
Why do Transformers work so well?

Fabio James Fehr

18 MARCH 2024



Outline

Personal Background

A Brief History of NLP and Deep Learning

The Attention Mechanism

My PhD Research

Personal Background

A Brief History of NLP and Deep Learning

The Attention Mechanism

My PhD Research

Who is this guy?

UCT

- **BBusSci: Analytics** (2015-2018)
Thesis: Natural Language Processing
- Stefan Britz



Who is this guy?

UCT

- **BBusSci:** Analytics (2015-2018)
Thesis: Natural Language Processing
- Stefan Britz
- **MSc:** Statistics (2019-2020)
Thesis: Nonparametric methods vs
deep learning - Allan Clark



Who is this guy?

UCT

- **BBusSci:** Analytics (2015-2018)
Thesis: Natural Language Processing
- Stefan Britz
- **MSc:** Statistics (2019-2020)
Thesis: Nonparametric methods vs
deep learning - Allan Clark

EPFL & Idiap - Switzerland

- **PhD** Electrical Engineering
(2021-2025)
Nonparametric methods for NLP



Personal Background

A Brief History of NLP and Deep Learning

The Attention Mechanism

My PhD Research

Why should we care?

The “AI Revolution”

- ChatGPT (Text)
- MidJourney (Images)
- AlphaGo (Games)
- Siri (Audio)
- etc ...



Why should we care?

The “AI Revolution”

- ChatGPT (Text)
- MidJourney (Images)
- AlphaGo (Games)
- Siri (Audio)
- etc ...



Why should we care?

The “AI Revolution”

- ChatGPT (Text)
- MidJourney (Images)
- AlphaGo (Games)
- Siri (Audio)
- etc ...

What is the secret sauce?



Why should we care?

The “AI Revolution”

- ChatGPT (Text)
- MidJourney (Images)
- AlphaGo (Games)
- Siri (Audio)
- etc ...

What is the secret sauce?

- The attention mechanism (NLP)
- Large-scale pretraining (Deep Learning)



Problem: Machine Translation

Given text examples of French and English language pairs, translate the following:

Il m'a entarté = ???
French *English*

How do we solve translation? (2000s)

Tokenization

- Break into smaller units (words)

Je suis un chat. I am a cat.
└─┘ └─┘ └─┘ └─┘ └─┘ └─┘ └─┘ └─┘

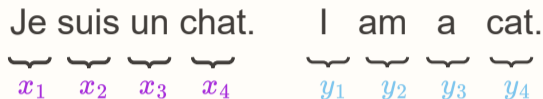
How do we solve translation? (2000s)

Tokenization

- Break into smaller units (words)
- Build vocabulary

Je suis un chat. I am a cat.

x_1 x_2 x_3 x_4 y_1 y_2 y_3 y_4



How do we solve translation? (2000s)

Tokenization

- Break into smaller units (words)
- Build vocabulary

Vectorisation

- Make into numbers (0 or 1)

Je suis un chat.

x_1	x_2	x_3	x_4
1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1
\vdots	\vdots	\vdots	\vdots

I am a cat.

y_1	y_2	y_3	y_4
1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1
\vdots	\vdots	\vdots	\vdots

How do we solve translation? (2000s)

Tokenization

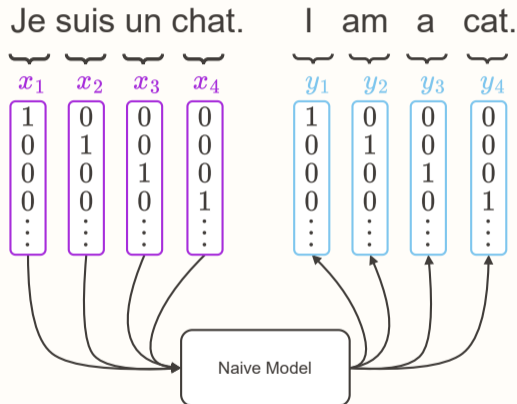
- Break into smaller units (words)
- Build vocabulary

Vectorisation

- Make into numbers (0 or 1)

Classification

- Multinomial logistic regression (per word) or Naive Bayes



How do we solve translation? (2000s)

Tokenization

- Break into smaller units (words)
- Build vocabulary

Vectorisation

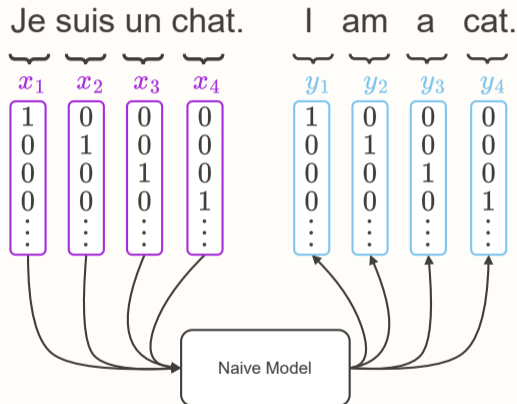
- Make into numbers (0 or 1)

Semantics?

