

# **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

- Subs
- Pow
- Rem

# 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?

CDT-2

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 integrate
- 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

Forbes

≡ WIRED

SUBSCRIBE

## 'Dangerous' AI offers to write fake news

By Jane Wakefield  
Technology reporter

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

mputing power required for AI

rd  
ess  
N

utperform  
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: Let's mainly use attention mechanisms!

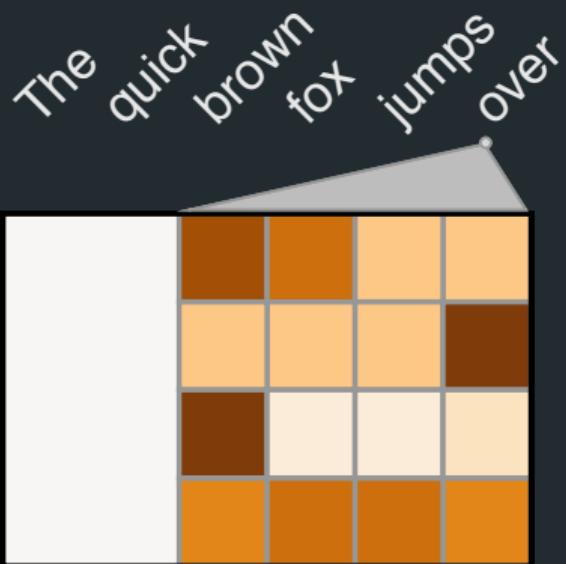


# A bit of context on transformers

What if... Attention is all you need?

Key idea: Let's mainly use attention mechanisms!

- Do attention with multiple heads (i.e. attention mechanisms in parallel)

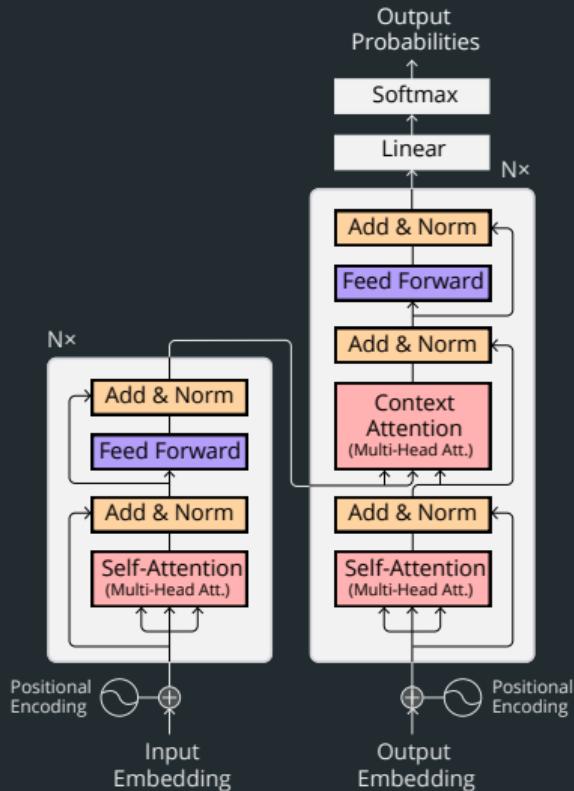


# A bit of context on transformers

What if... Attention is all you need?

Key idea: Let's mainly use attention mechanisms!

- Do attention with multiple heads (i.e. attention mechanisms in parallel)
- ... and do it through several layers

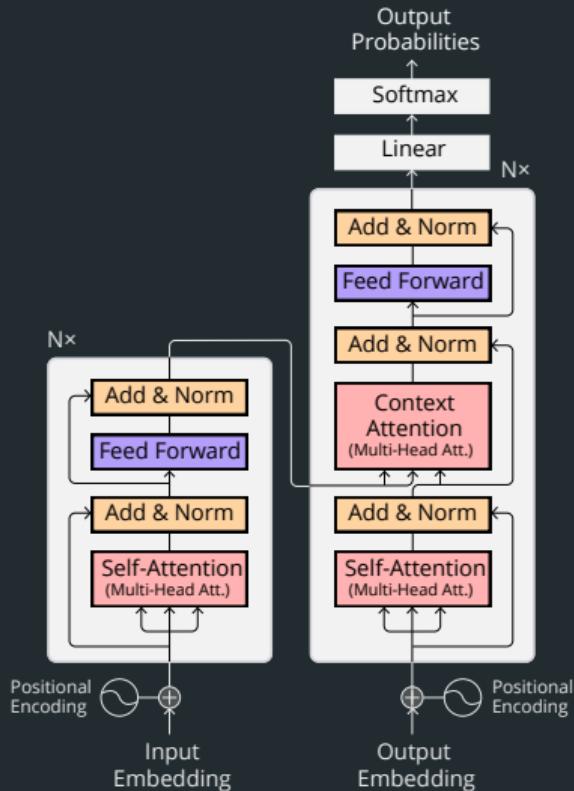


# A bit of context on transformers

What if... Attention is all you need?

Key idea: Let's mainly use attention mechanisms!

- Do attention with multiple heads (i.e. attention mechanisms in parallel)
- ... and do it through several layers
- Inspiration for big general-purpose models like BERT and GPT-3!

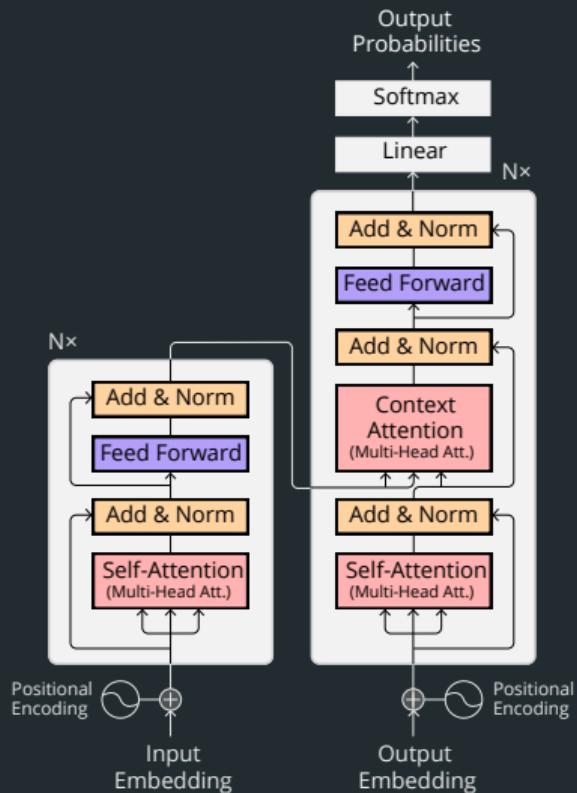


# A bit of context on transformers

What if... Attention is all you need?

Key idea: Let's mainly use attention mechanisms!

- Do attention with multiple heads (i.e. attention mechanisms in parallel)
- ... and do it through several layers
- Inspiration for big general-purpose models like BERT and GPT-3!

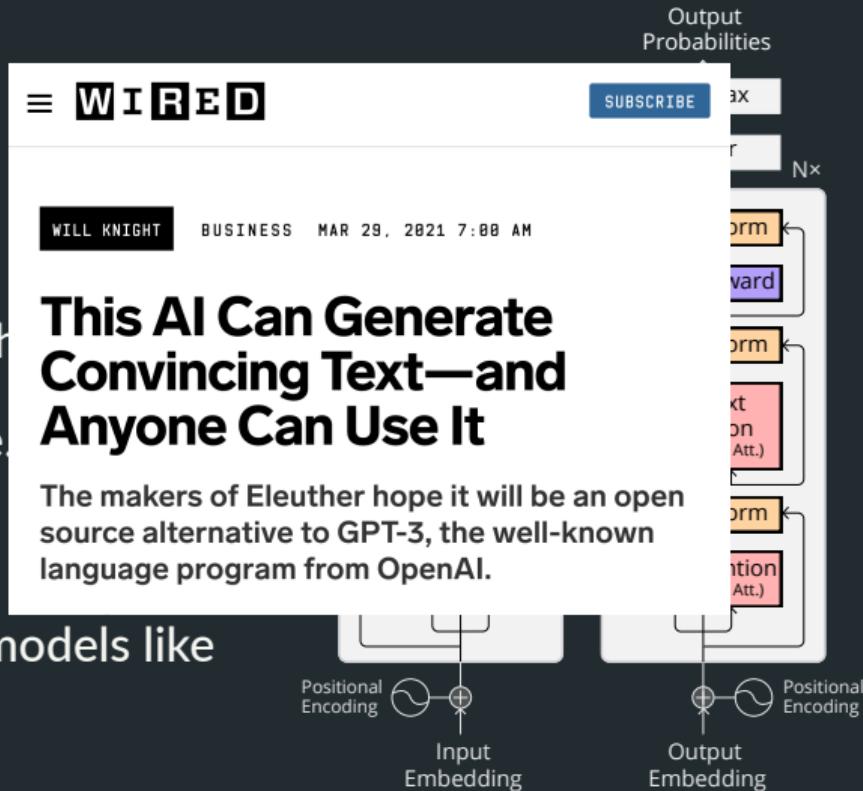


# A bit of context on transformers

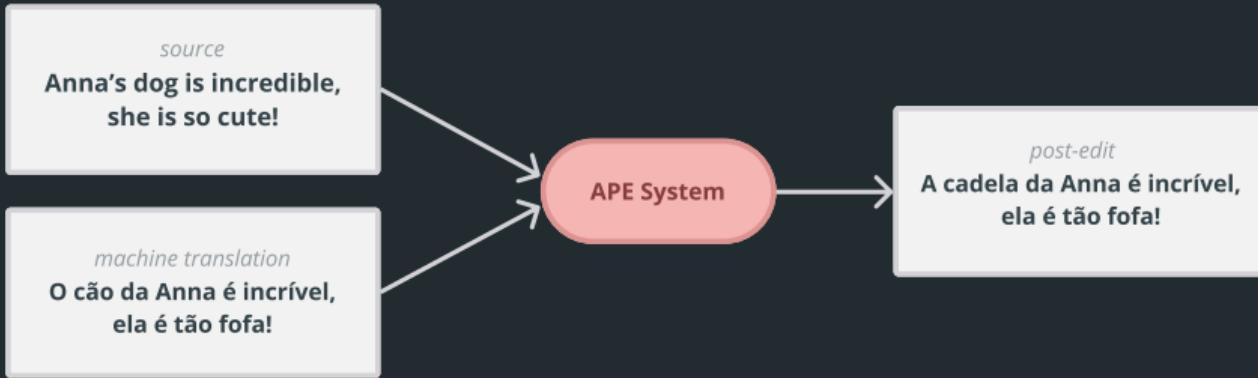
What if... Attention is all you need?

Key idea: Let's mainly use attention mechanisms

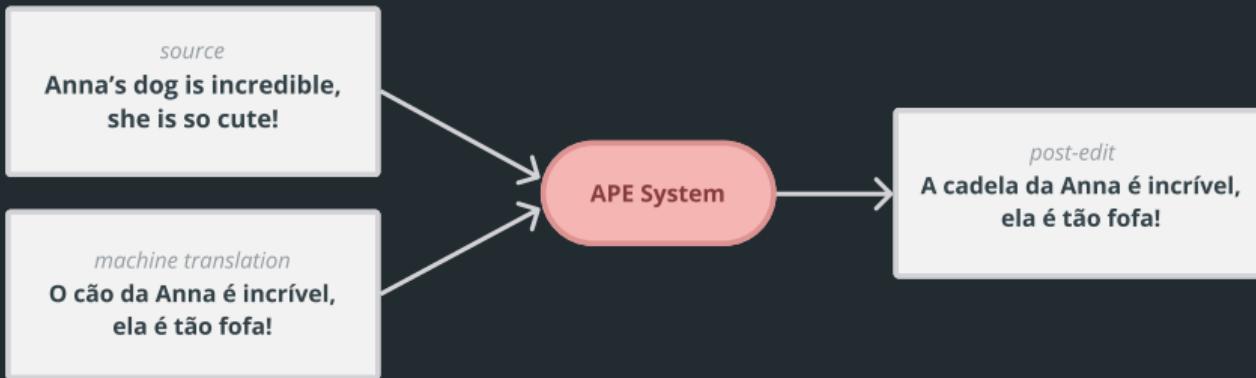
- Do attention with multiple heads (i.e. mechanisms in parallel)
- ... and do it through several layers
- Inspiration for big general-purpose models like BERT and GPT-3!



