

# Mechanistic Interpretability

Fundamentals  
WT 2024/25

Frederick Riemenschneider



24.10.2024

Recap

A Mathematical  
Framework

Preliminaries

One-Layer Attention-Only  
Transformers

References

Recap

A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

References

# Recap

# attention transformers

"the transformer"

sequence-to-sequence learning

BERT

T5

encoder

GPT

decoder

LSTM

## Recap

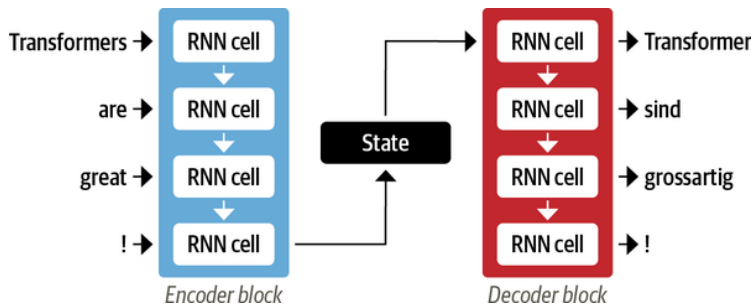
### A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

### References

# Encoder-decoder Architecture



**Figure 1:** Tunstall et al. 2022, p. 4.

## Recap

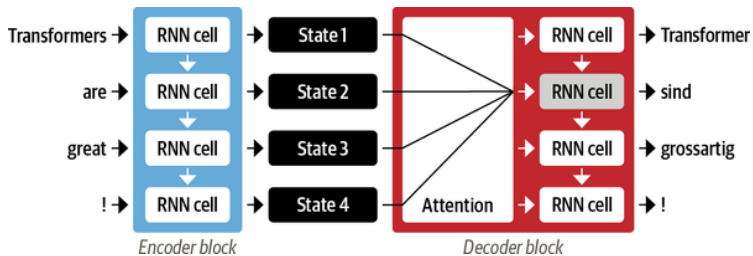
### A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

### References

# Attention



**Figure 2:** Tunstall et al. 2022, p. 5.

## Recap

### A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

### References

# Attention

- context as a weighted average:

$$\mathbf{c}_i = \sum_j^{T^e} \alpha_{i,j} \mathbf{h}_j^e$$

- normalization via softmax:

$$\alpha_{i,j} = \frac{\exp(e_{i,j})}{\sum_k^{T^e} \exp(e_{i,k})}$$

- importance of  $\mathbf{h}_j^e$  for  $\mathbf{h}_{i-1}^d$ :

$$e_{i,j} = a(\mathbf{h}_{i-1}^d, \mathbf{h}_j^e)$$

## Recap

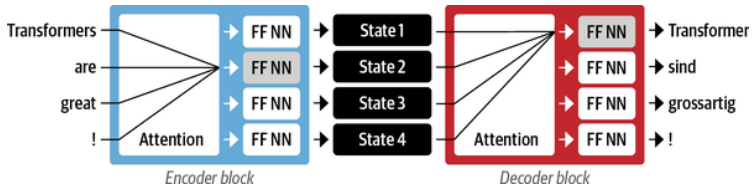
### A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

### References

# Self-Attention



**Figure 3:** Tunstall et al. 2022, p. 6.

## Recap

### A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

### References

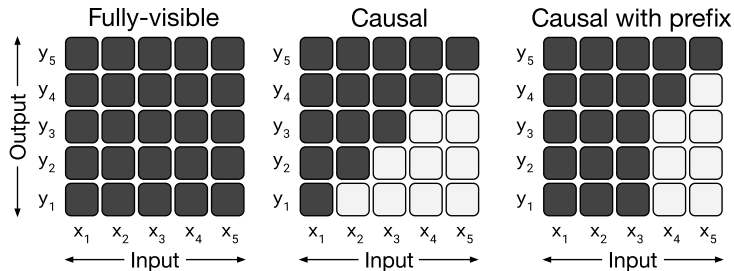
## Recap

## A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

## References



**Figure 4:** Raffel et al. 2020.



Recap

**A Mathematical Framework**

Preliminaries

One-Layer Attention-Only Transformers

References

Which pre-training objectives do you know?

