



EACL 2024

Tutorial on *Transformer-specific Interpretability*, Part 2:

Measures of Context Mixing

Hosein Mohebbi, Jaap Jumelet, Michael Hanna,
Afra Alishahi, and Willem Zuidema



Understanding Society

March 21, 2024

Malta



UNIVERSITEIT VAN AMSTERDAM

Waltz of Transformers

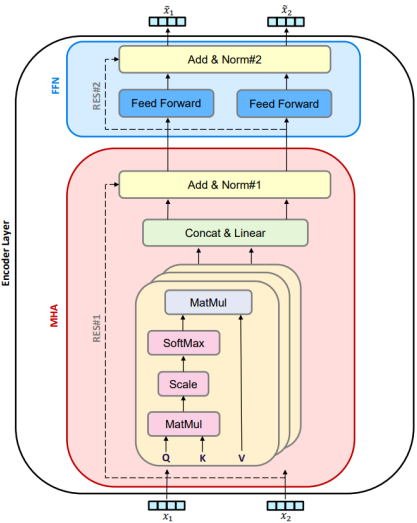
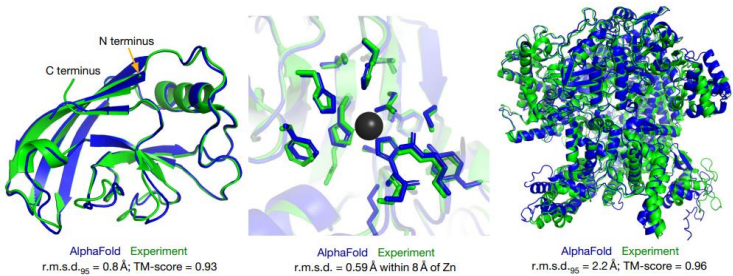
(ChatGPT)



(Seamless)



(AlphaFold)



(Stable Diffusion)



(Sora)



English transcription

Ask not what your country can do for ...

Ask not what your country can do for ...

Any-to-English speech translation

El rápido zorro marrón salta sobre ...

The quick brown fox jumps over ...

Non-English transcription

언덕 위에 올라 내려다보면 너무나 넓고 넓은 ...

언덕 위에 올라 내려다보면 너무나 넓고 넓은

No speech

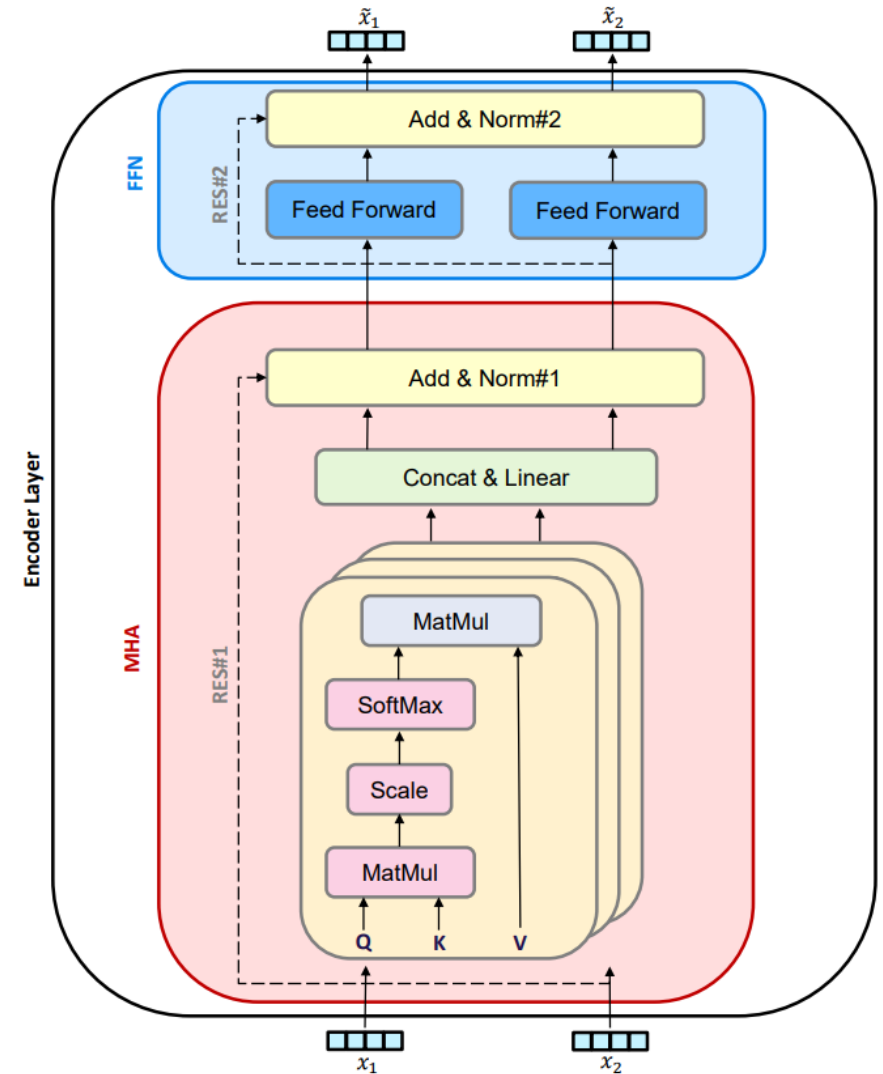
(background music playing)