Classification

- Multinomial logistic regression (per word) or Naive Bayes

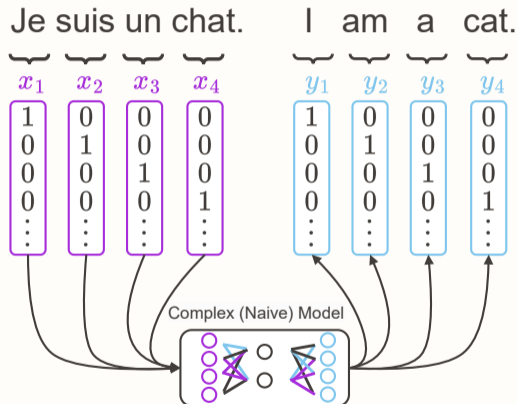
Simplistic? Independence?



How do we solve translation? (<2014)

Models are too simple?

- Neural networks



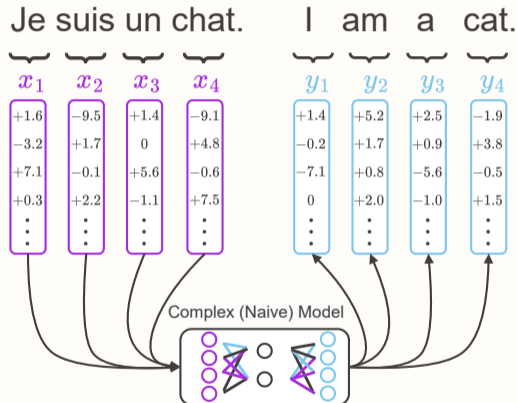
How do we solve translation? (<2014)

Models are too simple?

- Neural networks

Contextualised representations?

- Word2Vec [[Mikolov et al., 2013](#)]



How do we solve translation? (<2014)

Models are too simple?

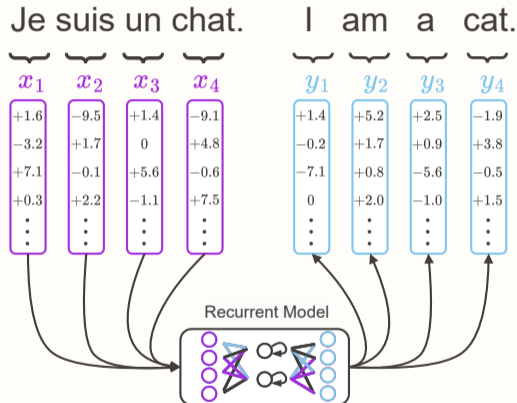
- Neural networks

Contextualised representations?

- Word2Vec [Mikolov et al., 2013]

Models assume independence?

- Recurrent Neural Networks



How do we solve translation? (<2014)

Models are too simple?

- Neural networks

Contextualised representations?

- Word2Vec [Mikolov et al., 2013]

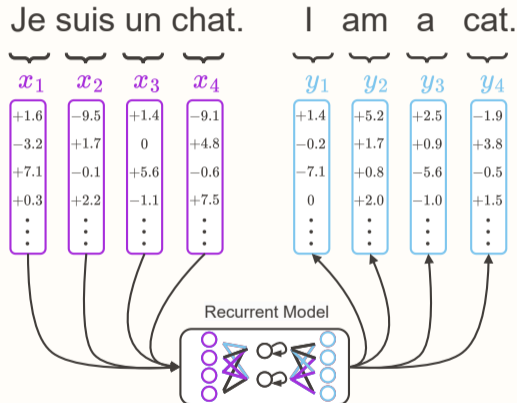
Models assume independence?

- Recurrent Neural Networks

Long term dependencies?

Recurrence on GPUs?

Vanishing gradients?



How do we solve translation? (2014 and beyond)

Long term dependencies?

- Attention for translation
[Bahdanau et al., 2014]

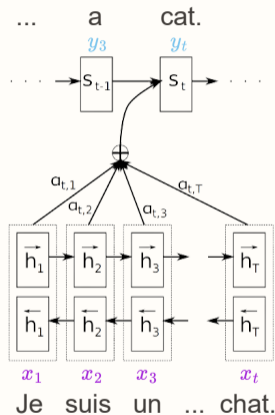


Figure: Bidirectional LSTM with attention

How do we solve translation? (2014 and beyond)

Long term dependencies?

- Attention for translation
[Bahdanau et al., 2014]

Recurrence?

- Transformers - only attention
[Vaswani et al., 2017]

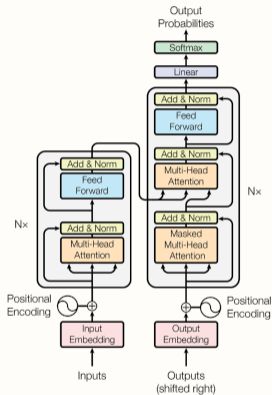


Figure: Transformer

How do we solve translation? (2014 and beyond)

Long term dependencies?

- Attention for translation
[Bahdanau et al., 2014]

Recurrence?

