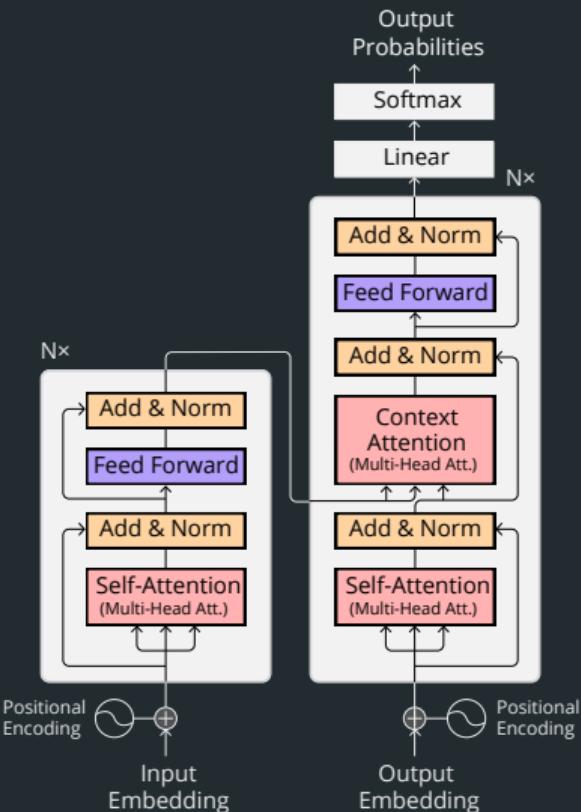


A bit of context on transformers

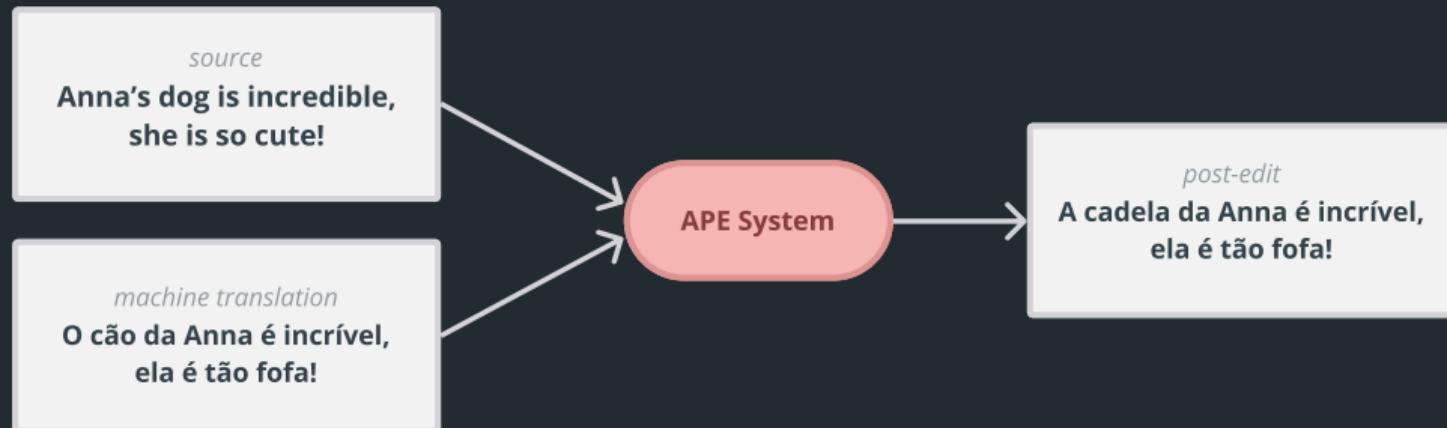
What if... Attention is all you need?

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

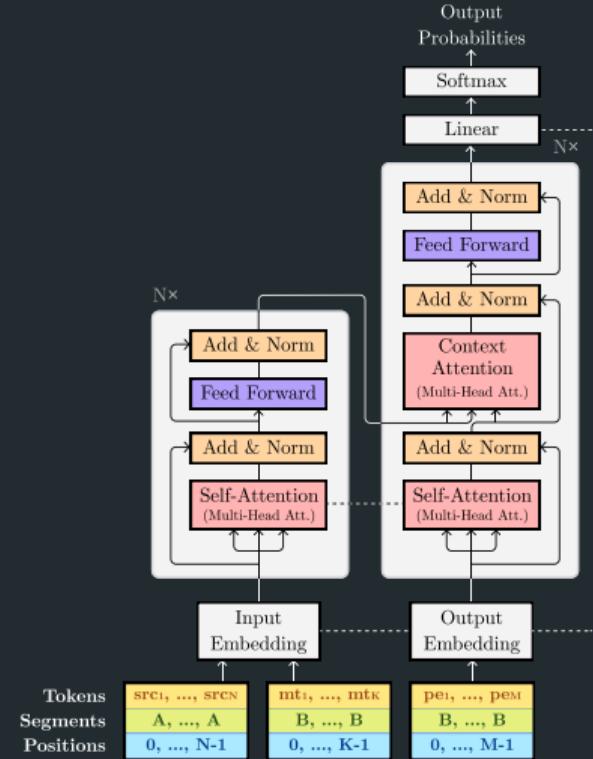
- In place of the RNNs, use self-attention
- Do this with multiple heads (i.e. attention mechanisms in parallel)
- ... and do it through several layers
- Inspiration for big general-purpose models like BERT!



What is APE?