Recap

**A Mathematical  
Framework**

Preliminaries

One-Layer Attention-Only  
Transformers

References

# A Mathematical Framework for Transformer Circuits

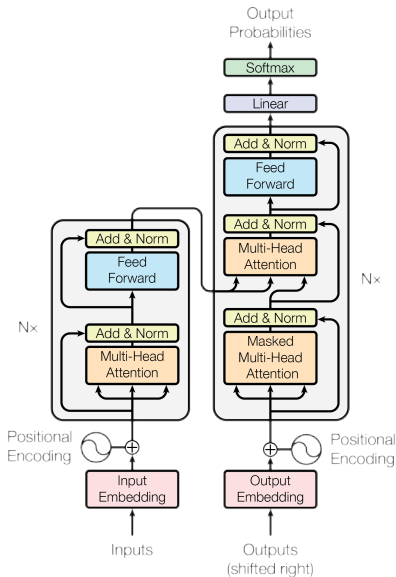
## Recap

## A Mathematical Framework

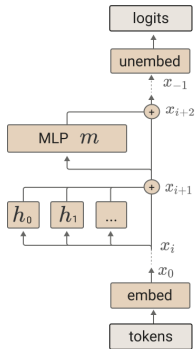
Preliminaries

One-Layer Attention-Only Transformers

## References



**Figure 5:** Vaswani et al. 2017.



The final logits are produced by applying the unembedding.

$$T(t) = W_U x_{-1}$$

An MLP layer,  $m$ , is run and added to the residual stream.

$$x_{i+2} = x_{i+1} + m(x_{i+1})$$

Each attention head,  $h$ , is run and added to the residual stream.

$$x_{i+1} = x_i + \sum_{h \in H_i} h(x_i)$$

Token embedding.

$$x_0 = W_E t$$

One  
residual  
block

**Figure 6:** Elhage et al. 2021.

## Recap

### A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

### References

# Residual Stream as Communication Channel

- main idea: attention heads and MLPs *add information* to the residual stream
- compositionality, responsibility splitting
- residual stream hardly interpretable, components adding to it may be interpretable
- **We (probably) don't want to interpret the residual stream!**

## Recap

### A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

## References

# Attention Heads are Independent and Additive

1) This is our input sentence\*

2) We embed each word\*

3) Split into 8 heads.  
We multiply  $X$  or  $R$  with weight matrices

4) Calculate attention using the resulting  $Q/K/V$  matrices

5) Concatenate the resulting  $Z$  matrices, then multiply with weight matrix  $W^O$  to produce the output of the layer

Thinking  
Machines



\* In all encoders other than #0, we don't need embedding. We start directly with the output of the encoder right below this one

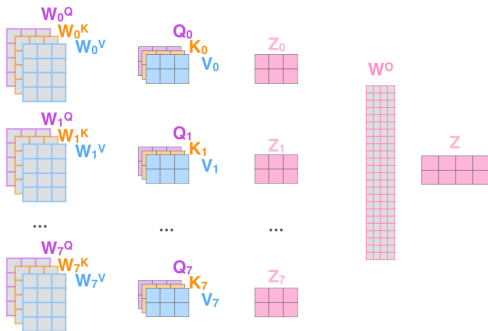


Figure 7: Alammar 2018.

## Recap

### A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

## References

# Attention Heads as Information Movement

Each vector receives three representations ("roles")

$$\begin{bmatrix} W_Q \end{bmatrix} \times \begin{bmatrix} \bullet \\ \bullet \\ \bullet \end{bmatrix} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \end{bmatrix}$$

**Query:** vector **from** which the attention is looking

"Hey there, do you have this information?"

$$\begin{bmatrix} W_K \end{bmatrix} \times \begin{bmatrix} \bullet \\ \bullet \\ \bullet \end{bmatrix} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \end{bmatrix}$$

**Key:** vector **at** which the query looks to compute weights

"Hi, I have this information - give me a large weight!"

$$\begin{bmatrix} W_V \end{bmatrix} \times \begin{bmatrix} \bullet \\ \bullet \\ \bullet \end{bmatrix} = \begin{bmatrix} \bullet \\ \bullet \\ \bullet \end{bmatrix}$$

**Value:** their weighted sum is attention output

"Here's the information I have!"

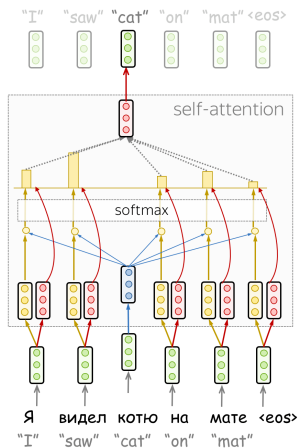


Figure 8: Voita 2023.