- Transformers - only attention
[Vaswani et al., 2017]

Scale

- The bitter lesson [Sutton, 2019]



noooooo you can't just scale up pure connectionist models on Internet data without inductive biases and modularization and expect them to learn real-world knowledge and grammar from form, or arithmetic and logical reasoning and causal inference—that's just memorization and superficial pattern-matching like Eliza, you need grounding in real-world communication with intent and social dynamics and multimodal robotic embodiment which can foster disentangled learning from guided exploration and self-directed goals expressed in Bayesian programs and probabilistic graphical models which are interpretable and pin down a unique semantics which can be debiased and expressed with uncertainty, and learned efficiently on tiny academic datasets. the next only show how this is a dead-end, we need to stop chasing WGA and avoid the complexity of the brain and consider the social context to motivate AI's structural learning for 1000 World memories...



haha gpus go bitterrrr

How do we solve translation? (2014 and beyond)

Long term dependencies?

- Attention for translation
[Bahdanau et al., 2014]

Recurrence?

- Transformers - only attention
[Vaswani et al., 2017]

Scale

- The bitter lesson [Sutton, 2019]

Pretraining

- “Large” foundation models
(2017-2021)

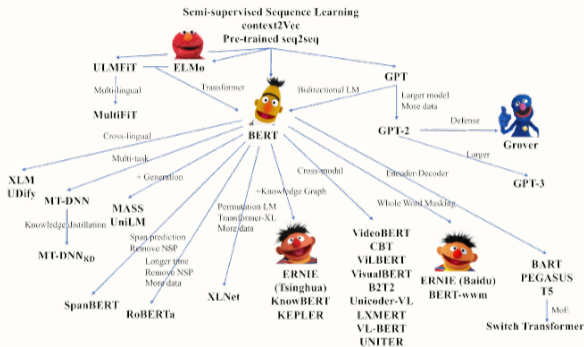


Figure: Pretrained models [Han et al., 2021]

How do we solve translation? (2014 and beyond)

Long term dependencies?

- Attention for translation
[Bahdanau et al., 2014]

Recurrence?

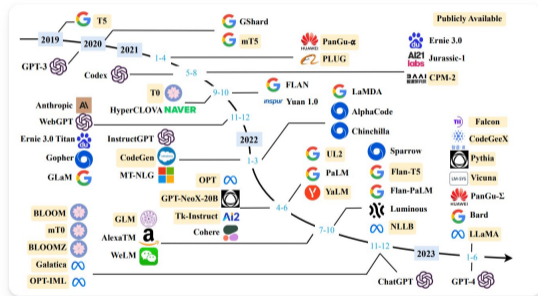
- Transformers - only attention
[Vaswani et al., 2017]

Scale

- The bitter lesson [Sutton, 2019]

Pretraining

- “Large” foundation models (2017-2021)
- Large Language Models (2023)



Takeaways:

The secret sauce:

1. The attention mechanism
2. Scaling data and model size

Takeaways:

The secret sauce:

1. The attention mechanism
2. Scaling data and model size

Il m'a entarté = ???
French *English*

Takeaways:

The secret sauce:

1. The attention mechanism
2. Scaling data and model size

Il m'a entarté = ???
French *English*



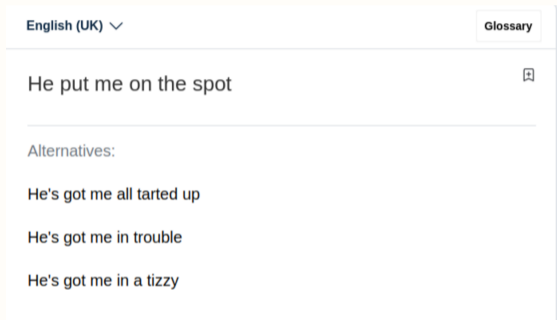
Figure: Google translate

Takeaways:

The secret sauce:

1. The attention mechanism
2. Scaling data and model size

Il m'a entarté = ???
French *English*



English (UK) ▾ Glossary

He put me on the spot 🔖

Alternatives:

He's got me all tarted up

He's got me in trouble

He's got me in a tizzy

Figure: DeepL

Takeaways:

The secret sauce:

1. The attention mechanism
2. Scaling data and model size

Il m'a entarté = ???
French *English*

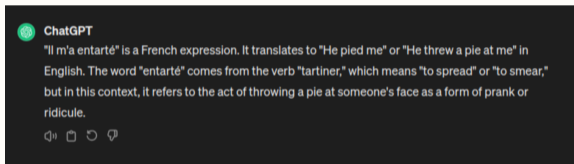


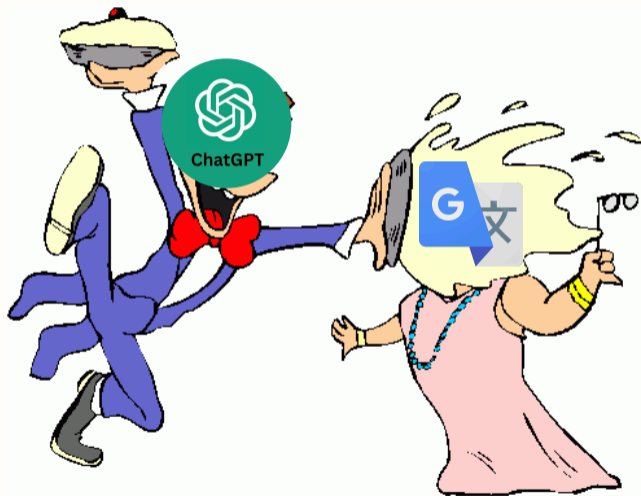
Figure: ChatGPT

Takeaways:

The secret sauce:

1. The attention mechanism
2. Scaling data and model size

Il m'a entarté = ???
French *English*



Personal Background

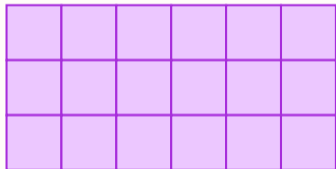
A Brief History of NLP and Deep Learning

The Attention Mechanism

My PhD Research

Data

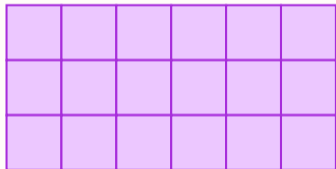
$$\mathbf{Z} \in \mathbb{R}^{N \times P}$$



- N observations, P dimensional

Data

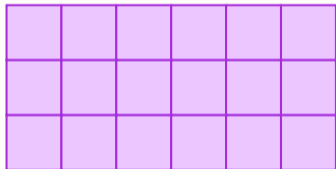
$$\mathbf{Z} \in \mathbb{R}^{N \times P}$$



- N observations, P dimensional
- No particular order

Data

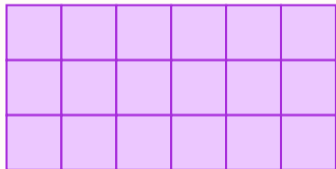
$$\mathbf{Z} \in \mathbb{R}^{N \times P}$$



- N observations, P dimensional
- No particular order
- Common across domains

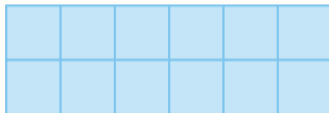
Data

$$\mathbf{Z} \in \mathbb{R}^{N \times P}$$



- N observations, P dimensional
- No particular order
- Common across domains

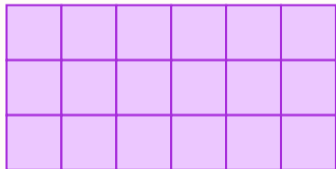
$$\mathbf{U} \in \mathbb{R}^{M \times P}$$



- M observations, P dimensional

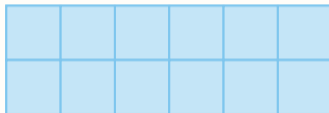
Data

$$\mathbf{Z} \in \mathbb{R}^{N \times P}$$



- N observations, P dimensional
- No particular order
- Common across domains

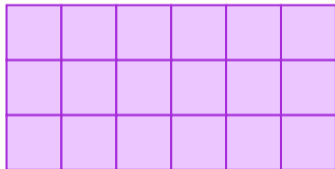
$$\mathbf{U} \in \mathbb{R}^{M \times P}$$



- M observations, P dimensional
- No particular order

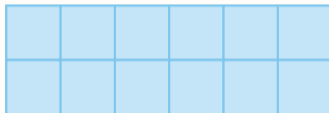
Data

$$\mathbf{Z} \in \mathbb{R}^{N \times P}$$



- N observations, P dimensional
- No particular order
- Common across domains

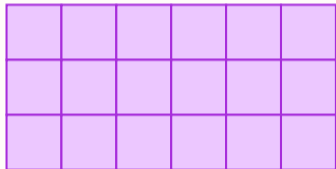
$$\mathbf{U} \in \mathbb{R}^{M \times P}$$



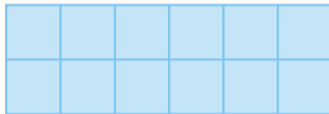
- M observations, P dimensional
- No particular order
- $\mathbf{U} = \mathbf{Z}$ or $\mathbf{U} \neq \mathbf{Z}$

How do we get Z and U to interact?

$$Z \in \mathbb{R}^{N \times P}$$

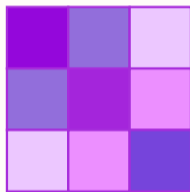


$$U \in \mathbb{R}^{M \times P}$$



How do we get Z and U to interact?

$$\mathbf{Z}\mathbf{Z}^T \in \mathbb{R}_+^{N \times N}$$



$$\mathbf{U}\mathbf{Z}^T \in \mathbb{R}^{M \times N}$$

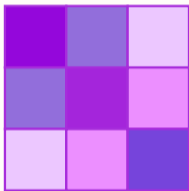


- $U = Z$: $Variance(\mathbf{Z}, \mathbf{Z})$

- $U \neq Z$: $Covariance(\mathbf{U}, \mathbf{Z})$

How do we project this interaction forward?