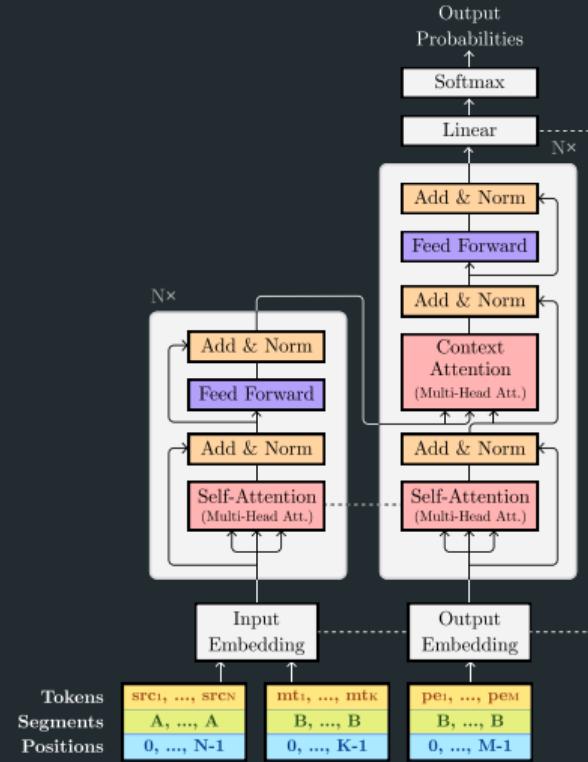


BERT for APE



BERT for APE

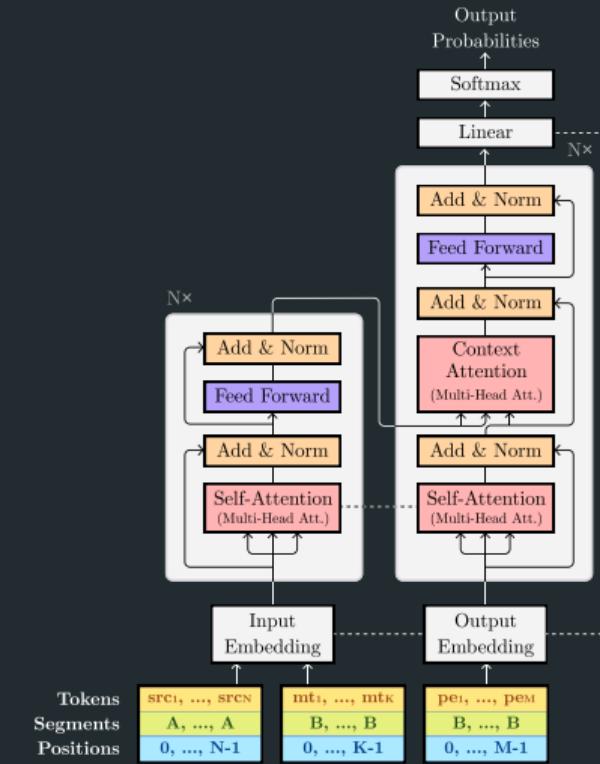
Key idea: Use BERT to do APE



BERT for APE

Key idea: Use BERT to do APE

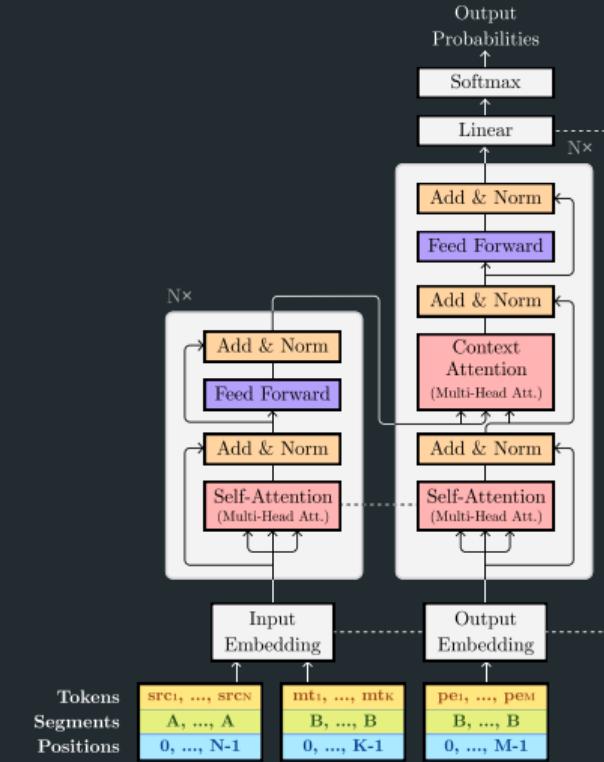
- Prior to this work, BERT was mainly used for simple classification tasks



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

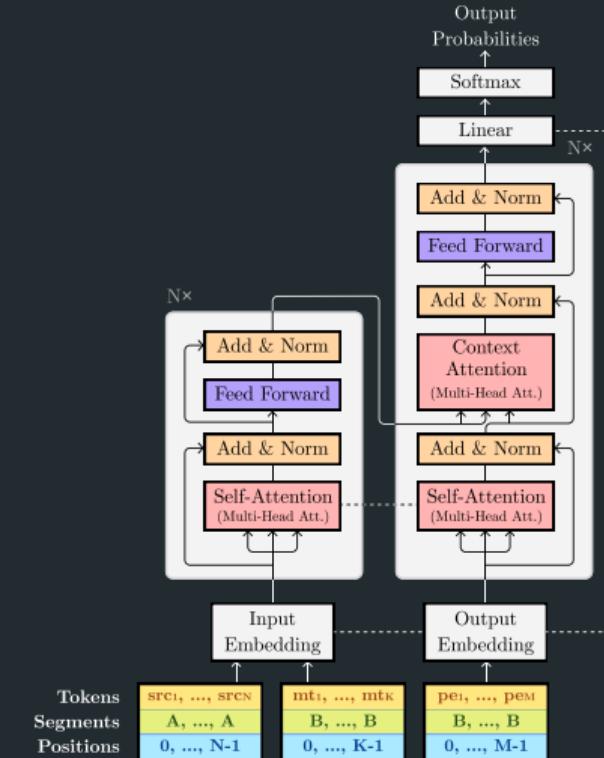


Table of Contents

A Simple and Effective Approach to APE with Transfer Learning

Adaptively Sparse Transformers

Efficient Marg. of Discrete Latent Variables via Sparsity

Conclusions

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**.