ø

(Whisper)

Mathematics in Transformer

$$(\mathbf{x}_1, \dots, \mathbf{x}_n)$$

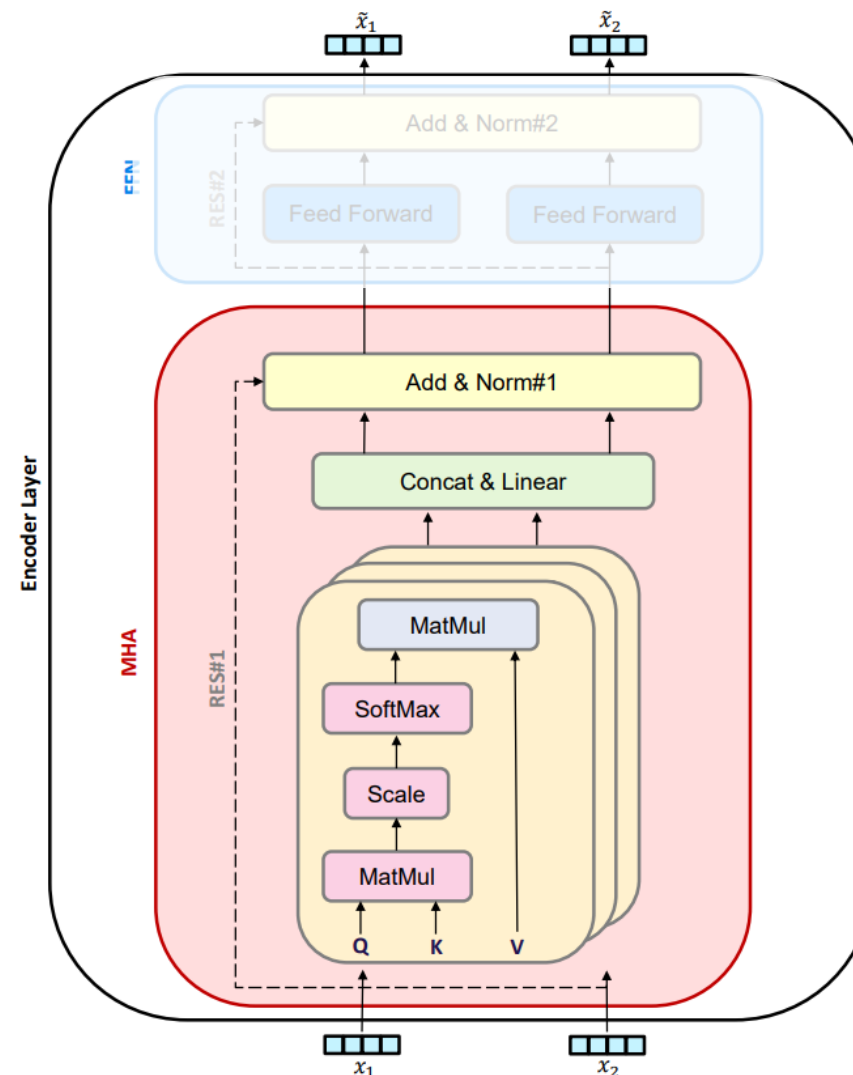


Mathematics in Transformer

$$\begin{aligned}
 & (\mathbf{x}_1, \dots, \mathbf{x}_n) \\
 & \left. \begin{aligned} \mathbf{q}_i^h &= \mathbf{x}_i \mathbf{W}_Q^h + \mathbf{b}_Q^h \\ \mathbf{k}_i^h &= \mathbf{x}_i \mathbf{W}_K^h + \mathbf{b}_K^h \\ \mathbf{v}_i^h &= \mathbf{x}_i \mathbf{W}_V^h + \mathbf{b}_V^h \end{aligned} \right\} \alpha_{i,j} = \operatorname{softmax}_{\mathbf{x}_j \in \mathcal{X}} \left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d}} \right) \in \mathbb{R} \\
 & \mathbf{z}_i^h = \sum_{j=1}^n \alpha_{i,j}^h \mathbf{v}_j^h \\
 & \mathbf{z}_i = \operatorname{CONCAT}(\mathbf{z}_i^1, \dots, \mathbf{z}_i^H) \mathbf{W}_O \\
 & \mathbf{z}_i = \operatorname{LN}_{\text{MHA}}(\mathbf{z}_i + \mathbf{x}_i) \\
 & \tilde{\mathbf{x}}_i = \max(0, \mathbf{z}_i \mathbf{W}_1 + \mathbf{b}_1) \mathbf{W}_2 + \mathbf{b}_2 \\
 & \tilde{\mathbf{x}}_i = \operatorname{LN}_{\text{FFN}}(\tilde{\mathbf{x}}_i + \mathbf{z}_i)
 \end{aligned}$$