$$\mathbf{Z}\mathbf{Z}^T \in \mathbb{R}_+^{N \times N}$$



- $U = \mathbf{Z}$: *Variance*(\mathbf{Z}, \mathbf{Z})

$$\mathbf{U}\mathbf{Z}^T \in \mathbb{R}^{M \times N}$$



- $U \neq \mathbf{Z}$: *Covariance*(\mathbf{U}, \mathbf{Z})

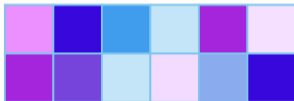
How do we project this interaction forward?

$$\mathbf{Z}\mathbf{Z}^T\mathbf{Z} \in \mathbb{R}^{N \times P}$$



- $U = Z$: *Variance*(Z, Z)
- Project by Z

$$\mathbf{U}\mathbf{Z}^T\mathbf{Z} \in \mathbb{R}^{M \times P}$$



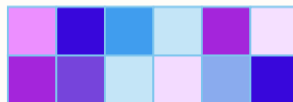
- $U \neq Z$: *Covariance*(U, Z)
- Project by Z

What are the problems with this?

$$\mathbf{Z}\mathbf{Z}^T\mathbf{Z} \in \mathbb{R}^{N \times P}$$



$$\mathbf{U}\mathbf{Z}^T\mathbf{Z} \in \mathbb{R}^{M \times P}$$

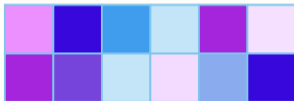


What are the problems with this?

$$\mathbf{Z}\mathbf{Z}^T\mathbf{Z} \in \mathbb{R}^{N \times P}$$



$$\mathbf{U}\mathbf{Z}^T\mathbf{Z} \in \mathbb{R}^{M \times P}$$



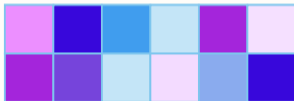
- Strictly positive multiplications?

What are the problems with this?

$$\mathbf{Z}\mathbf{Z}^T\mathbf{Z} \in \mathbb{R}^{N \times P}$$



$$\mathbf{U}\mathbf{Z}^T\mathbf{Z} \in \mathbb{R}^{M \times P}$$



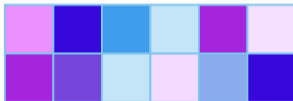
- Strictly positive multiplications?
- Normalisations and scaling?

What are the problems with this?

$$\mathbf{Z}\mathbf{Z}^T\mathbf{Z} \in \mathbb{R}^{N \times P}$$



$$\mathbf{U}\mathbf{Z}^T\mathbf{Z} \in \mathbb{R}^{M \times P}$$



- Strictly positive multiplications?
- Normalisations and scaling?

- Single interaction value?

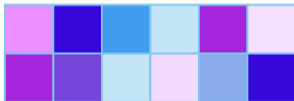
What are the problems with this?

$$\mathbf{Z}\mathbf{Z}^T\mathbf{Z} \in \mathbb{R}^{N \times P}$$



- Strictly positive multiplications?
- Normalisations and scaling?

$$\mathbf{U}\mathbf{Z}^T\mathbf{Z} \in \mathbb{R}^{M \times P}$$



- Single interaction value?
- Quadratic?

The attention mechanism

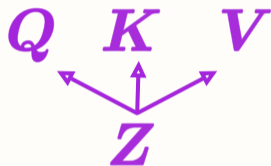


Figure: Self attention

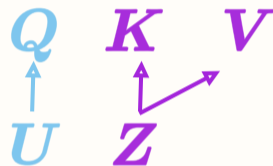


Figure: Cross attention

- $W^Q, W^K, W^V \in \mathbb{R}^{P \times d}$

The attention mechanism

- Cross attention

$$QK^T$$

$$\in \mathbb{R}^{M \times N}$$



The attention mechanism

- Cross attention
- Scaling

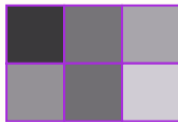
$$\left(\frac{QK^T}{\sqrt{d}} \right) \in \mathbb{R}^{M \times N}$$



The attention mechanism

- Cross attention
- Scaling
- Normalisation

$$\text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right) \in \mathbb{R}^{M \times N}$$



The attention mechanism

- Cross attention
- Scaling
- Normalisation
- Projection

$$\text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V \in \mathbb{R}^{M \times d}$$



Why is attention so cool?

Advantages:

- Simple layer
- No particular order
- Unbounded inputs
- Probabilistic interpretation

$$\text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V \in \mathbb{R}^{M \times d}$$



Why is attention so cool?

Advantages:

- Simple layer
- No particular order
- Unbounded inputs
- Probabilistic interpretation

Problems solved:

- Strictly positive multiplications? ✓

$$\text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V \in \mathbb{R}^{M \times d}$$



Why is attention so cool?