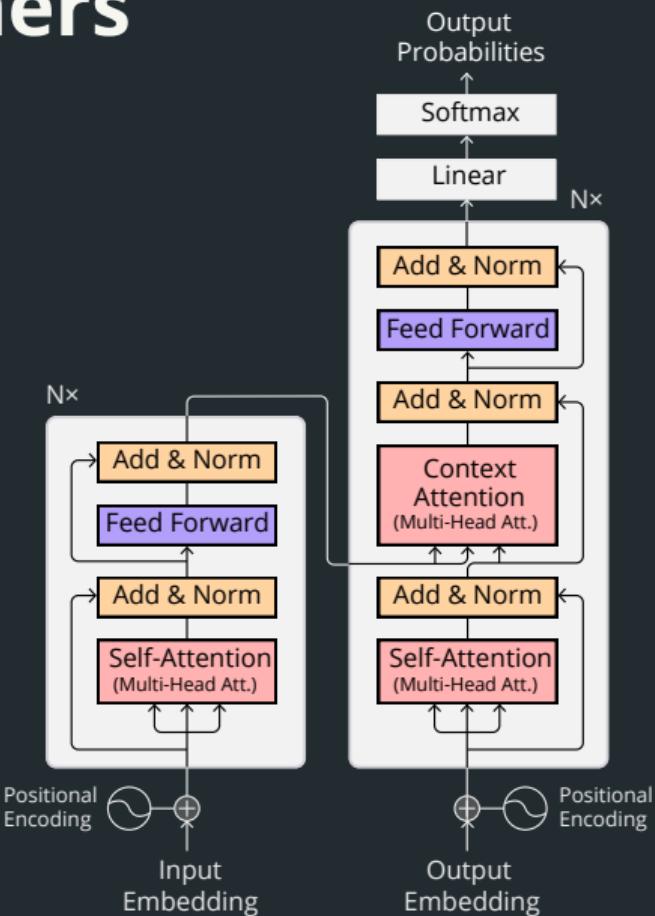
Our solution is to bet on **sparsity**:

- for interpretability
- for discovering linguistic structure
- for efficiency

Transformers

In each attention head:

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



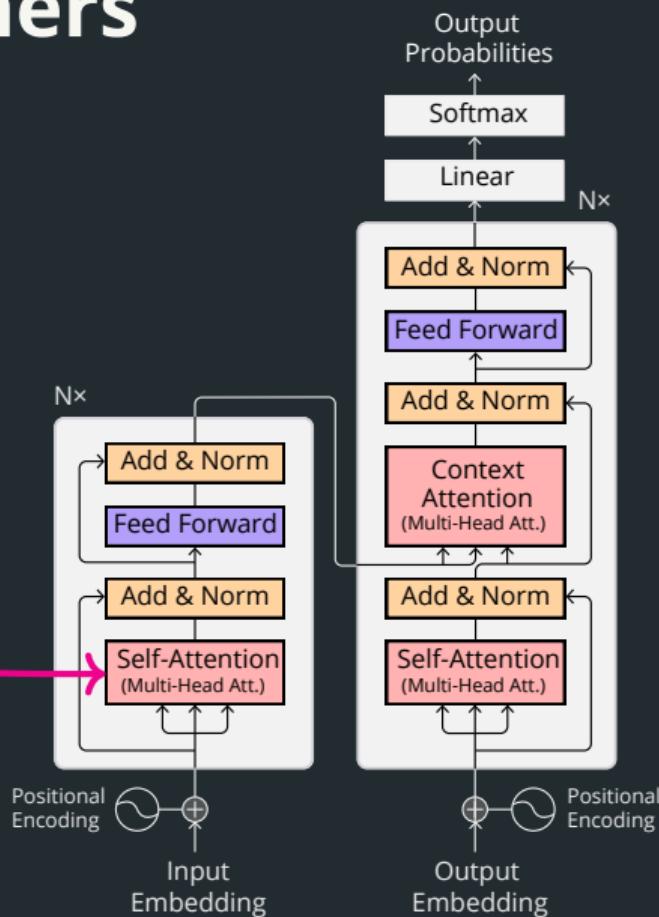
Transformers

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



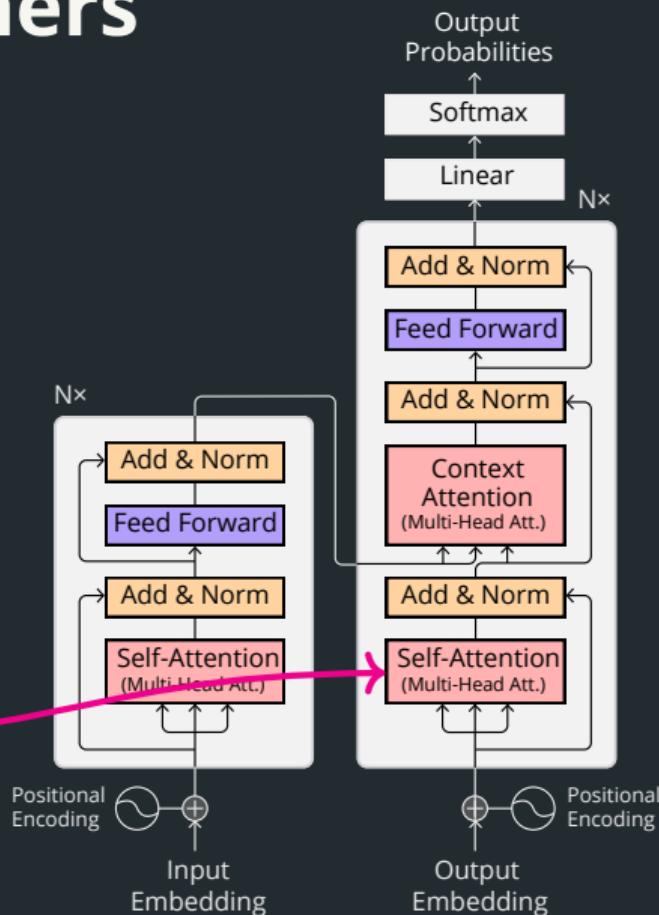
Transformers

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



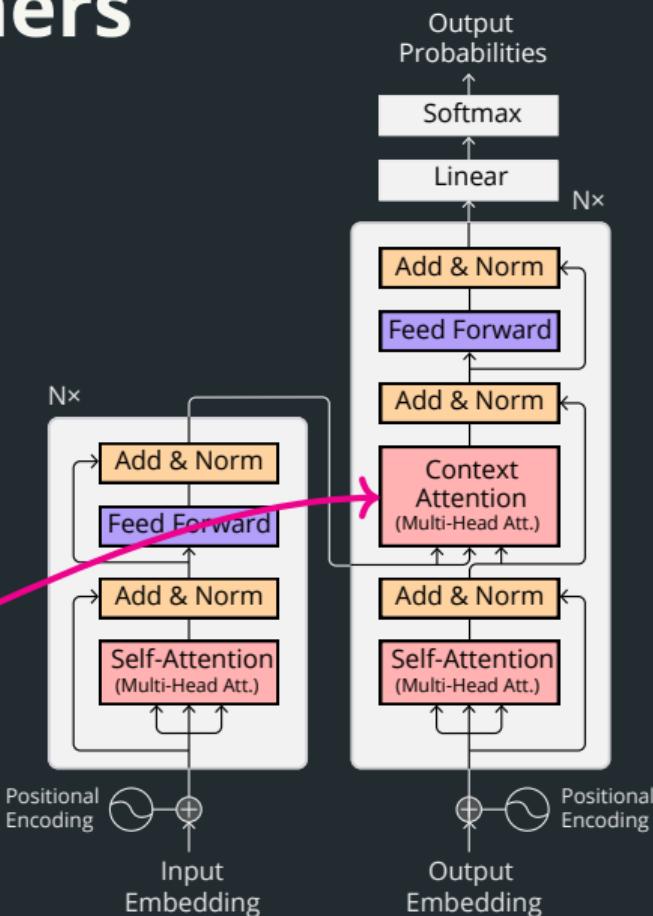
Transformers

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 α ! 

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

α -Entmax

Parametrized by $\alpha \geq 0$:

α -Entmax

Parametrized by $\alpha \geq 0$:

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

α -Entmax

Parametrized by $\alpha \geq 0$:

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

α -Entmax

Parametrized by $\alpha \geq 0$:

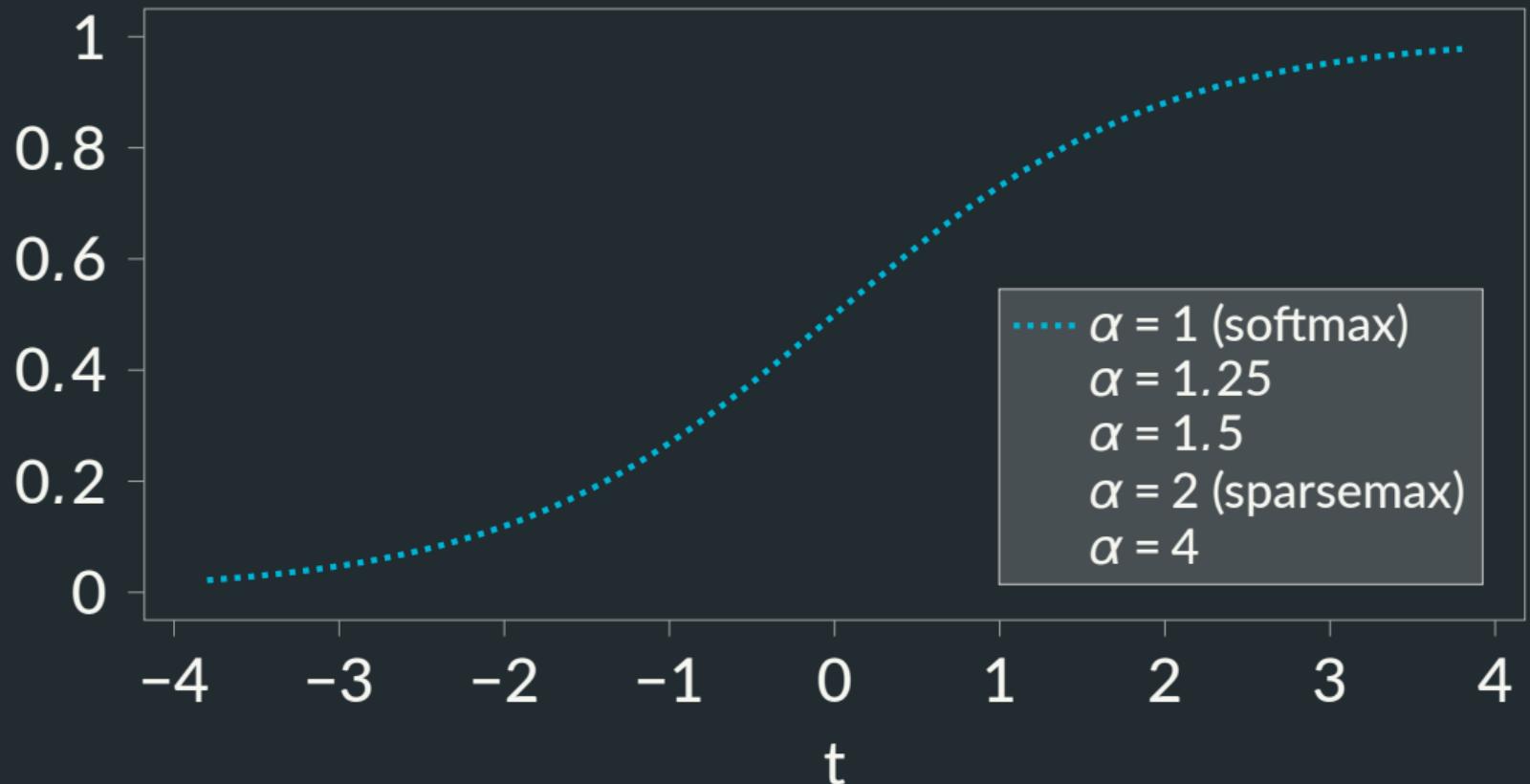
- **Argmax** corresponds to $\alpha \rightarrow \infty$
- **Softmax** amounts to $\alpha \rightarrow 1$
- **Sparsemax** amounts to $\alpha = 2$.