# What is APE?

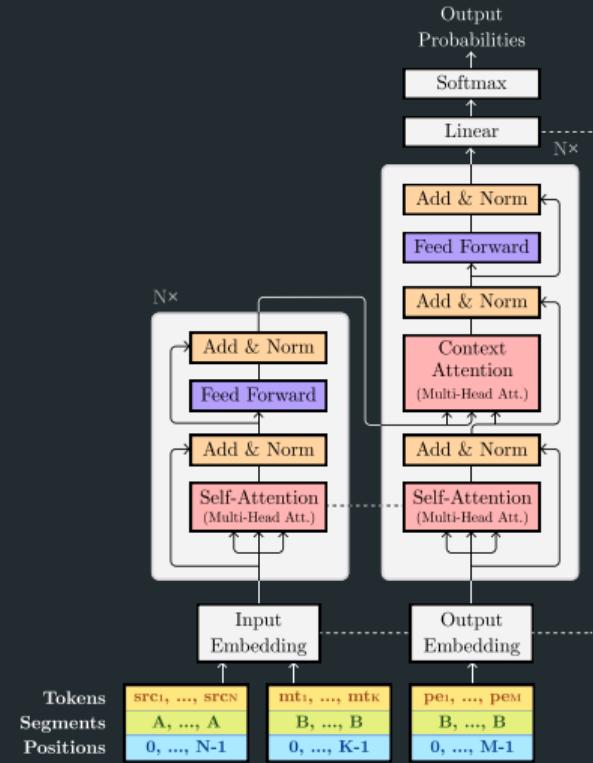


# What is APE?



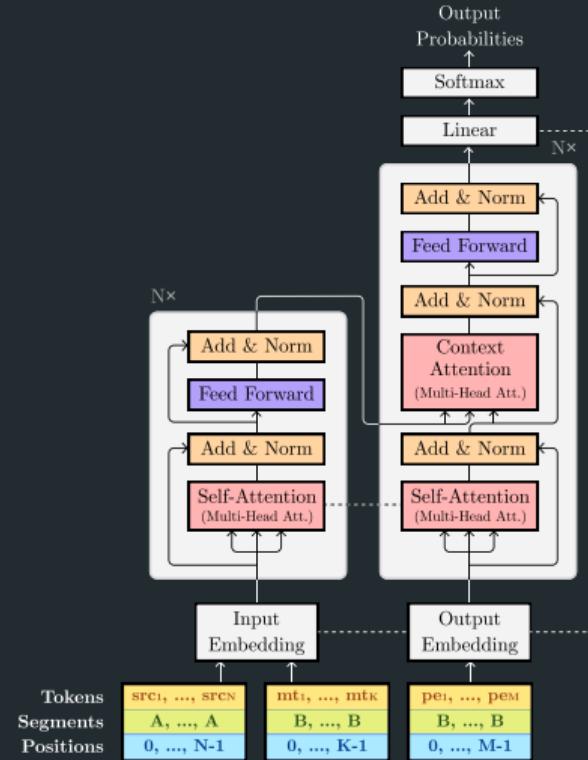
**Challenge:** APE data is very scarce! Need to create artificial data.

# BERT for APE



# BERT for APE

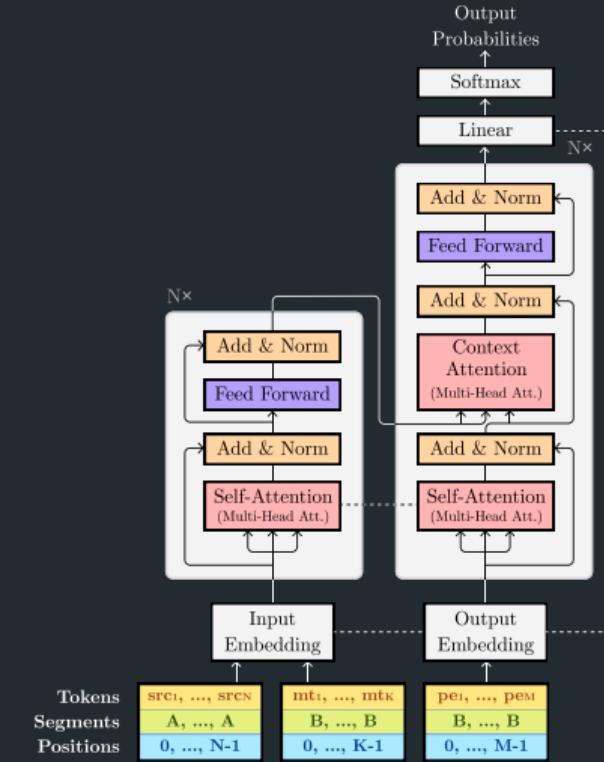
Key idea: Use BERT to do APE



# BERT for APE

Key idea: Use BERT to do APE

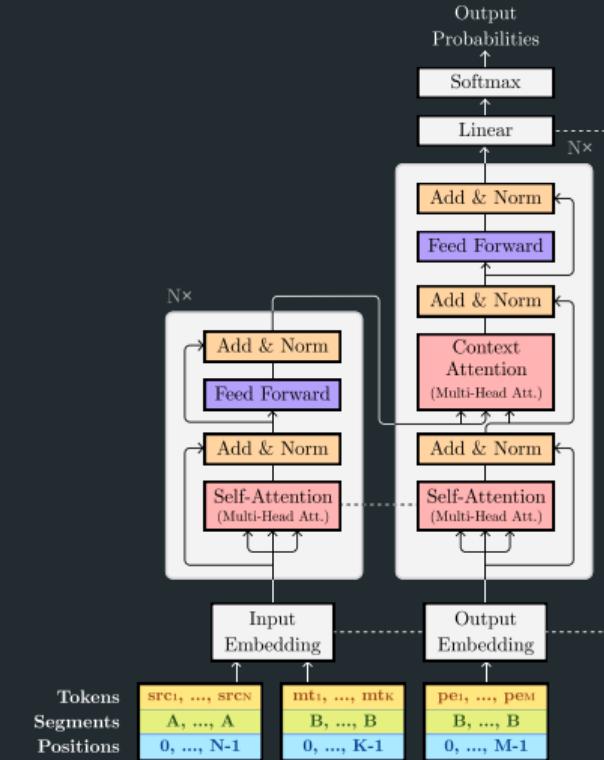
- Prior to this work, BERT was mainly used for simple classification tasks



# BERT for APE

Key idea: Use BERT to do APE

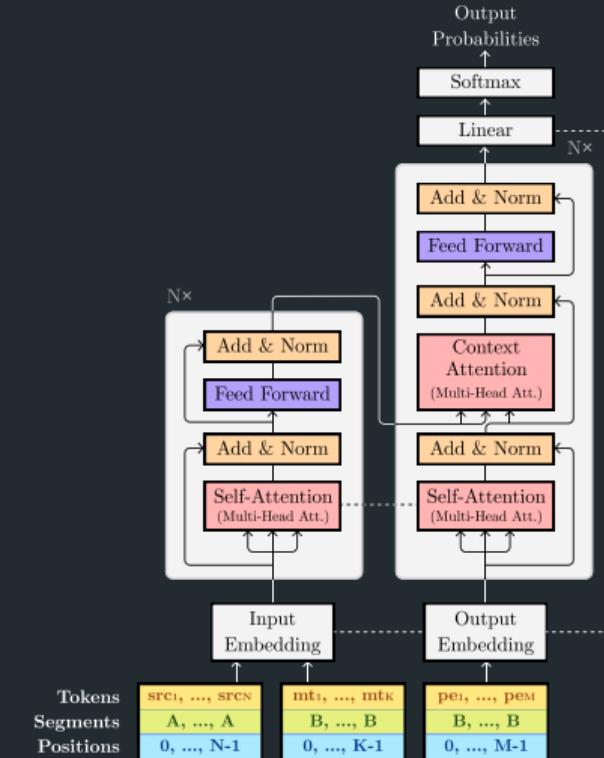
- Prior to this work, BERT was mainly used for simple classification tasks
- We introduced an effective method to use BERT in a generation task (APE)



# 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



# Key results

model (data size)	TER↓	BLEU↑
mt baseline	24.48	62.49

# Key results

---

model (data size)	TER↓	BLEU↑
mt baseline	24.48	62.49
dual-source transformer (8M)	18.10	71.72

---

# Key results

---

model (data size)	TER↓	BLEU↑
mt baseline	24.48	62.49
dual-source transformer (8M)	18.10	71.72
dual-source transformer (23K)	27.73	59.78

---

# Key results

model (data size)	TER↓	BLEU↑
mt baseline	24.48	62.49
dual-source transformer (8M)	18.10	71.72
dual-source transformer (23K)	27.73	59.78
ours (23K)	19.03	70.66

# Key results

model (data size)	TER↓	BLEU↑
mt baseline	24.48	62.49
dual-source transformer (8M)	18.10	71.72
dual-source transformer (23K)	27.73	59.78
ours (23K)	19.03	70.66
ours (8M)	17.26	73.42

# Conclusions and impact

# Conclusions and impact

- One of pioneers in using pre-trained transformer encoders for a generation task

# Conclusions and impact

- One of pioneers in using pre-trained transformer encoders for a generation task
- Massive improvement in low-resource scenario (**data-efficiency**)

# Conclusions and impact

- One of pioneers in using pre-trained transformer encoders for a generation task
- Massive improvement in low-resource scenario (**data-efficiency**)
- Steered SOTA of APE towards using **weak supervision** through pre-trained models

# Conclusions and impact

- One of pioneers in using pre-trained transformer encoders for a generation task
- Massive improvement in low-resource scenario (**data-efficiency**)
- Steered SOTA of APE towards using **weak supervision** through pre-trained models
- Inspired other works that use scarce data (e.g., dialogue with metadata)