Advantages:

- Simple layer
- No particular order
- Unbounded inputs
- Probabilistic interpretation

Problems solved:

- Strictly positive multiplications? ✓
- Normalisations and scaling? ✓

$$\text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V \in \mathbb{R}^{M \times d}$$



Why is attention so cool?

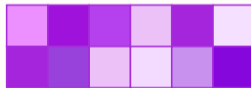
Advantages:

- Simple layer
- No particular order
- Unbounded inputs
- Probabilistic interpretation

Problems solved:

- Strictly positive multiplications? ✓
- Normalisations and scaling? ✓
- Single interaction value? ✓

$$\text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V \in \mathbb{R}^{M \times d}$$



Why is attention so cool?

Advantages:

- Simple layer
- No particular order
- Unbounded inputs
- Probabilistic interpretation

Problems solved:

- Strictly positive multiplications? ✓
- Normalisations and scaling? ✓
- Single interaction value? ✓
- Quadratic? ✗

$$\text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V \in \mathbb{R}^{M \times d}$$



Personal Background

A Brief History of NLP and Deep Learning

The Attention Mechanism

My PhD Research

Intuition

Attention mechanism:

- No particular order
- Unbounded inputs
- Probabilistic interpretation

$$\text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V \in \mathbb{R}^{M \times d}$$



Intuition

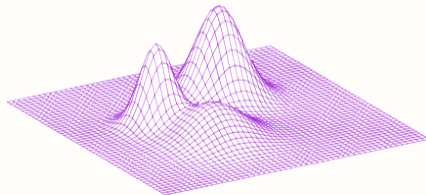
Attention mechanism:

- No particular order
- Unbounded inputs
- Probabilistic interpretation

Nonparametric distributions

- Exchangable set
- Theoretically infinite
- Mixture of distributions

$$\text{Softmax} \left(\frac{QK^T}{\sqrt{d}} \right) V \in \mathbb{R}^{M \times d}$$



Nonparametric latent variable modelling

Regularisation:

- Generalisation
- Sparse representations



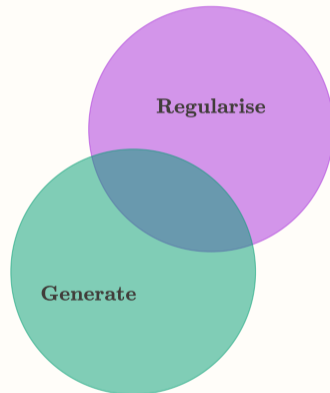
Nonparametric latent variable modelling

Regularisation:

- Generalisation
- Sparse representations

Generation:

- Generative modelling



Nonparametric latent variable modelling

Regularisation:

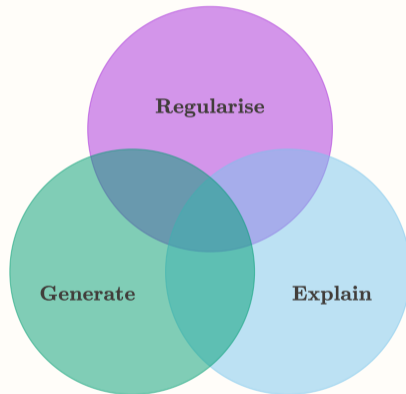
- Generalisation
- Sparse representations

Generation:

- Generative modelling

Explainability:

- Disentanglement



Nonparametric latent variable modelling

Regularisation:

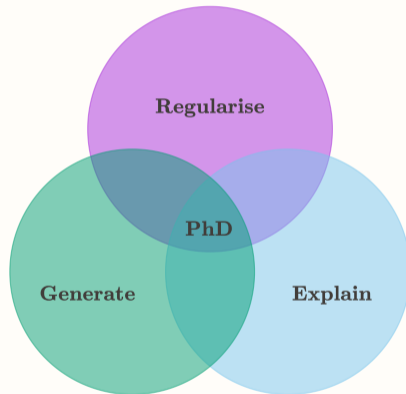
- Generalisation
- Sparse representations

Generation:

- Generative modelling

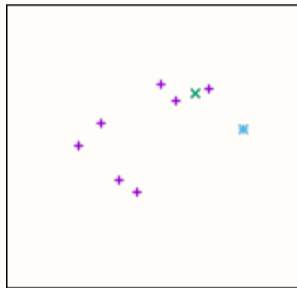
Explainability:

- Disentanglement



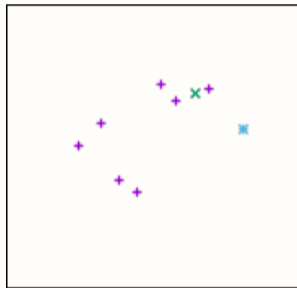
The denoising attention mechanism

$$\text{Attn}(QKV) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V$$



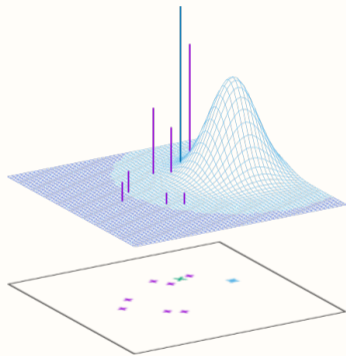
The denoising attention mechanism

$$\text{Attn}(\mathbf{U} \mathbf{Z}) = \text{Softmax} \left(\frac{\mathbf{U} \mathbf{Z}^T}{\sqrt{d}} \right) \mathbf{Z}$$