MHA

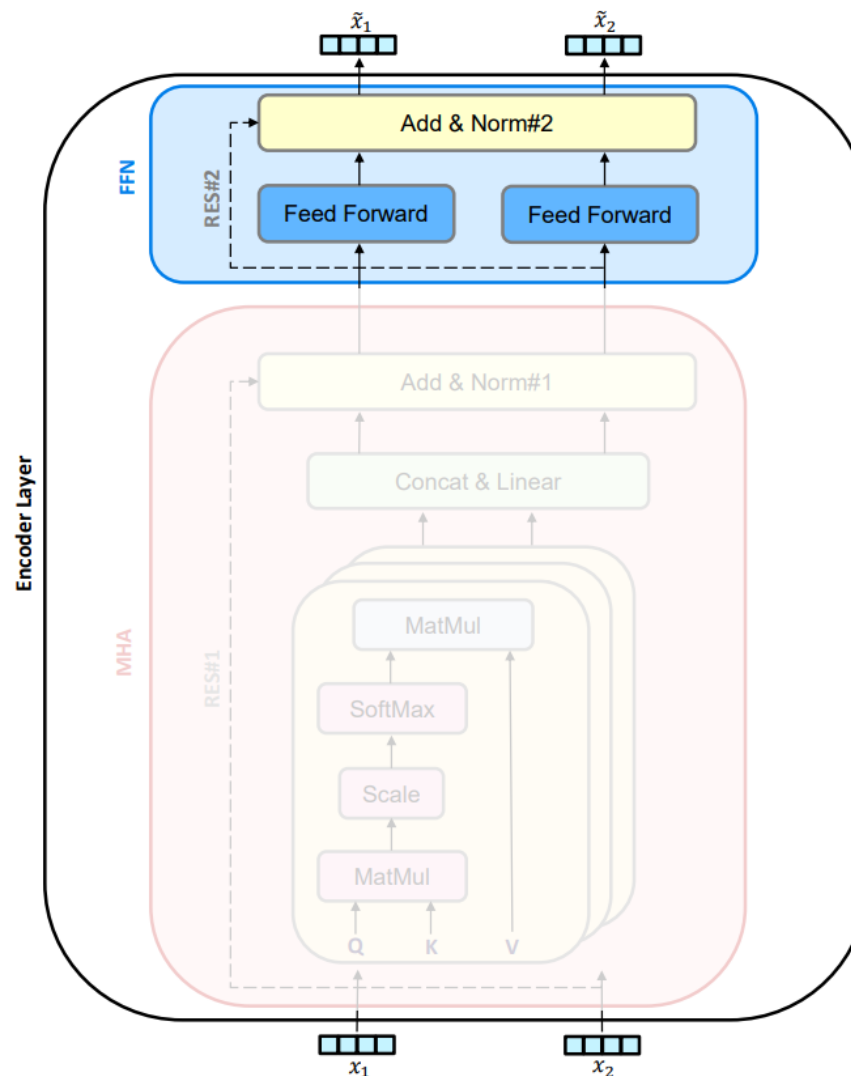
FFN



Mathematics in Transformer

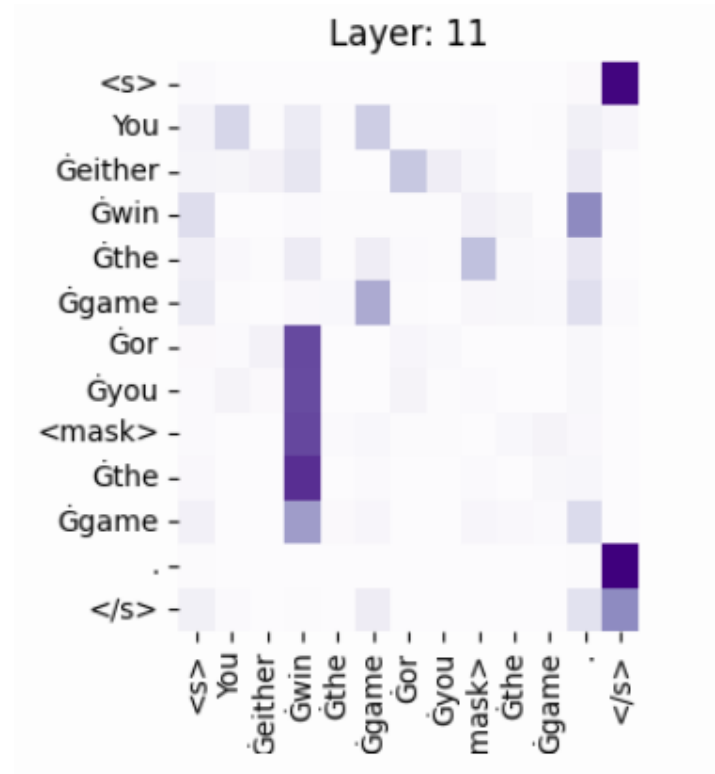
$$\begin{aligned}
 & (\mathbf{x}_1, \dots, \mathbf{x}_n) \\
 & \left. \begin{aligned} \mathbf{q}_i^h &= \mathbf{x}_i \mathbf{W}_Q^h + \mathbf{b}_Q^h \\ \mathbf{k}_i^h &= \mathbf{x}_i \mathbf{W}_K^h + \mathbf{b}_K^h \\ \mathbf{v}_i^h &= \mathbf{x}_i \mathbf{W}_V^h + \mathbf{b}_V^h \end{aligned} \right\} \alpha_{i,j} = \operatorname{softmax}_{\mathbf{x}_j \in \mathcal{X}} \left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d}} \right) \in \mathbb{R} \\