# 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**.

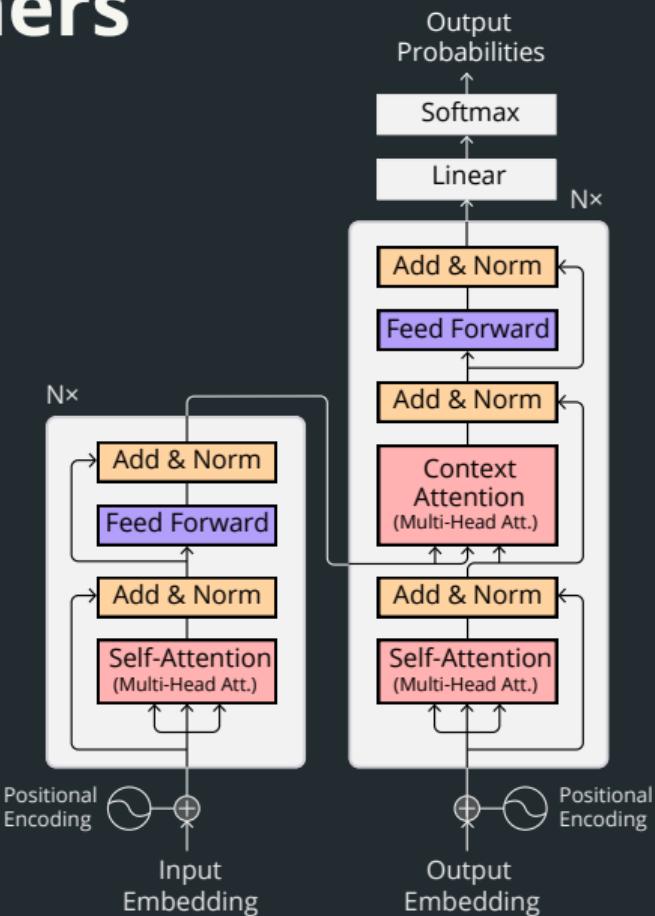
Our solution is to bet on **sparsity**:

- for interpretability
- for discovering linguistic structure
- for efficiency

# Transformers

In each attention head:

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



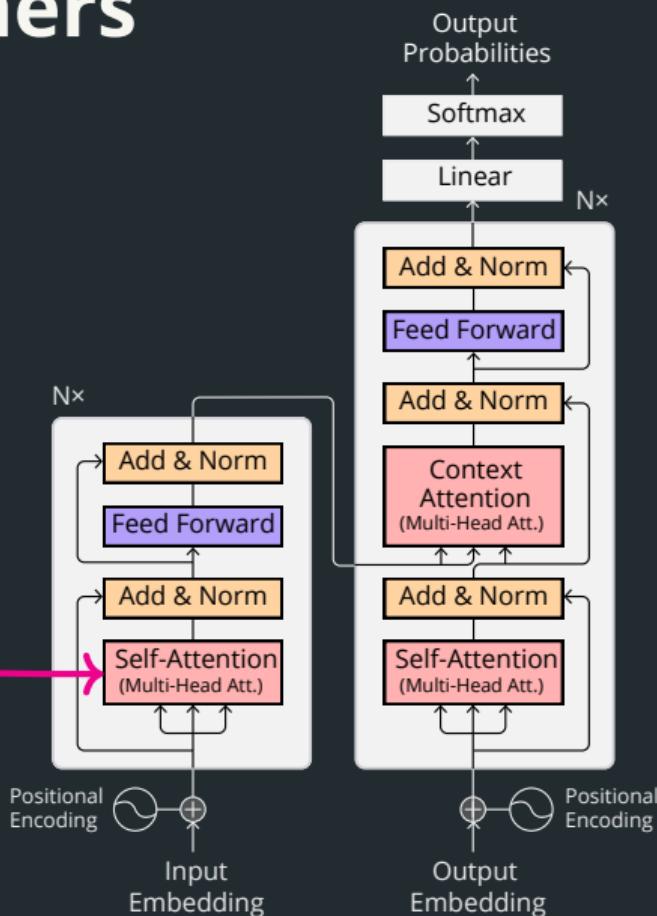
# Transformers

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



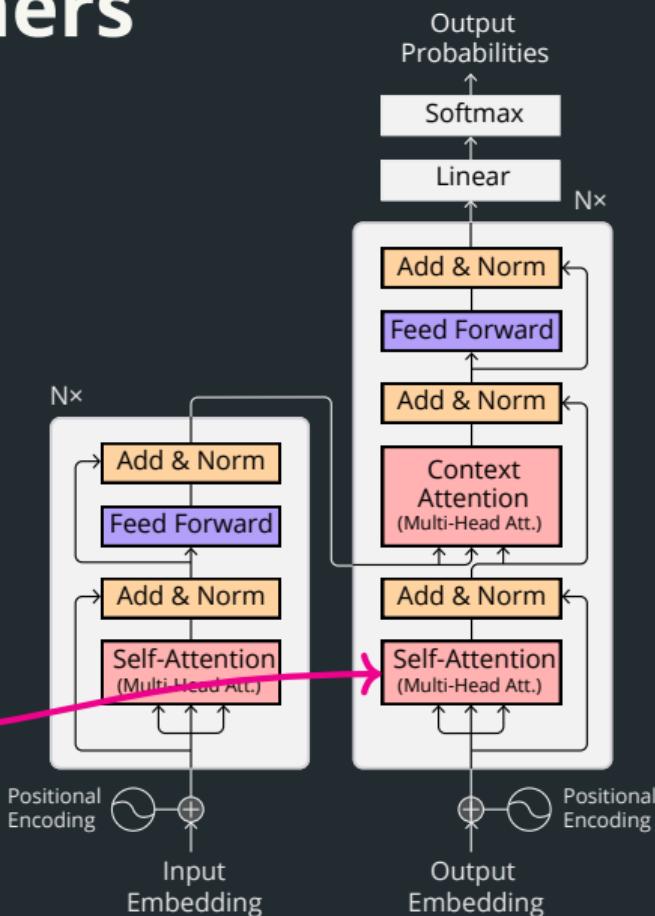
# Transformers

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



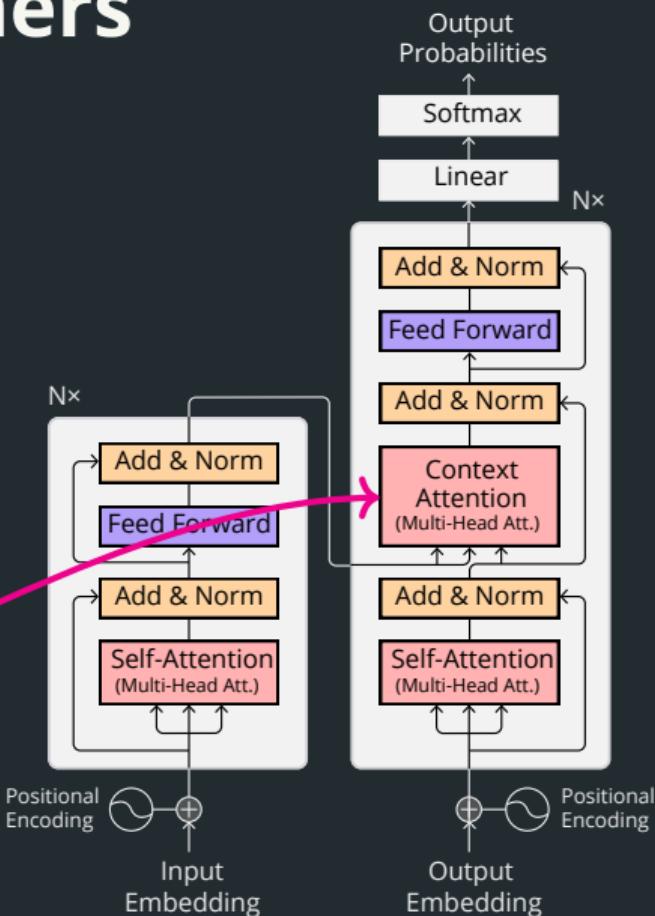
# Transformers

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  $\alpha$ ! 

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

# $\alpha$ -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  $\alpha \geq 0$ :

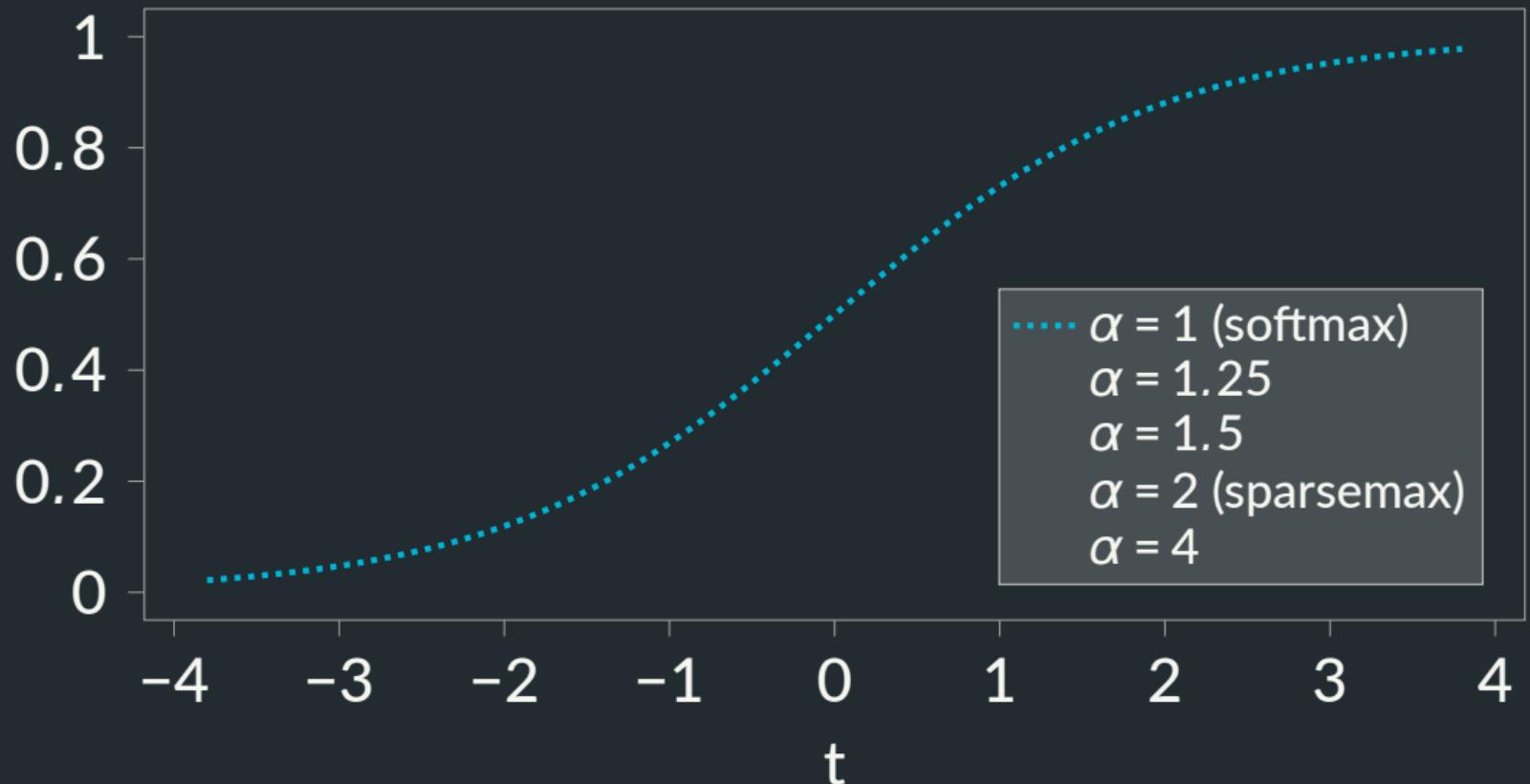
- **Argmax** corresponds to  $\alpha \rightarrow \infty$
- **Softmax** amounts to  $\alpha \rightarrow 1$
- **Sparsemax** amounts to  $\alpha = 2$ .