The denoising attention mechanism

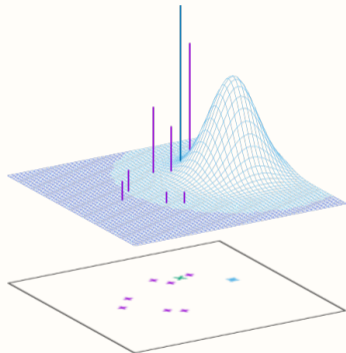
$$\text{Attn}(\mathbf{U}, \mathbf{Z}) = \text{Softmax}\left(\frac{\mathbf{U}\mathbf{Z}^T}{\sqrt{d}}\right)\mathbf{Z}$$



The denoising attention mechanism

$$\text{Attn}(\mathbf{U} \mathbf{Z}) = \text{Softmax} \left(\frac{\mathbf{U} \mathbf{Z}^T}{\sqrt{d}} \right) \mathbf{Z}$$

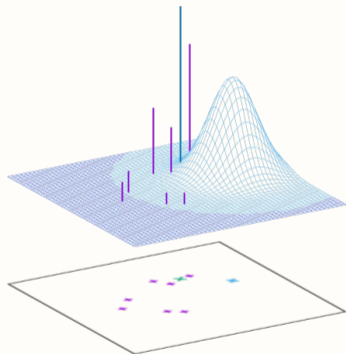
- Prior information \mathbf{Z}



The denoising attention mechanism

$$\text{Attn}(\mathbf{U}, \mathbf{Z}) = \text{Softmax}\left(\frac{\mathbf{U}\mathbf{Z}^T}{\sqrt{d}}\right)\mathbf{Z}$$

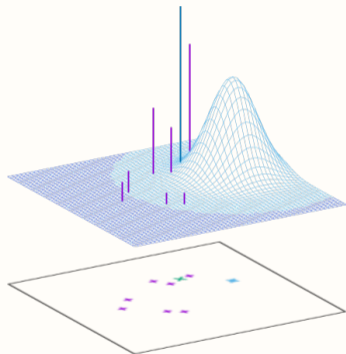
- Prior information \mathbf{Z}
- Noisy query \mathbf{U}



The denoising attention mechanism

$$\text{Attn}(\mathbf{U} \mathbf{Z}) = \text{Softmax} \left(\frac{\mathbf{U} \mathbf{Z}^T}{\sqrt{d}} \right) \mathbf{Z}$$

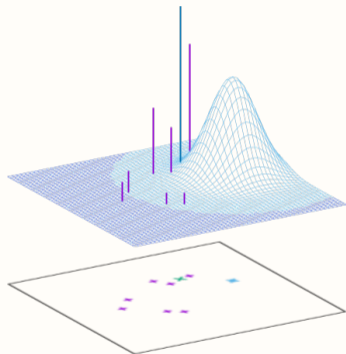
- Prior information \mathbf{Z}
- Noisy query \mathbf{U}
- Posterior update $\text{Attn}(\mathbf{U} \mathbf{Z})$



The denoising attention mechanism

$$\text{Attn}(\mathbf{U} \mathbf{Z}) = \text{Softmax} \left(\frac{\mathbf{U} \mathbf{Z}^T}{\sqrt{d}} \right) \mathbf{Z}$$

- Prior information \mathbf{Z}
- Noisy query \mathbf{U}
- Posterior update $\text{Attn}(\mathbf{U} \mathbf{Z})$
- Proof of concept [[Henderson and Fehr, 2023](#)]



Current research

Theory:

- A VAE for Transformers
[[Henderson and Fehr, 2023](#)]

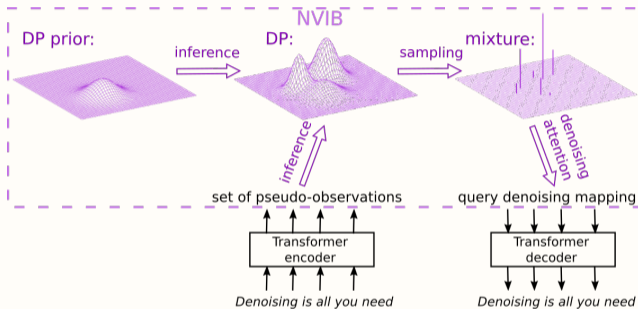


Figure: Nonparametric Variational Information Bottleneck (NVIB)

Current research

Theory:

- A VAE for Transformers [Henderson and Fehr, 2023]

Regularisation:

- Post-training regularisation [Fehr and Henderson, 2023]