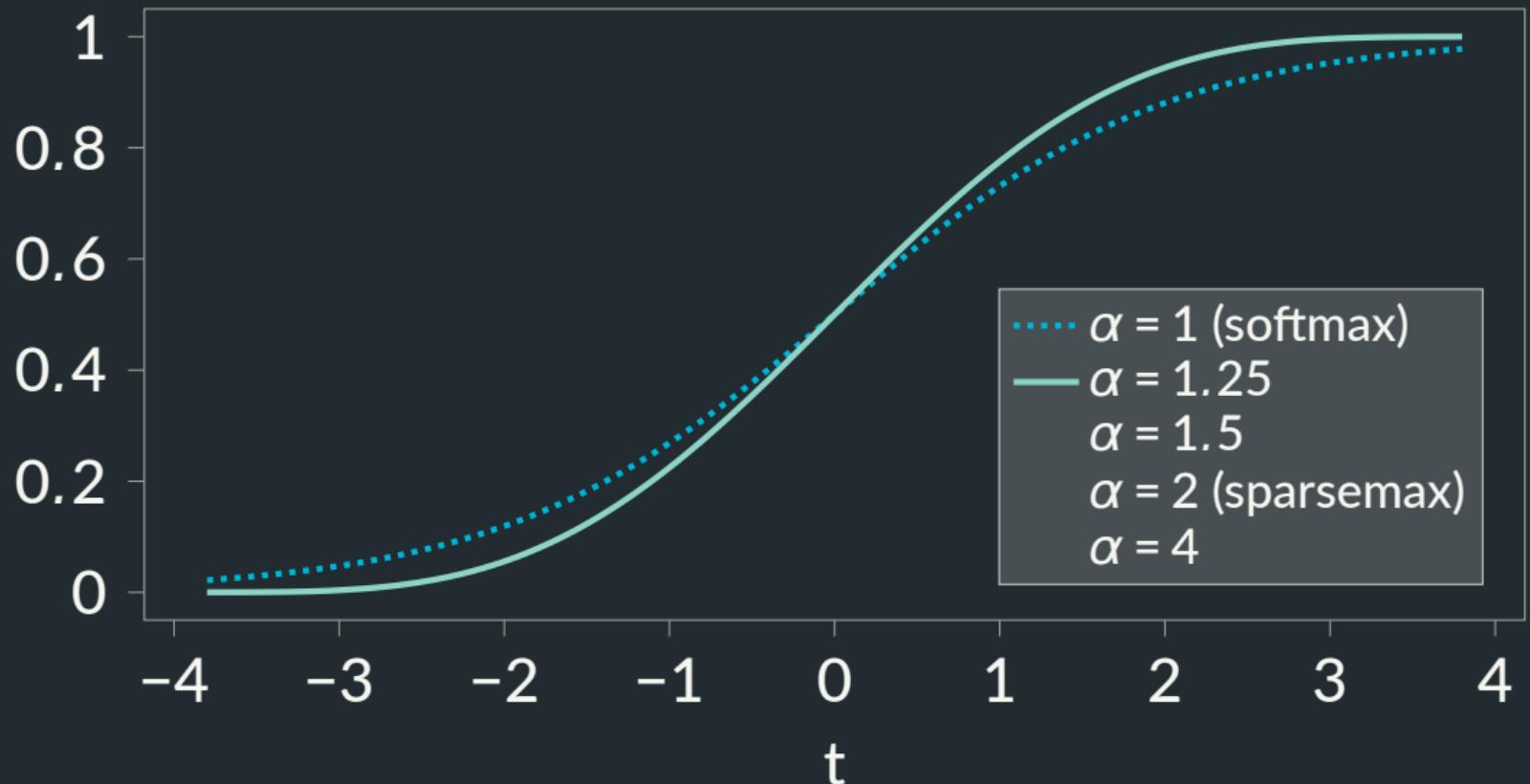
α -Entmax

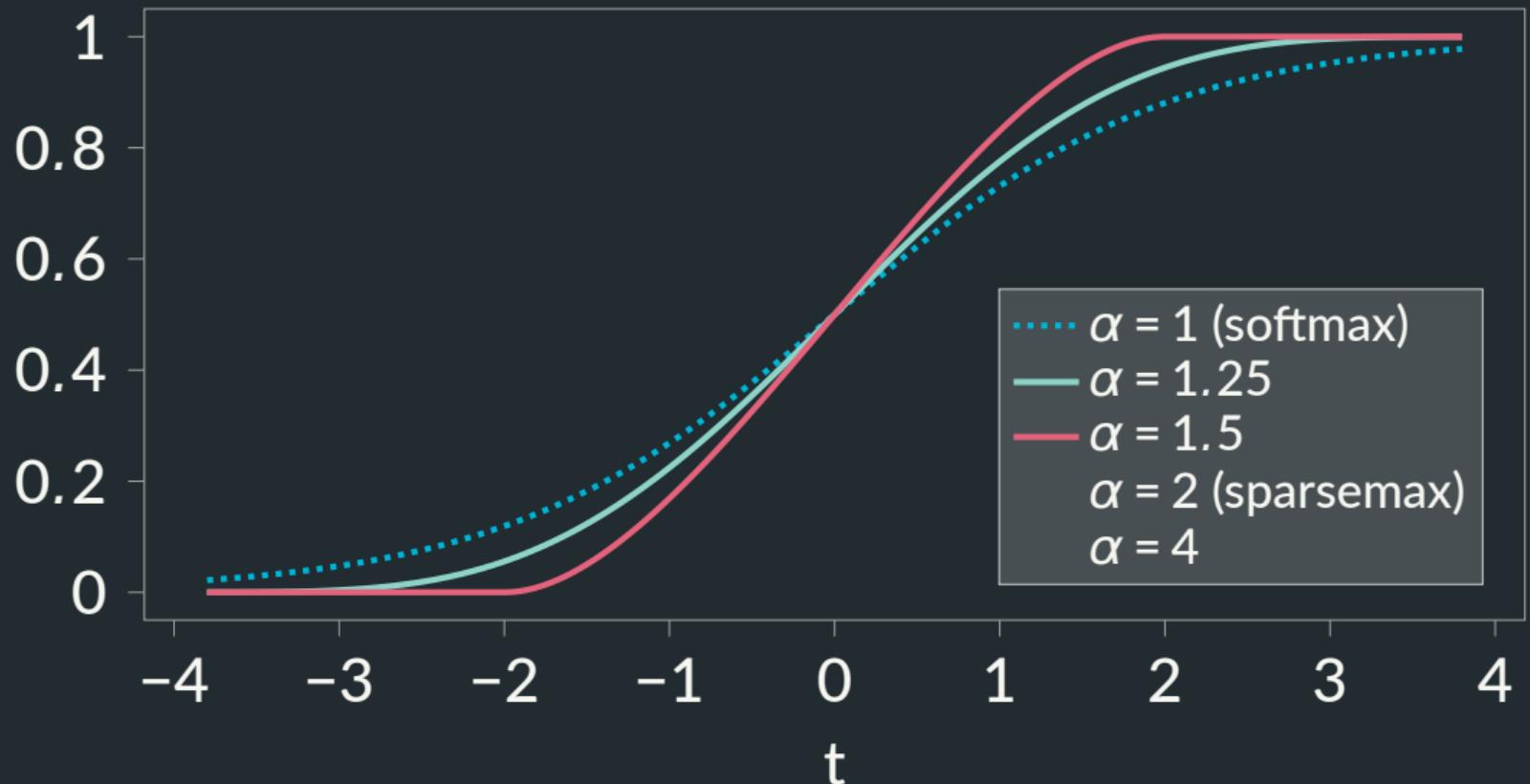
Parametrized by $\alpha \geq 0$:

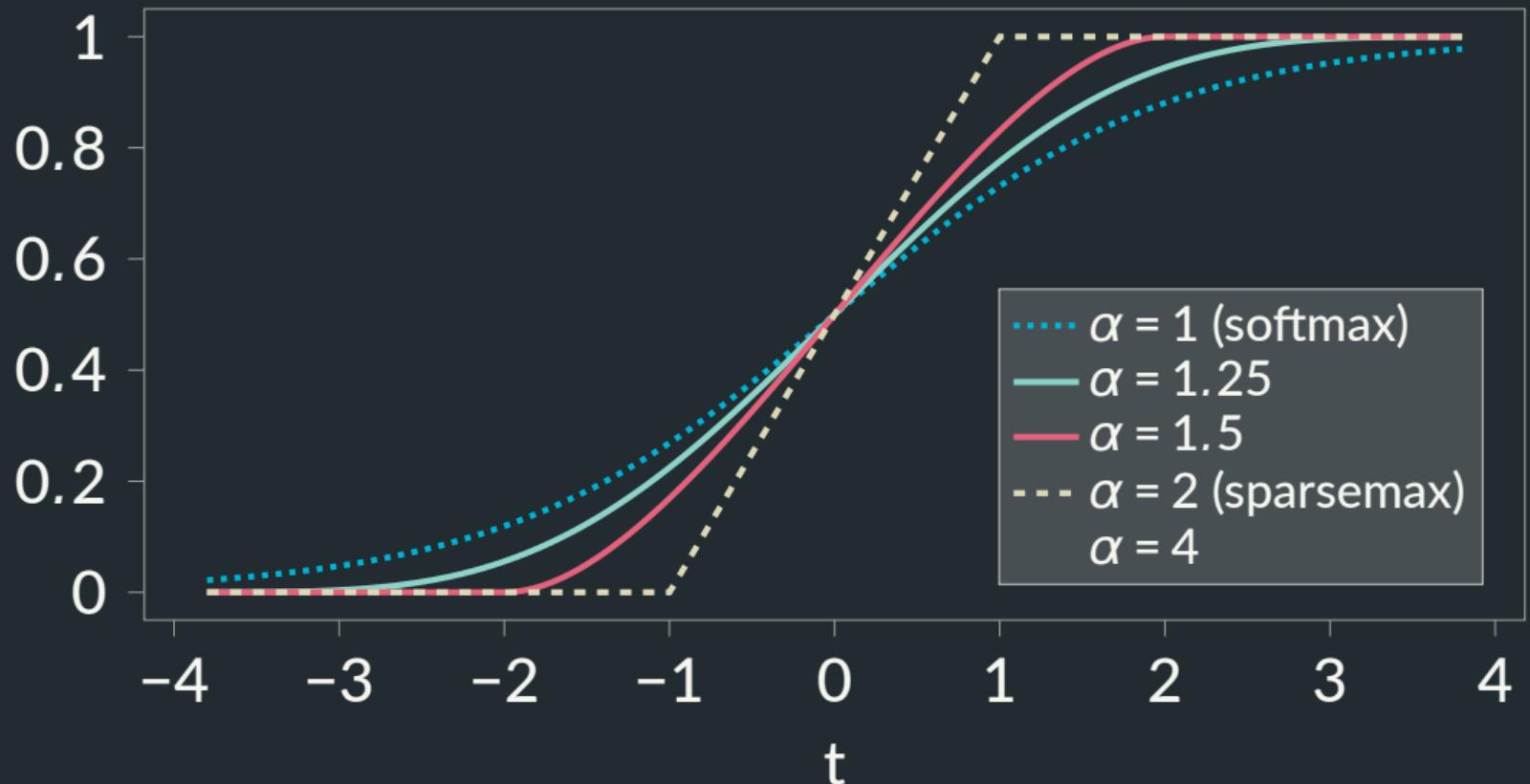
- **Argmax** corresponds to $\alpha \rightarrow \infty$
- **Softmax** amounts to $\alpha \rightarrow 1$
- **Sparsemax** amounts to $\alpha = 2$.