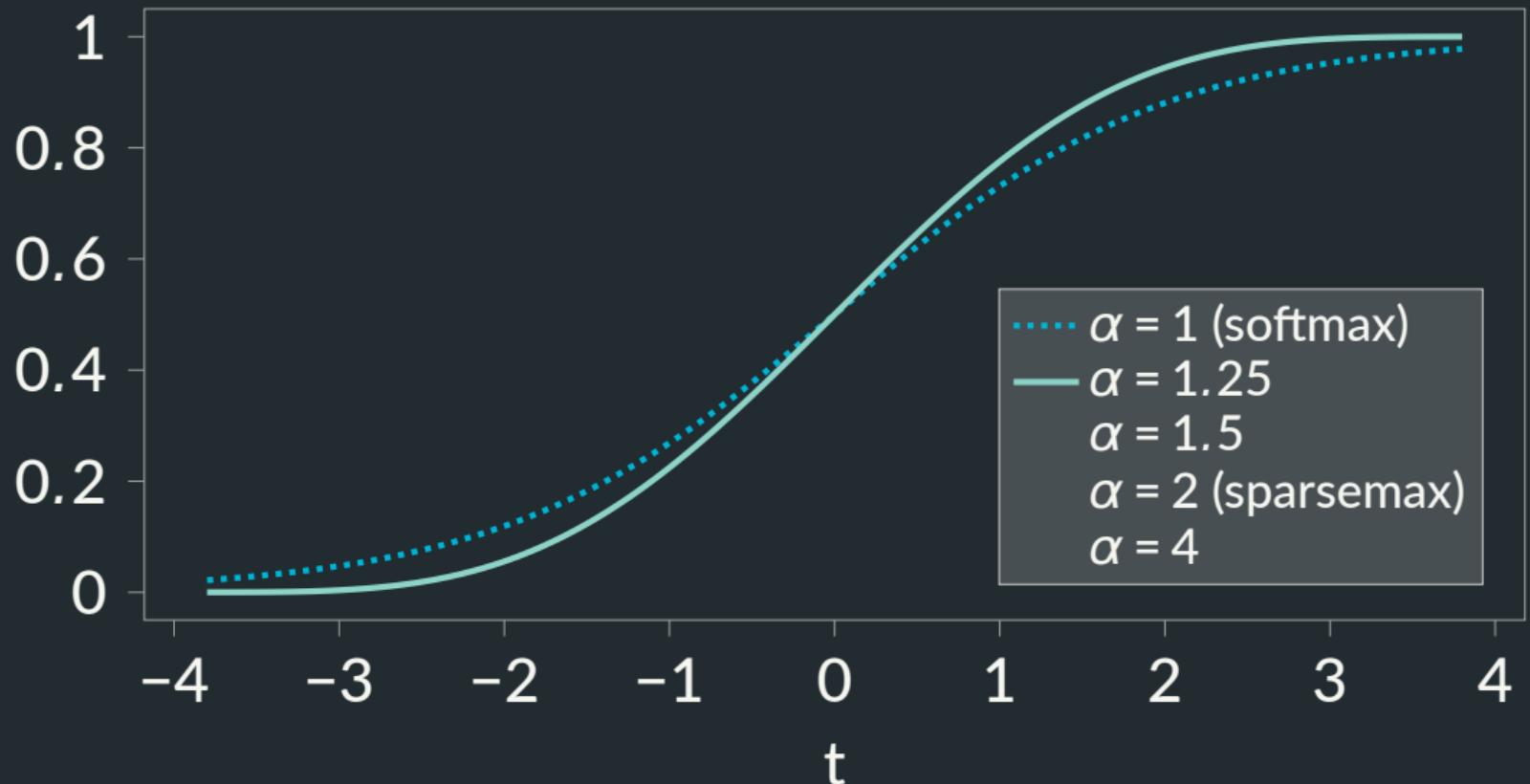
# $\alpha$ -entmax

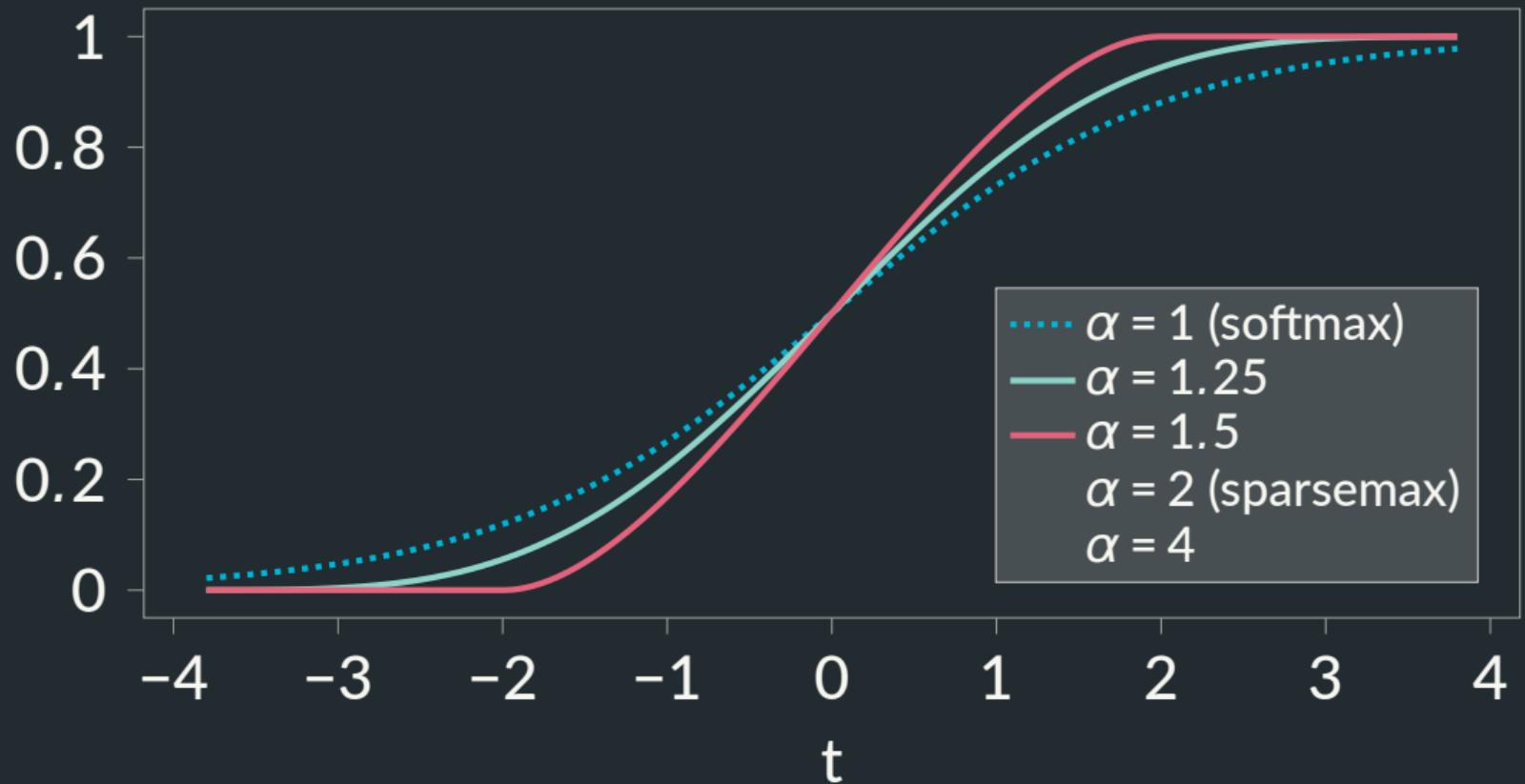
Parametrized by  $\alpha \geq 0$ :

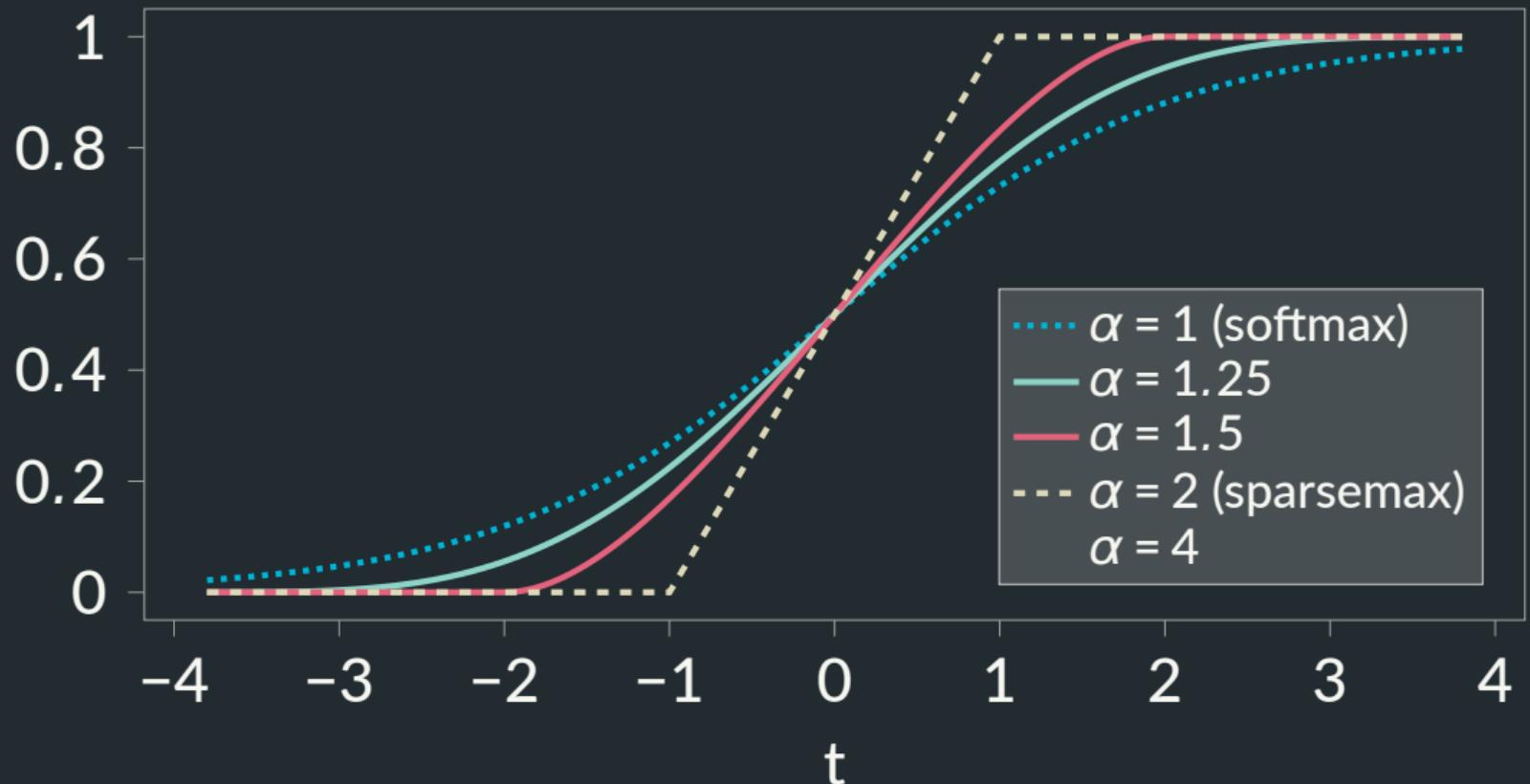
- **Argmax** corresponds to  $\alpha \rightarrow \infty$
- **Softmax** amounts to  $\alpha \rightarrow 1$
- **Sparsemax** amounts to  $\alpha = 2$ .