## Recap

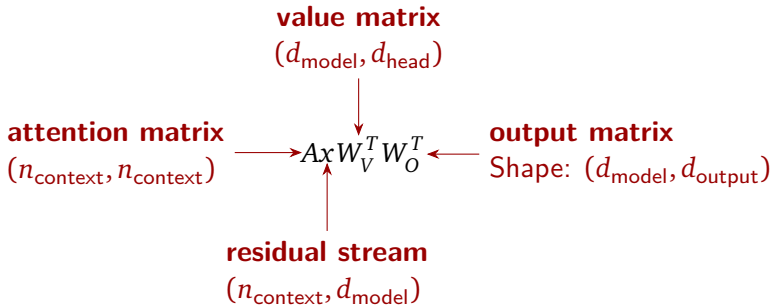
## A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

## References

# Attention Heads as Information Movement



## Recap

## A Mathematical Framework

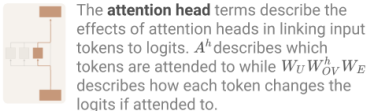
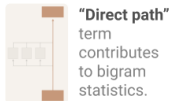
Preliminaries

One-Layer Attention-Only Transformers

## References



$$T = \text{Id} \otimes W_U W_E + \sum_{h \in H} A^h \otimes (W_U W_{OV}^h W_E)$$



$$T(x) = W_U W_E x^T + \sum_{h \in H} A^h x (W_U W_{OV}^h W_E)^T$$

## Recap

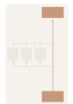
## A Mathematical Framework

Preliminaries

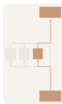
One-Layer Attention-Only Transformers

## References

$$T = \text{Id} \otimes W_U W_E + \sum_{h \in H} A^h \otimes (W_U W_{OV}^h W_E)$$



“Direct path” term contributes to bigram statistics.



The **attention head** terms describe the effects of attention heads in linking input tokens to logits.  $A^h$  describes which tokens are attended to while  $W_U W_{OV}^h W_E$  describes how each token changes the logits if attended to.

$$T(x) = W_U W_E x^T + \sum_{h \in H} A^h x (W_U W_{OV}^h W_E)^T$$

We can look at each attention head independently. Attention is the only communication possibility, enabling skip-trigrams.

## Recap

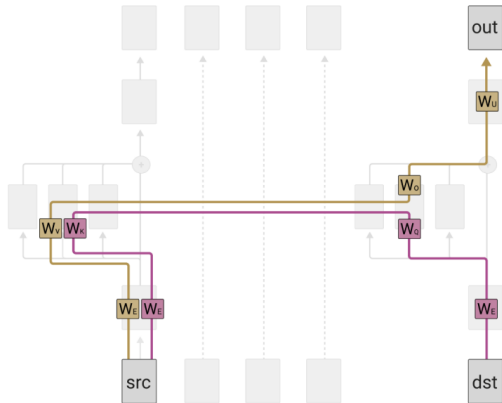
## A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

## References

# Attention Heads as Information Movement



The OV ("output-value") circuit determines how attending to a given token affects the logits.

$$W_U W_O W_V W_E$$

The QK ("query-key") circuit controls which tokens the head prefers to attend to.

$$W_E^T W_Q^T W_K W_E$$

## Recap

### A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

### References

# Skip-Trigrams

## Some examples of large entries QK/OV circuit

Source Token	Destination Token	Out Token	Example Skip Tri-grams
"perfect"	"are", "looks", "is", "provides"	"perfect", "super", "absolute", "pure"	"perfect... are perfect", "perfect... looks super"
"large"	"contains", "using", "specify", "contain"	"large", "small", "very", "huge"	"large... using large", "large... contains small"
"two"	"One", "\n ", "has", "\r\n ", "One"	"two", "three", "four", "five", "one"	"two... One two", "two... has three"
"lambda"	"\$\\", "}{\\", "+\\"", "(\\", "\${\\"	"lambda", "sorted", "lambda", "operator"	"lambda... \$\lambda", "lambda... +\lambda"
"nbsp"	"&", "\&", "}&", ">&", "=&"	"nbsp", "01", "gt", "00012", "nbs", "quot"	"nbsp... &nbsp", "nbsp... >&nbsp"
"Great"	"The", "The", "the", "contains", "/"	"Great", "great", "poor", "Every"	"Great... The Great", "Great... the great"