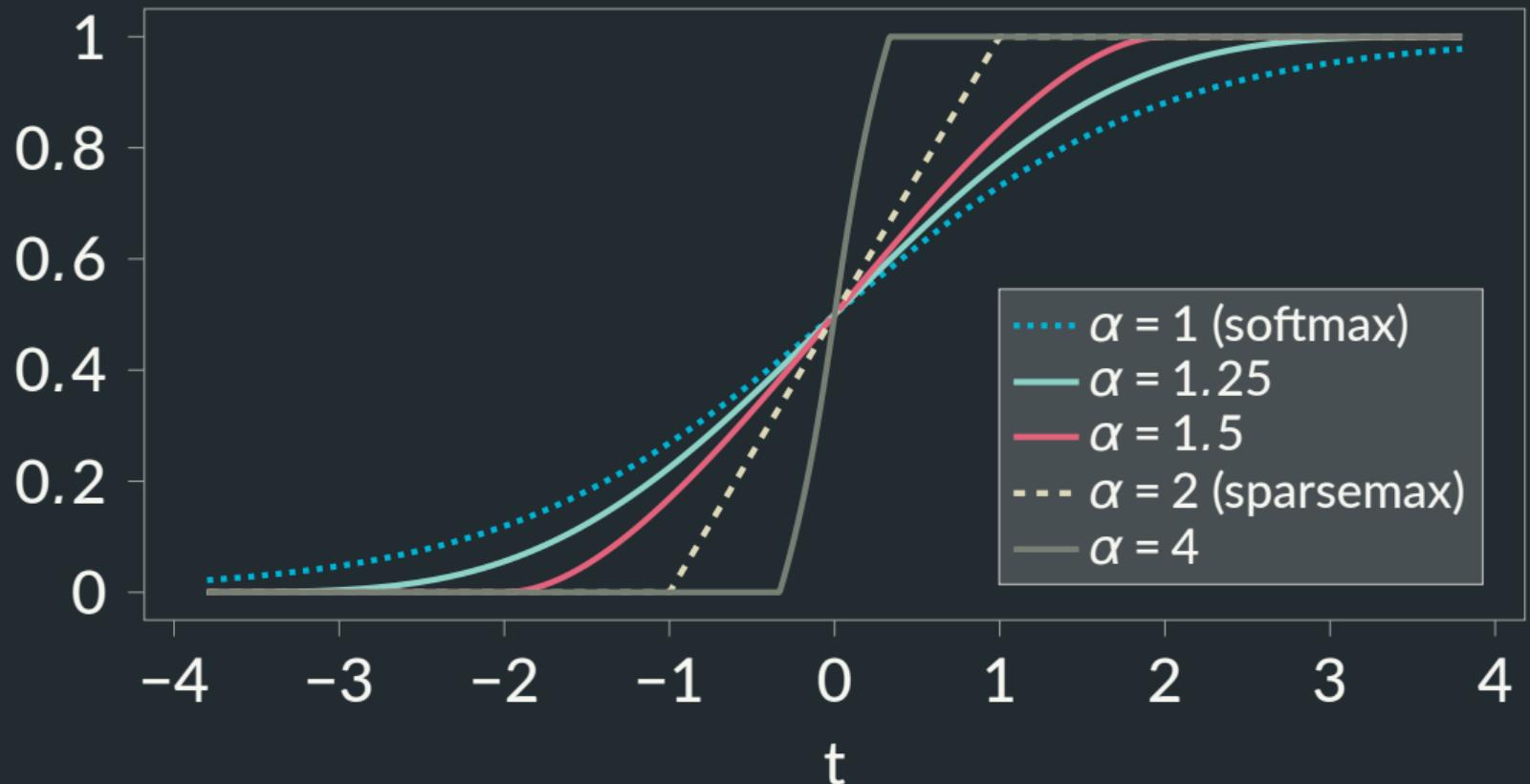
Key result: can be sparse for $\alpha > 1$, propensity for sparsity increases with α .











Learning α

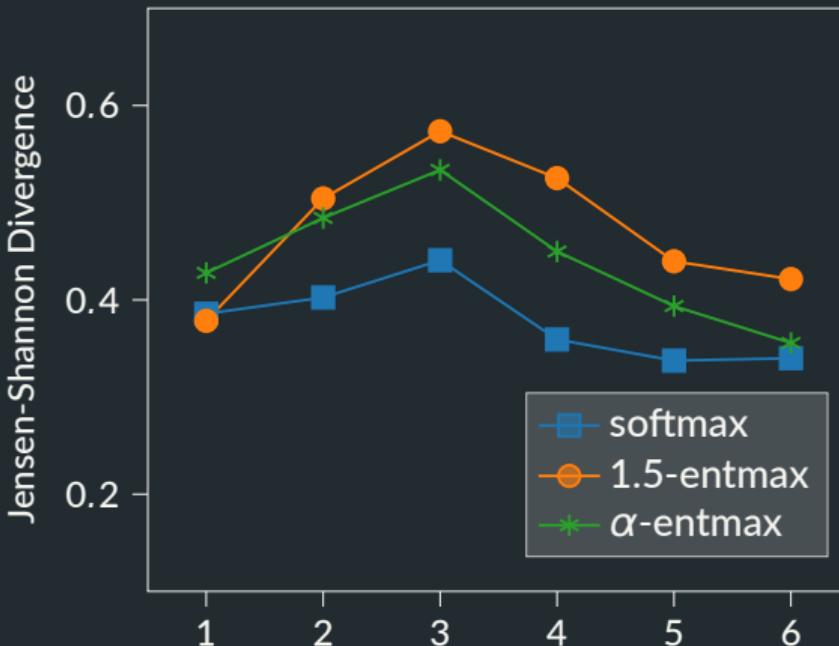
Learning α

Key contribution:

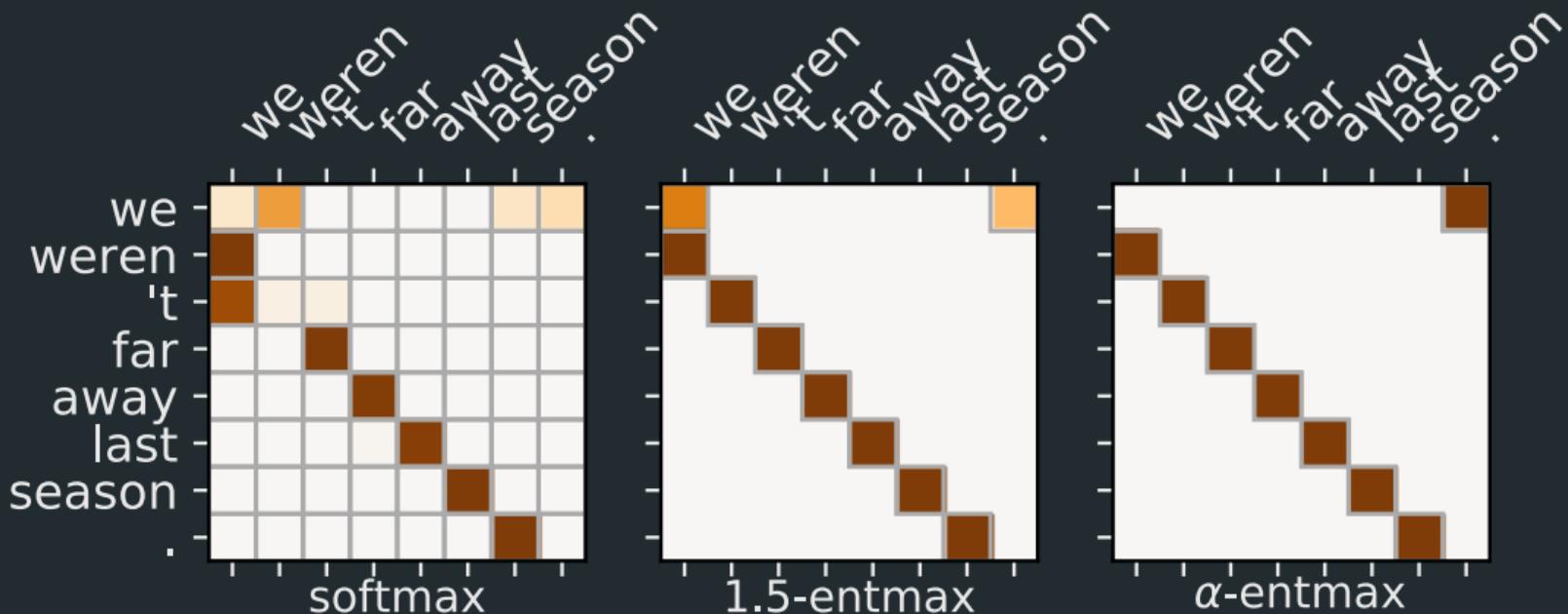
a closed-form expression for $\frac{\partial \alpha\text{-entmax}(\mathbf{z})}{\partial \alpha}$