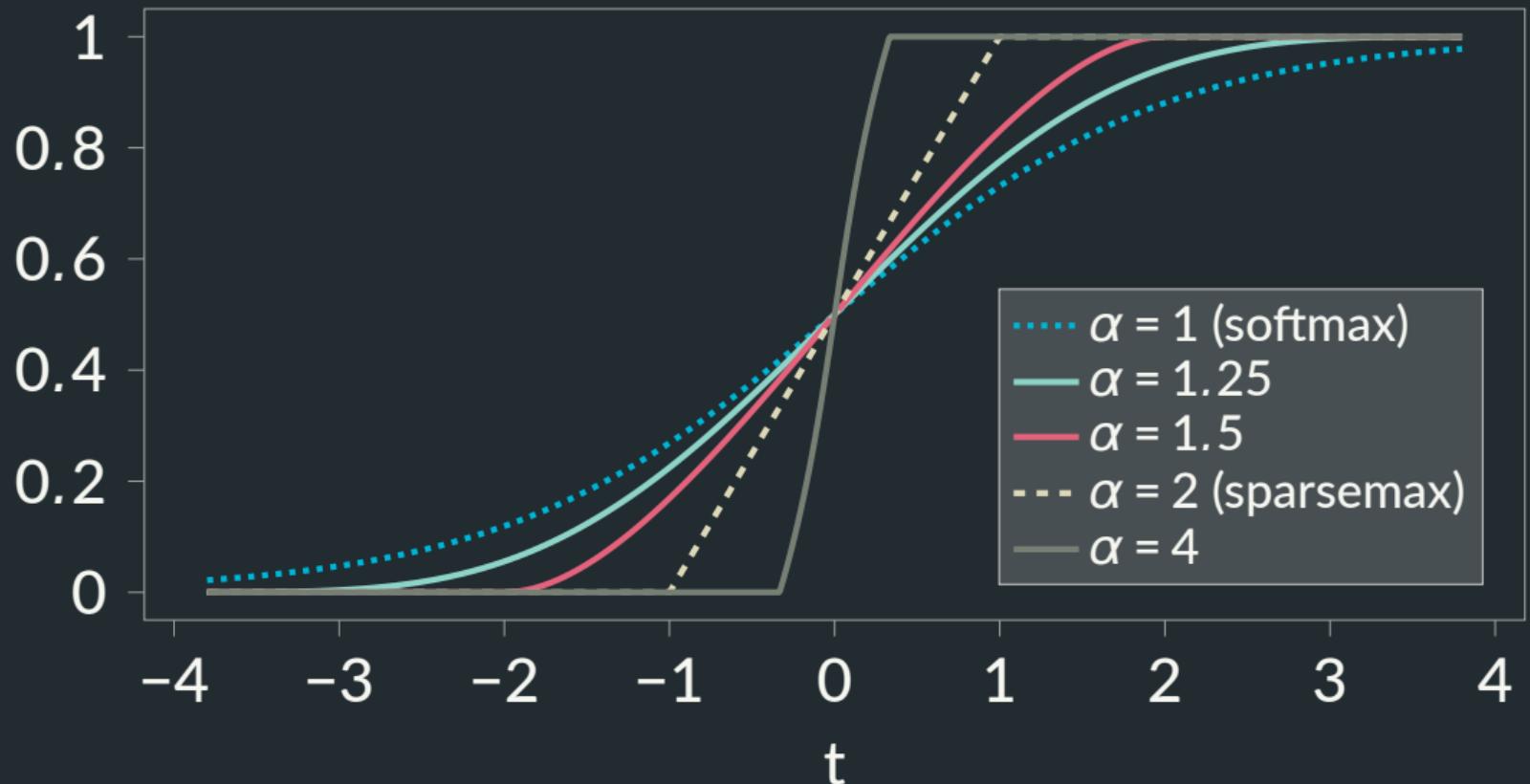
Key result: can be sparse for  $\alpha > 1$ , propensity for sparsity increases with  $\alpha$ .











# Learning $\alpha$

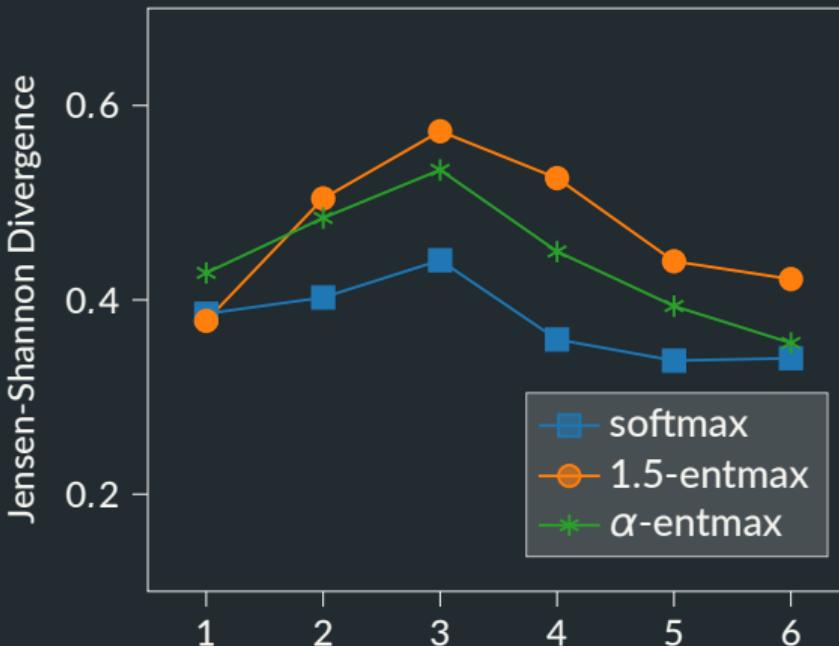
# Learning $\alpha$

Key contribution:

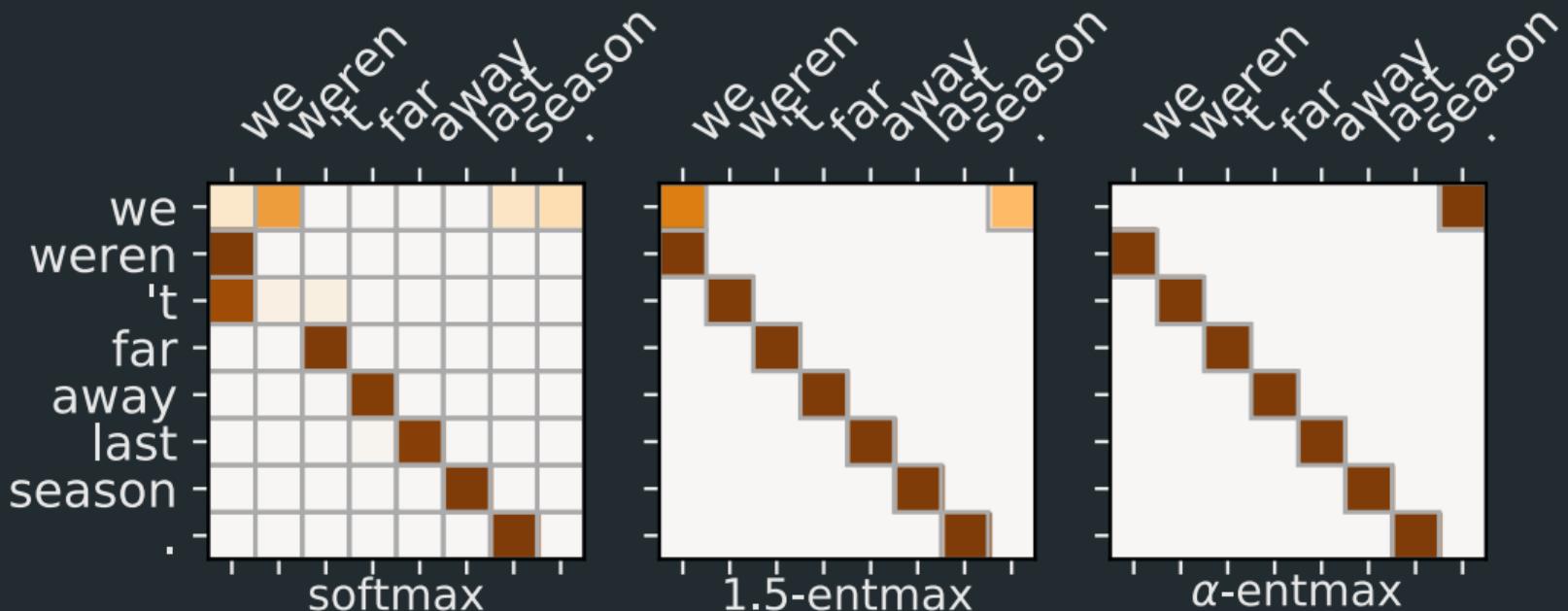
a closed-form expression for  $\frac{\partial \alpha\text{-entmax}(\mathbf{z})}{\partial \alpha}$



# Head diversity per layer

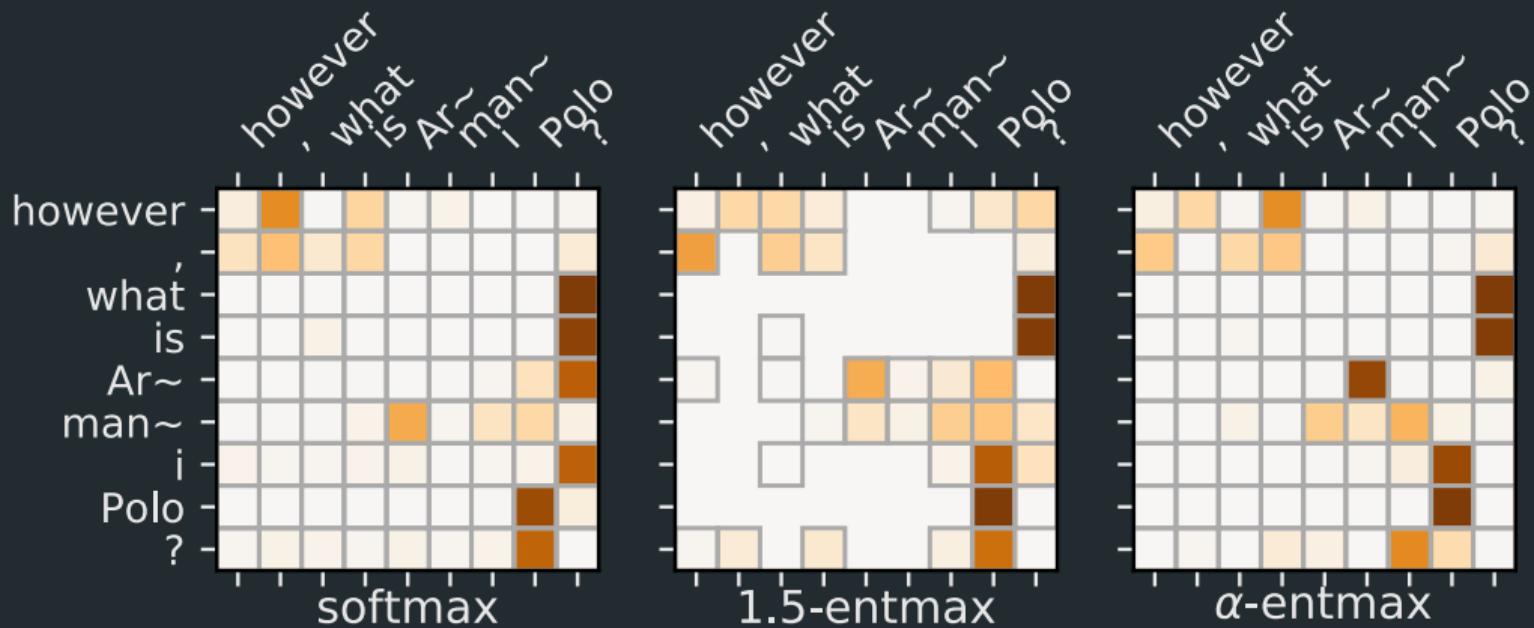


# Previous position head



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

# Interrogation-detecting head



Learned  $\alpha = 1.05$ .