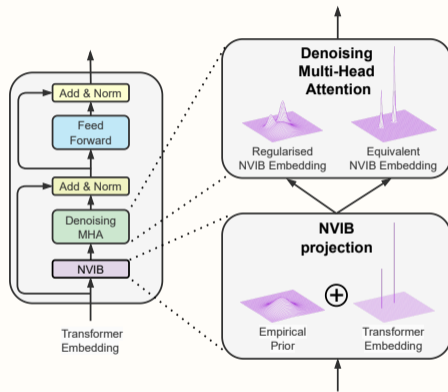


Figure: Nonparametric Variational Regularisation

Current research

Theory:

- A VAE for Transformers [Henderson and Fehr, 2023]

Regularisation:

- Post-training regularisation [Fehr and Henderson, 2023]

Explainability:

- Abstraction [Behjati et al., 2023]

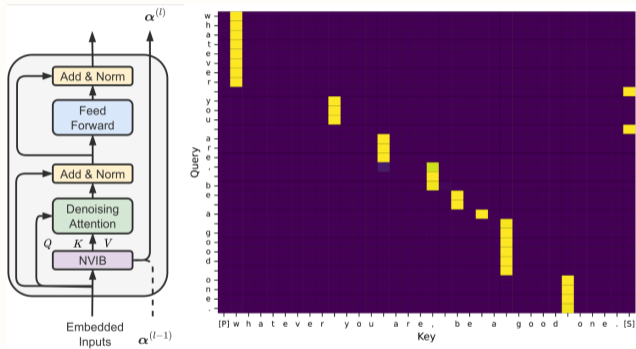


Figure: Learned layer-wise abstraction

Current research

Theory:

- A VAE for Transformers
[Henderson and Fehr, 2023]

Regularisation:

- Post-training regularisation
[Fehr and Henderson, 2023]

Explainability:

- Abstraction
[Behjati et al., 2023]

Generation:

- Coming soon 2024!

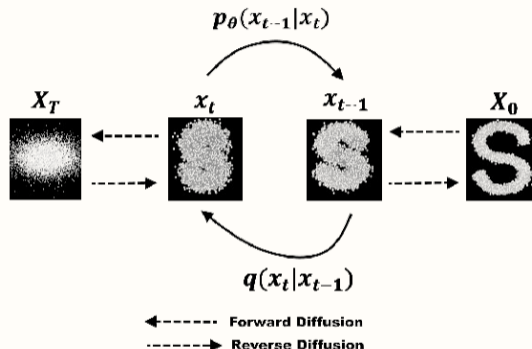


Figure: Latent attention-based diffusion for text

Fin

Personal Background

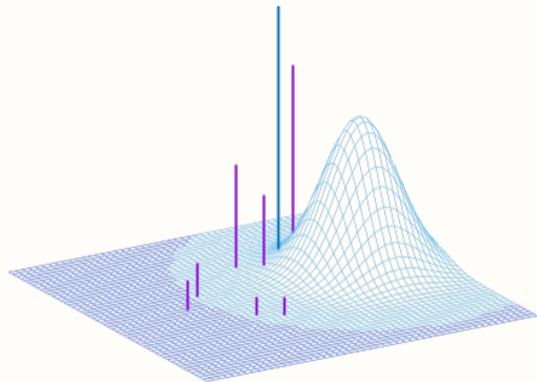
A Brief History of NLP and Deep Learning

The Attention Mechanism

My PhD Research

Fin

Personal Background
A Brief History of NLP and Deep Learning
The Attention Mechanism
My PhD Research



Thank you for your attention!

References I



Bahdanau, D., Cho, K., and Bengio, Y. (2014).

Neural machine translation by jointly learning to align and translate.
[arXiv preprint arXiv:1409.0473](https://arxiv.org/abs/1409.0473).



Behjati, M., Fehr, F. J., and Henderson, J. (2023).

Learning to abstract with nonparametric variational information bottleneck.
In [The 2023 Conference on Empirical Methods in Natural Language Processing](#).



Fehr, F. and Henderson, J. (2023).

Nonparametric variational regularisation of pretrained transformers.



Han, X., Zhang, Z., Ding, N., Gu, Y., Liu, X., Huo, Y., Qiu, J., Zhang, L., Han, W., Huang, M., Jin, Q., Lan, Y., Liu, Y., Liu, Z., Lu, Z., Qiu, X., Song, R., Tang, J., rong Wen, J., Yuan, J., Zhao, W. X., and Zhu, J. (2021).

Pre-trained models: Past, present and future.
[AI Open](#), 2:225–250.



Henderson, J. and Fehr, F. (2023).

A VAE for Transformers with Nonparametric Variational Information Bottleneck.
In [International Conference on Learning Representations](#).



Mikolov, T., Chen, K., Corrado, G. S., and Dean, J. (2013).

Efficient estimation of word representations in vector space.
In [International Conference on Learning Representations](#).

References II



Sutton, R. (2019).

The bitter lesson.

[March, 13:2019.](#)



Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. (2017).

Attention is all you need.

In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, [Advances in Neural Information Processing Systems](#), volume 30. Curran Associates, Inc.