Head Diversity per Layer

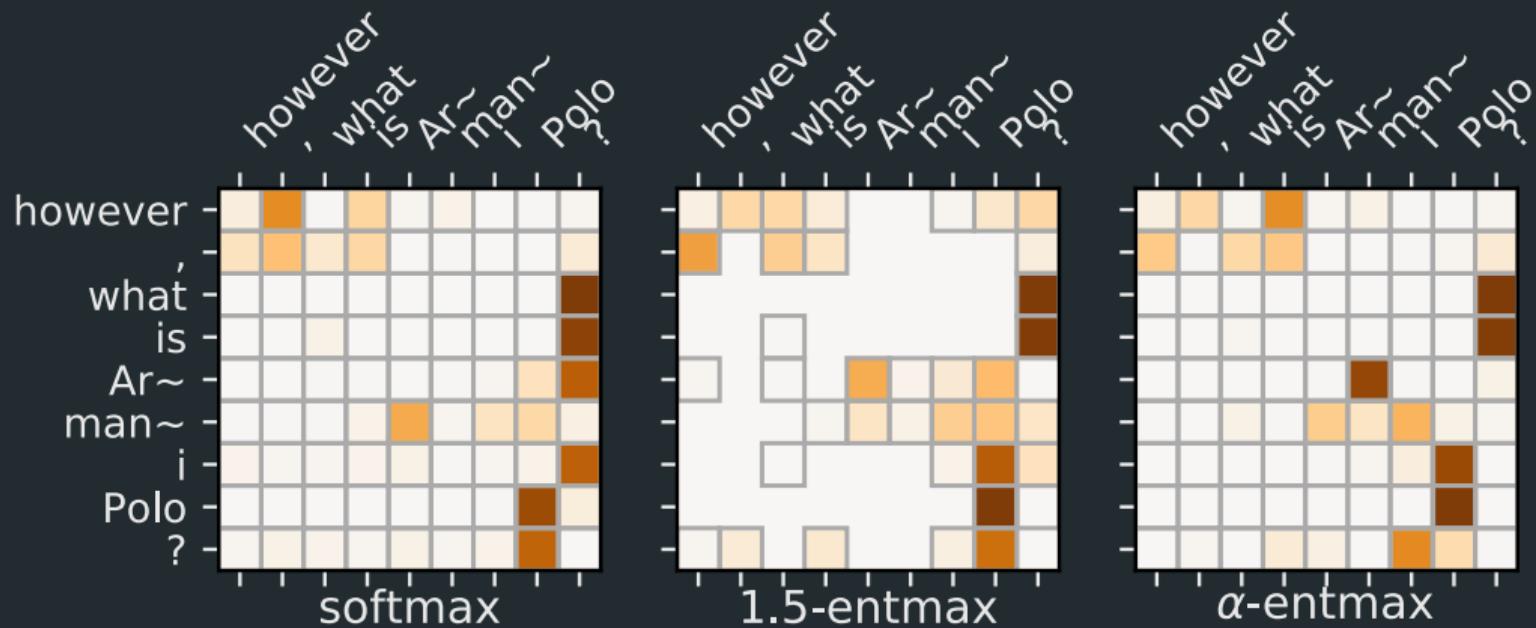


Previous Position Head



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 APE with Transfer Learning

Adaptively Sparse Transformers

Efficient Marg. of Discrete Latent Variables via Sparsity

Conclusions

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

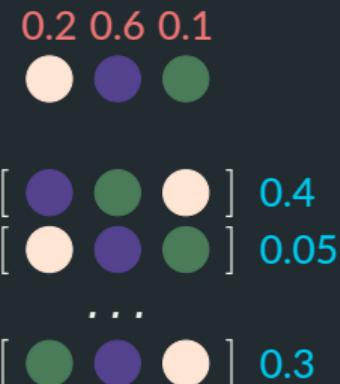


Training Discrete or Structured Latent Variable Models

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, (...)



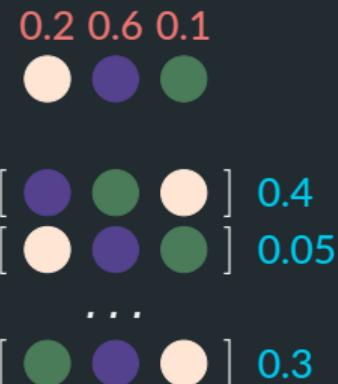
Training Discrete or Structured Latent Variable Models

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

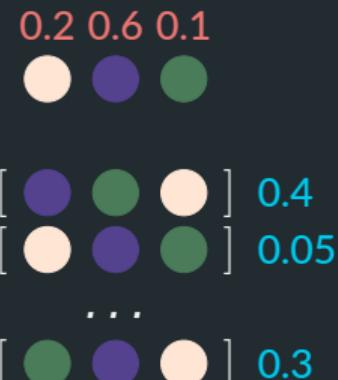
Latent variable z can be **discrete** or **structured**

$\pi(z|x, \theta)$: 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 π 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!**

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

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 APE with Transfer Learning

Adaptively Sparse Transformers

Efficient Marg. of Discrete Latent Variables via Sparsity

Conclusions

References I

-  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](https://doi.org/10.1186/1687-5281-2013-55).
-  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](https://doi.org/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*.
-  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*.

References II

-  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*.