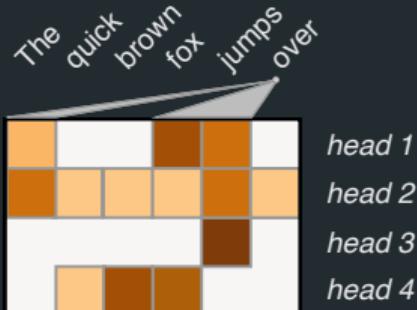
# Conclusions and impact

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

# Conclusions and impact

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

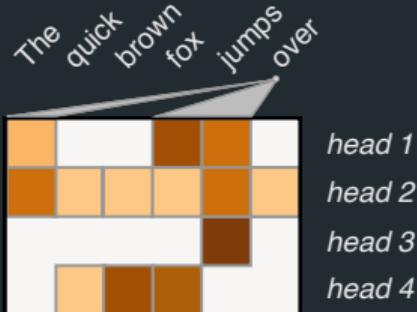
*adaptive sparsity*



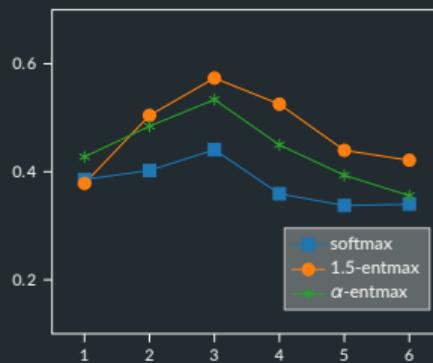
# Conclusions and impact

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

*adaptive sparsity*



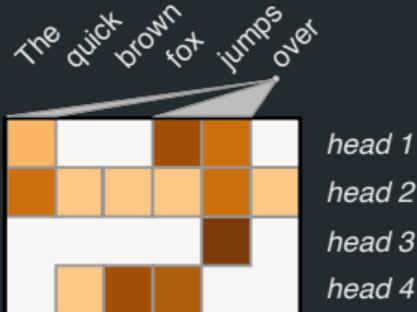
*reduced head redundancy*



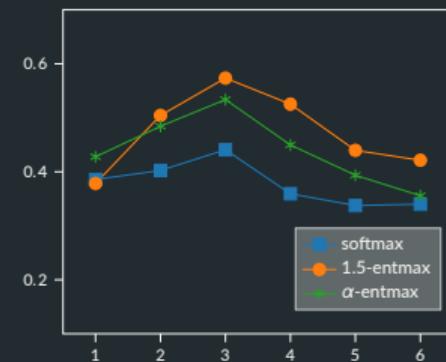
# Conclusions and impact

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

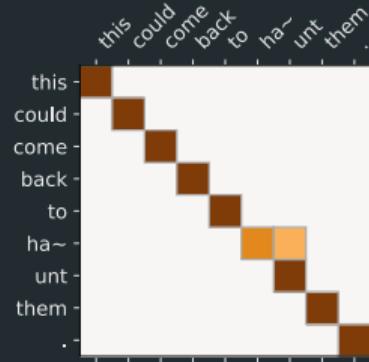
*adaptive sparsity*



*reduced head redundancy*



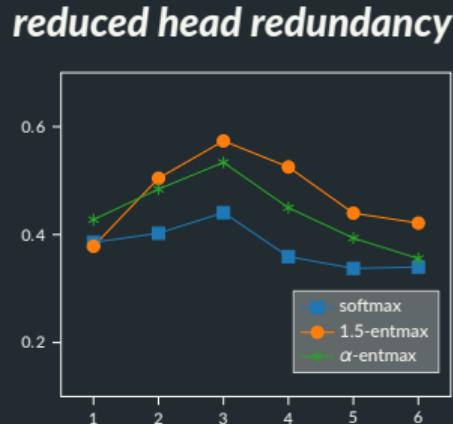
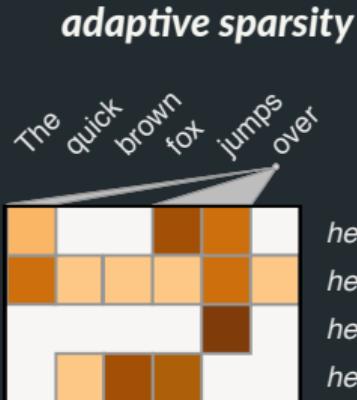
*clearer head roles*



# Conclusions and impact

Subsequent work has:

- focused on taking computational advantage of sparsity
- proposed other sparse activations (e.g., ReLU)
- incorporated fixed attention patterns that we found.



# 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  $z$

0.2 0.6 0.1  
● ● ●

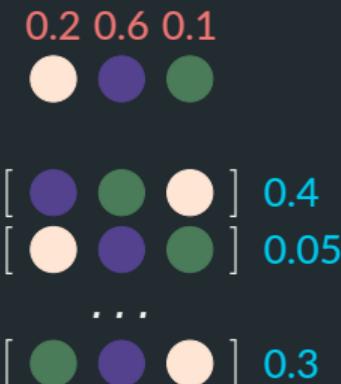
[ ● ● ● ]  
[ ● ● ● ]

...  
[ ● ● ● ]

# Training Discrete or Structured Latent Variable Models

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

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

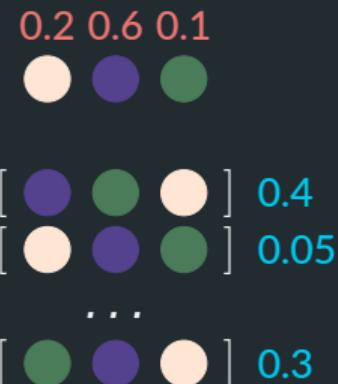


# Training Discrete or Structured Latent Variable Models

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, (...)



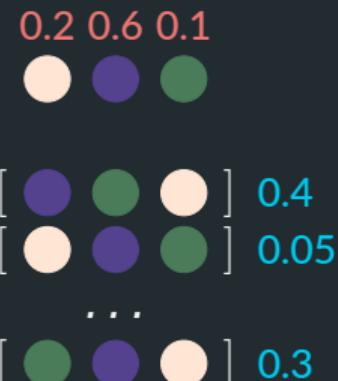
# Training Discrete or Structured Latent Variable Models

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

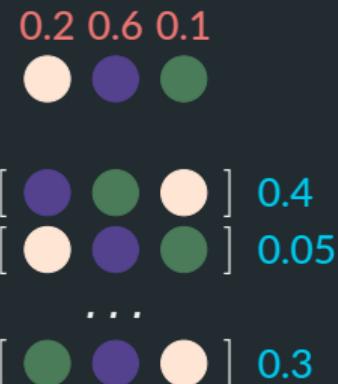
Latent variable  $z$  can be **discrete** or **structured**

$\pi(z|x, \theta)$ : 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

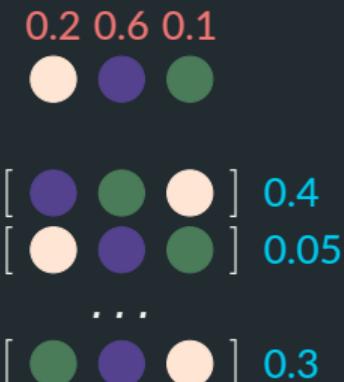
Latent variable  $z$  can be **discrete** or **structured**

$\pi(z|x, \theta)$ : 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!

# Key Takeaways

We introduce a new method  
to train latent variable models.

# Key Takeaways

We introduce a new method  
to train latent variable models.

*discrete and structured*

0.2 0.6 0.1  


[  ] 0.4  
[  ] 0.05

...

[  ] 0.3

# Key Takeaways

We introduce a new method  
to train latent variable models.

*discrete and structured*

0.2 0.6 0.1  


[    ] 0.4  
[    ] 0.05  
...  
[    ] 0.3

*deterministic, yet efficient*

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

# Key Takeaways

We introduce a new method  
to train latent variable models.

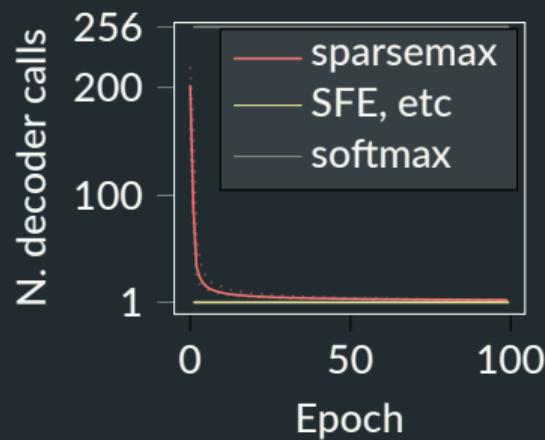
*discrete and structured*

0.2	0.6	0.1			
[				]	0.4
[				]	0.05
...					
[				]	0.3

*deterministic, yet efficient*

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

*sparse, as needed*



# 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 (2020a). "Efficient Marginalization of Discrete and Structured Latent Variables via Sparsity". In: *Proc. NeurIPS*. URL: <https://arxiv.org/abs/2007.01919>.
-  – (2020b). "Efficient marginalization of discrete and structured latent variables via sparsity". In: *Proc. NeurIPS*.
-  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*.
-  Junczys-Dowmunt, Marcin and Roman Grundkiewicz (2018). "MS-UEdin Submission to the WMT2018 APE Shared Task: Dual-Source Transformer for Automatic Post-Editing". In: *Proceedings of WMT18*.

# References II

-  Kodama, Takashi, Ryuichiro Higashinaka, Koh Mitsuda, Ryo Masumura, Yushi Aono, Ryuta Nakamura, Noritake Adachi, and Hidetoshi Kawabata (2020). "Generating Responses That Reflect Meta Information in User-Generated Question Answer Pairs". In: *Proceedings of LREC*.
-  Lazaridou, Angeliki, Alexander Peysakhovich, and Marco Baroni (2017). "Multi-agent cooperation and the emergence of (natural) language". In: *Proc. ICLR*.
-  Lee, Jihyung, WonKee Lee, Jaehun Shin, Baikjin Jung, Young-Kil Kim, and Jong-Hyeok Lee (2020). "POSTECH-ETRI's Submission to the WMT2020 APE Shared Task: Automatic Post-Editing with Cross-lingual Language Model". In: *Proceedings of WMT*.
-  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 and Mathieu Blondel (2017). "A Regularized Framework for Sparse and Structured Neural Attention". In: *arXiv preprint arXiv:1705.07704*.
-  Niculae, Vlad, André FT Martins, Mathieu Blondel, and Claire Cardie (2018). "SparseMAP: Differentiable sparse structured inference". In: *Proc. of ICML*.

# References III

-  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*.
-  Raganato, Alessandro, Yves Scherrer, and Jörg Tiedemann (2020). "Fixed Encoder Self-Attention Patterns in Transformer-Based Machine Translation". In: *Proceedings of EMNLP*.
-  Treviso, Marcos, António Góis, Patrick Fernandes, Erick Fonseca, and André F. T. Martins (2022). "Predicting Attention Sparsity in Transformers". In: *Proceedings of SPNLP*.
-  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*.
-  Zhang, Biao, Ivan Titov, and Rico Sennrich (2021). "Sparse Attention with Linear Units". In: *Proceedings of EMNLP*.