## Recap

## A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

## References

# The Tokenizer

## Recap

## A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

## References

### More examples of large entries QK/OV circuit

Source Token	Destination Token	Out Token	Example Skip Tri-grams
"indy"	"C", "C", "V", "V", "R", "c"	"indy", "obby", "INDY", "loyd"	"indy... Cindy", "indy... CINDY"
"Pike"	"P", "P", "V", "Sp", "V", "R"	"ike", "ikes", "ishing", "owler"	"Pike... Pike", "Pike... Spikes"
"Ralph"	"R", "R", "P", "P", "V", "r"	"alph", "ALPH", "obby", "erald"	"Ralph... Ralph", "Ralph... RALPH"
"Lloyd"	"L", "L", "P", "P", "R", "C"	"loyd", "alph", "\n ", "acman", ... "atherine"	"Lloyd... Lloyd", "Lloyd... Catherine"
"Pixmap"	"P", "Q", "P", "p", "U"	"ixmap", "Canvas", "Embed", "grade"	"Pixmap... Pixmap", "Pixmap... QCanvas"

## Recap

## A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

## References

## Limited Expressivity Can Create Bugs which Seem Strange from the Outside

Source Token	Destination Token	Out Token	"Correct" Skip Tri-grams	"Bug" Skip Tri-grams
"Pixmap"	"P", "Q", "P", "p", "U"	"ixmap", "Canvas", "Embed", "grade"	"Pixmap... Pixmap", "Pixmap... QCanvas"	"Pixmap... P <b>Canvas</b> "
Source Token	Destination Token	Out Token	"Correct" Skip Tri-grams	"Bug" Skip Tri-grams
"Lloyd"	"L", "L", "P", "P", "R", "C"	"loyd", "alph", "\n ", "acman", ... "atherine"	"Lloyd... Lloyd", "Lloyd... Catherine"	"Lloyd... <b>Cloyd</b> ", "Lloyd... <b>Latherine</b> "
Source Token	Destination Token	Out Token	"Correct" Skip Tri-grams	"Bug" Skip Tri-grams
"keep"	"in", "at", "out", "under", "off"	"bay", ... "mind", ... "wraps"	"keep... in mind", "keep... at bay", "keep... under wraps"	"keep... in <b>bay</b> ", "keep... at <b>wraps</b> ", "keep... under <b>mind</b> "

# Induction Heads

## Recap

## A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

## References

### Induction Head - Example 1

Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
 Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
 Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
 Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
 Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
 Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the  
 Mr and Mrs Dursley, of ... such nonsense. Mr Dursley was the

Present Token

Attention

Logit Effect

### Induction Head - Example 2

the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they  
 the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they  
 the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they  
 the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they  
 the Potters. Mrs ... the Potters arrived ... the Potters had ... keeping the Potters away; they

# Wrapping Up

- How would you conceptualize the residual stream?
- What does attention essentially do?
  - What is the meaning of query, key, and value?
- What does an MLP essentially do?
- What can a One-Layer Attention-Only Transformer do?
- What do we enable with two layers?
- What is an Induction Head?

## Recap

### A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

## References



Recap

**A Mathematical  
Framework**

Preliminaries

One-Layer Attention-Only  
Transformers


References

## References


# References

 Alammar, Jay (2018). *The Illustrated Transformer*.


<https://jalammar.github.io/illustrated-transformer/>.  
[Blog post].

 Elhage, Nelson et al. (2021). “A Mathematical Framework for Transformer Circuits”. In: *Transformer Circuits Thread*.

<https://transformer-circuits.pub/2021/framework/index.html>.

 Raffel, Colin et al. (2020). “Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer”. In: *Journal of Machine Learning Research* 21.140, pp. 1–67. url:

<http://jmlr.org/papers/v21/20-074.html>.

 Tunstall, Lewis, Leandro von Werra, and Thomas Wolf (2022). *Natural Language Processing with Transformers*. Sebastopol: O'Reilly.

## Recap

### A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

## References

# References



Vaswani, Ashish et al. (2017). “Attention is All you Need”. In: *Advances in Neural Information Processing Systems*. Ed. by I. Guyon et al. Vol. 30. Curran Associates, Inc. url: [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).



Voita, Elena (Nov. 2023). *NLP Course For You*. url: [https://lena-voita.github.io/nlp\\_course.html](https://lena-voita.github.io/nlp_course.html).

## Recap

### A Mathematical Framework

Preliminaries

One-Layer Attention-Only Transformers

## References