# Parameter sharing analysis

	TER↓	BLEU↑
MT Baseline	24.76	62.11
Transformer	27.80	60.76
Transformer decoder	20.33	69.31
Pre-trained BERT <i>with CA ← SA</i>	20.83	69.11
<i>and SA ↔ Encoder SA</i>	<b>18.44</b>	<b>72.25</b>
<i>and CA ↔ SA</i>	18.75	71.83
<i>and FF ↔ Encoder FF</i>	19.04	71.53

# $\Omega$ -Regularized Argmax

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

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

# $\Omega$ -Regularized Argmax

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

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

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

# $\Omega$ -Regularized Argmax

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

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

- 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:

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

- 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

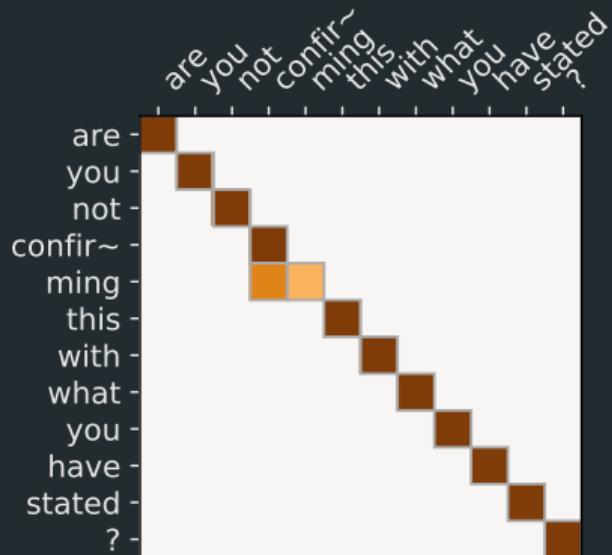
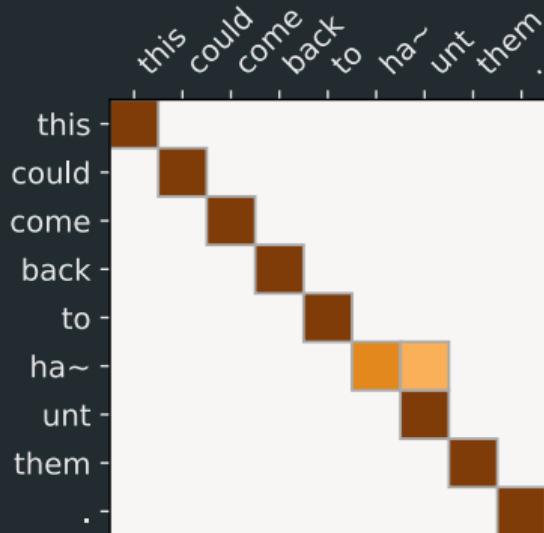
For convex  $\Omega$ , define the  $\Omega$ -regularized argmax transformation:

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

- 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?

# Subword-Merging Head



Learned  $\alpha = 1.91$ .

# Semi-Supervised VAE

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)\end{aligned}$$

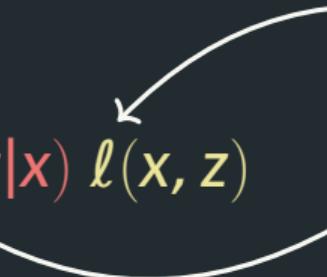
- Semi-Supervised VAE on MNIST:  $z$  is one of 10 categories

# Semi-Supervised VAE

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)\end{aligned}$$

Gaussian VAE

classification network



- Semi-Supervised VAE on MNIST:  $z$  is one of 10 categories

# Semi-Supervised VAE

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)\end{aligned}$$

sum over the 10 digits

Gaussian VAE

classification network

The diagram illustrates the decomposition of the loss function for a semi-supervised VAE. The total loss  $\mathcal{L}_x(\theta)$  is shown as a sum over all latent variables  $z \in \mathcal{Z}$ . Each term in the sum is the product of the Gaussian prior probability  $\pi(z|x)$  and the classification loss  $\ell(x, z)$ . Arrows point from the 'sum over the 10 digits' text to the summation symbol, and from the 'Gaussian VAE' and 'classification network' labels to the two components of the loss term.

- Semi-Supervised VAE on MNIST:  $z$  is one of 10 categories

# Semi-Supervised VAE

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)\end{aligned}$$

sum over  
the 10 digits

Gaussian VAE

classification network

The diagram illustrates the loss function for a semi-supervised VAE. It shows two components: a Gaussian VAE loss (sum over the 10 digits) and a classification network loss. The Gaussian VAE loss is represented by a curved arrow pointing to the sum term in the equation. The classification network loss is represented by a curved arrow pointing to the expectation term. The equation itself is  $\mathcal{L}_x(\theta) = \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) = \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)$ .

- Semi-Supervised VAE on MNIST:  $z$  is one of 10 categories
- Train this with 10% labeled data

# Semi-Supervised VAE

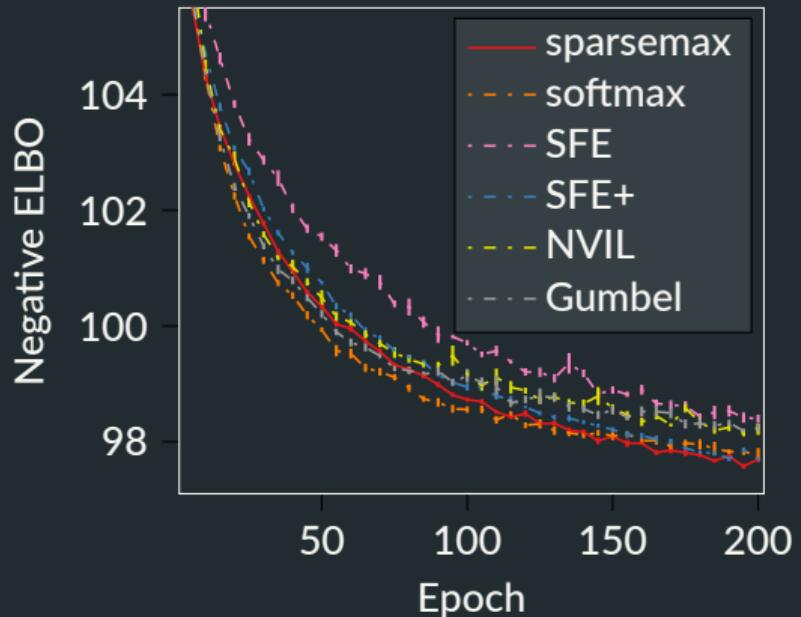
Method	Accuracy (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	$94.75 \pm .002$	1
SFE+	$96.53 \pm .001$	2
NVIL	$96.01 \pm .002$	1
Gumbel	$95.46 \pm .001$	1
<i>Marginalization</i>		
Dense	$96.93 \pm .001$	10

# Semi-Supervised VAE

Method	Accuracy (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	$94.75 \pm .002$	1
SFE+	$96.53 \pm .001$	2
NVIL	$96.01 \pm .002$	1
Gumbel	$95.46 \pm .001$	1
<i>Marginalization</i>		
Dense	$96.93 \pm .001$	10
Sparse	$96.87 \pm .001$	$1.01 \pm 0.01$

# Semi-Supervised VAE

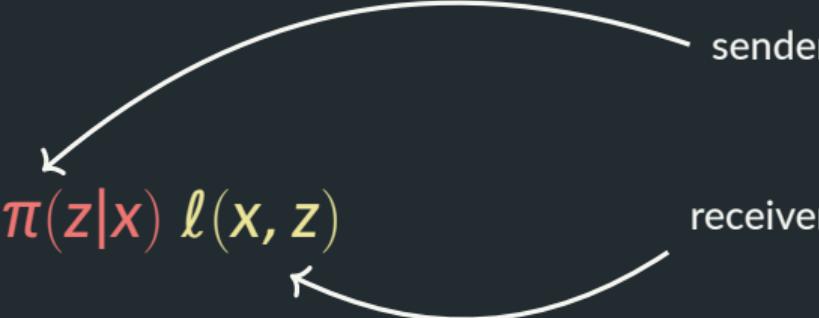
Method	Accuracy (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	$94.75 \pm .002$	1
SFE+	$96.53 \pm .001$	2
NVIL	$96.01 \pm .002$	1
Gumbel	$95.46 \pm .001$	1
<i>Marginalization</i>		
Dense	$96.93 \pm .001$	10
Sparse	$96.87 \pm .001$	$1.01 \pm .01$



# Emergent communication

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)\end{aligned}$$

# Emergent communication

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)\end{aligned}$$


- receiver picks image from a set  $\mathcal{V}$  based on message

# Emergent communication

sum over  
all possible messages  
in the vocabulary

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)\end{aligned}$$

- receiver picks image from a set  $\mathcal{V}$  based on message
- images come from ImageNet

# Emergent Communication

... but make it harder:  $|\mathcal{Z}| = 256$ ,  $|\mathcal{V}| = 16$

---

Method	success (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	33.05 $\pm 2.84$	1
SFE+	44.32 $\pm 2.72$	2
NVIL	37.04 $\pm 1.61$	1
Gumbel	23.51 $\pm 16.19$	1
ST Gumbel	27.42 $\pm 13.36$	1
<i>Marginalization</i>		

---

# Emergent Communication

... but make it harder:  $|\mathcal{Z}| = 256$ ,  $|\mathcal{V}| = 16$

Method	success (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	33.05 $\pm 2.84$	1
SFE+	44.32 $\pm 2.72$	2
NVIL	37.04 $\pm 1.61$	1
Gumbel	23.51 $\pm 16.19$	1
ST Gumbel	27.42 $\pm 13.36$	1
<i>Marginalization</i>		
Dense	93.37 $\pm 0.42$	256

# Emergent Communication

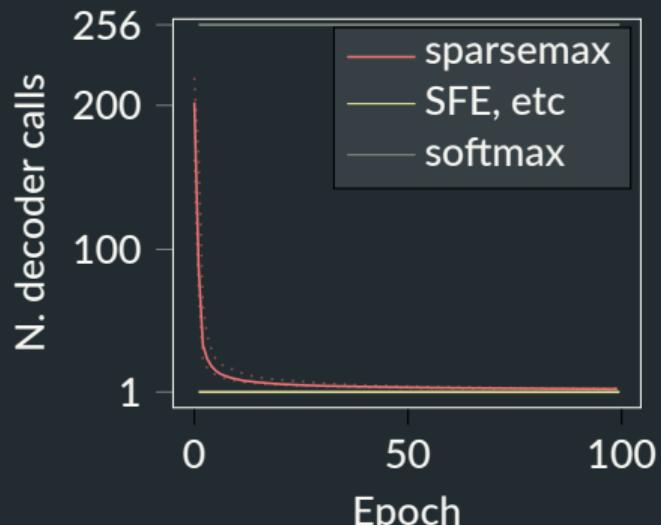
... but make it harder:  $|\mathcal{Z}| = 256$ ,  $|\mathcal{V}| = 16$

Method	success (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	33.05 $\pm 2.84$	1
SFE+	44.32 $\pm 2.72$	2
NVIL	37.04 $\pm 1.61$	1
Gumbel	23.51 $\pm 16.19$	1
ST Gumbel	27.42 $\pm 13.36$	1
<i>Marginalization</i>		
Dense	93.37 $\pm 0.42$	256
Sparse	93.35 $\pm 0.50$	3.13 $\pm 0.48$

# Emergent Communication

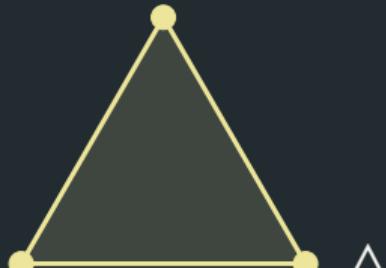
... but make it harder:  $|\mathcal{Z}| = 256$ ,  $|\mathcal{V}| = 16$

Method	success (%)	Dec. calls
<i>Monte Carlo</i>		
SFE	$33.05 \pm 2.84$	1
SFE+	$44.32 \pm 2.72$	2
NVIL	$37.04 \pm 1.61$	1
Gumbel	$23.51 \pm 16.19$	1
ST Gumbel	$27.42 \pm 13.36$	1
<i>Marginalization</i>		
Dense	$93.37 \pm 0.42$	256
Sparse	$93.35 \pm 0.50$	$3.13 \pm 0.48$



 $\Delta$  $\mathcal{M}$

$$\begin{aligned}\mathcal{M} &:= \text{conv} \left\{ \mathbf{a}_z : z \in \mathcal{Z} \right\} \\ &= \left\{ \mathbf{A}p : p \in \Delta \right\} \\ &= \left\{ \mathbb{E}_{Z \sim p} \mathbf{a}_Z : p \in \Delta \right\}\end{aligned}$$

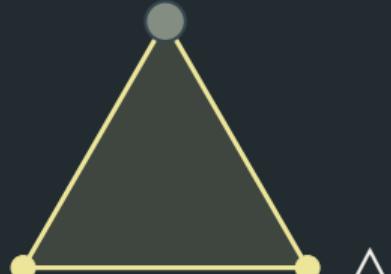


$\Delta$

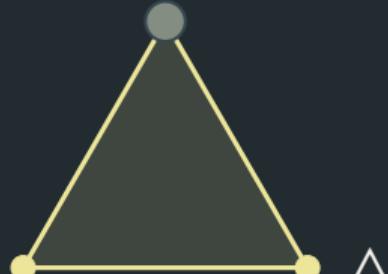


$\mathcal{M}$

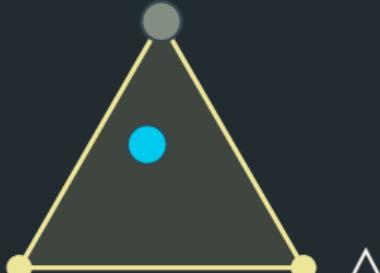
- **argmax**  $\arg \max_{p \in \Delta} p^T s$



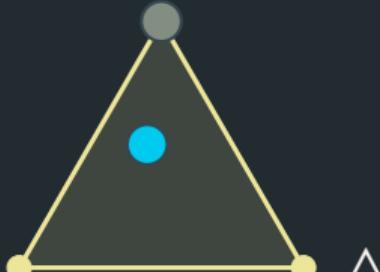
- **argmax**  $\arg \max_{p \in \Delta} p^T s$
- **MAP**  $\arg \max_{\mu \in \mathcal{M}} \mu^T t$



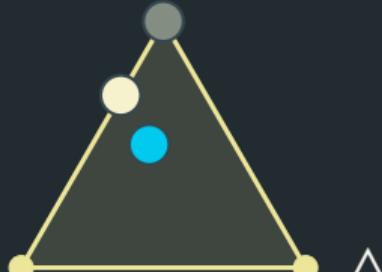
- **argmax**  $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{s}$
- **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{t}$
- **softmax**  $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{s} + H(\boldsymbol{p})$



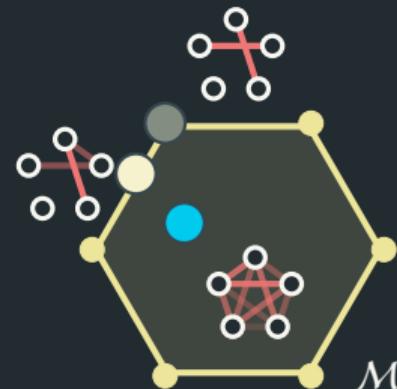
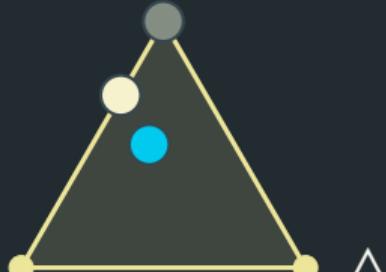
- **argmax**  $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{s}$
- **softmax**  $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{s} + H(\boldsymbol{p})$
- **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{t}$
- **marginals**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{t} + \tilde{H}(\boldsymbol{\mu})$



- **argmax**  $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{s}$
- **softmax**  $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{s} + H(\boldsymbol{p})$
- **sparsemax**  $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{s} - 1/2 \|\boldsymbol{p}\|^2$
- **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{t}$
- **marginals**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{t} + \tilde{H}(\boldsymbol{\mu})$



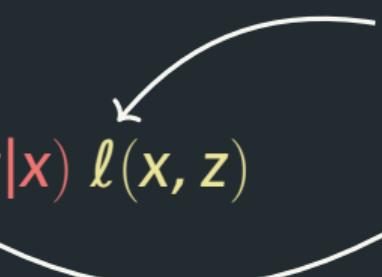
- **argmax**  $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{s}$
- **softmax**  $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{s} + H(\boldsymbol{p})$
- **sparsemax**  $\arg \max_{\boldsymbol{p} \in \Delta} \boldsymbol{p}^\top \boldsymbol{s} - 1/2 \|\boldsymbol{p}\|^2$
- **MAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{t}$
- **marginals**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{t} + \tilde{H}(\boldsymbol{\mu})$
- **SparseMAP**  $\arg \max_{\boldsymbol{\mu} \in \mathcal{M}} \boldsymbol{\mu}^\top \boldsymbol{t} - 1/2 \|\boldsymbol{\mu}\|^2$



# Bit-vector VAE

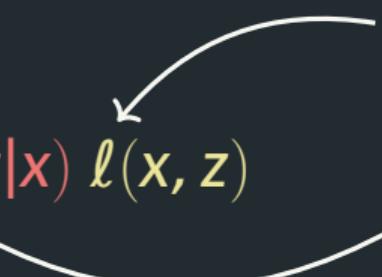
$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)\end{aligned}$$

# Bit-vector VAE

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)\end{aligned}$$


- VAE where  $z$  is a collection of  $D$  bits

# Bit-vector VAE

$$\begin{aligned}\mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z)\end{aligned}$$


- VAE where  $z$  is a collection of  $D$  bits
- Minimize the negative ELBO

# Bit-vector VAE

$$\begin{aligned} \mathcal{L}_x(\theta) &= \sum_{z \in \mathcal{Z}} \pi(z|x) \ell(x, z) \\ &= \mathbb{E}_{z \sim \pi(z|x)} \ell(x, z) \end{aligned}$$

sum over  
an exponentially large  
set of structures

generative network

inference network

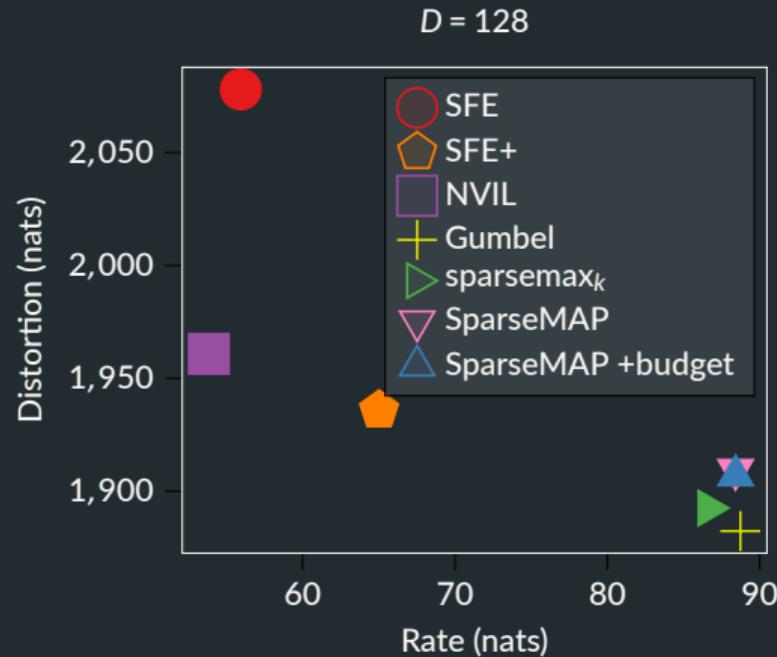
- VAE where  $z$  is a collection of  $D$  bits
- Minimize the negative ELBO

# Bit-vector VAE

Method	$D = 32$	$D = 128$
<i>Monte Carlo</i>		
SFE	3.74	3.77
SFE+	3.61	3.59
NVIL	3.65	3.60
Gumbel	3.57	3.49
<i>Marginalization</i>		
Top-k sparsemax	3.62	3.61
SparseMAP	3.72	3.67
SparseMAP (w/ budget)	3.64	3.66

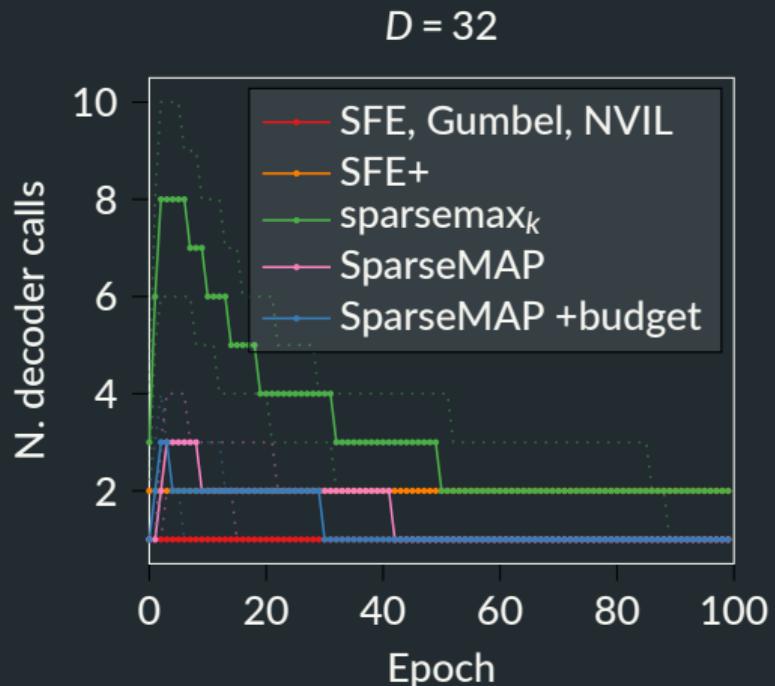
# Bit-vector VAE

Method	$D = 32$	$D = 128$
<i>Monte Carlo</i>		
SFE	3.74	3.77
SFE+	3.61	3.59
NVIL	3.65	3.60
Gumbel	3.57	3.49
<i>Marginalization</i>		
Top-k sparsemax	3.62	3.61
SparseMAP	3.72	3.67
SparseMAP (w/ budget)	3.64	3.66



# Bit-vector VAE

Method	$D = 32$	$D = 128$
<i>Monte Carlo</i>		
SFE	3.74	3.77
SFE+	3.61	3.59
NViL	3.65	3.60
Gumbel	3.57	3.49
<i>Marginalization</i>		
Top-k sparsemax	3.62	3.61
SparseMAP	3.72	3.67
SparseMAP (w/ budget)	3.64	3.66



# Bit-vector VAE

Method	$D = 32$	$D = 128$
<i>Monte Carlo</i>		
SFE	3.74	3.77
SFE+	3.61	3.59
NVIL	3.65	3.60
Gumbel	3.57	3.49
<i>Marginalization</i>		
Top-k sparsemax	3.62	3.61
SparseMAP	3.72	3.67
SparseMAP (w/ budget)	3.64	3.66