 & \mathbf{z}_i^h = \sum_{j=1}^n \alpha_{i,j}^h \mathbf{v}_j^h \\
 & \mathbf{z}_i = \operatorname{CONCAT}(\mathbf{z}_i^1, \dots, \mathbf{z}_i^H) \mathbf{W}_O \\
 & \mathbf{z}_i = \operatorname{LN}_{\text{MHA}}(\mathbf{z}_i + \mathbf{x}_i) \\
 & \left. \begin{aligned} \tilde{\mathbf{x}}_i &= \max(0, \mathbf{z}_i \mathbf{W}_1 + \mathbf{b}_1) \mathbf{W}_2 + \mathbf{b}_2 \\ \tilde{\mathbf{x}}_i &= \operatorname{LN}_{\text{FFN}}(\tilde{\mathbf{x}}_i + \mathbf{z}_i) \end{aligned} \right\}
 \end{aligned}$$

Encoder Layer

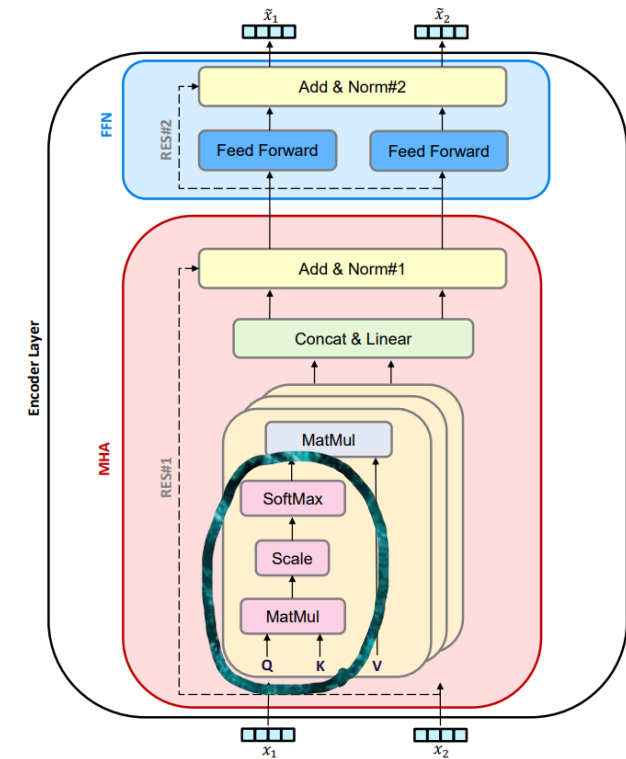
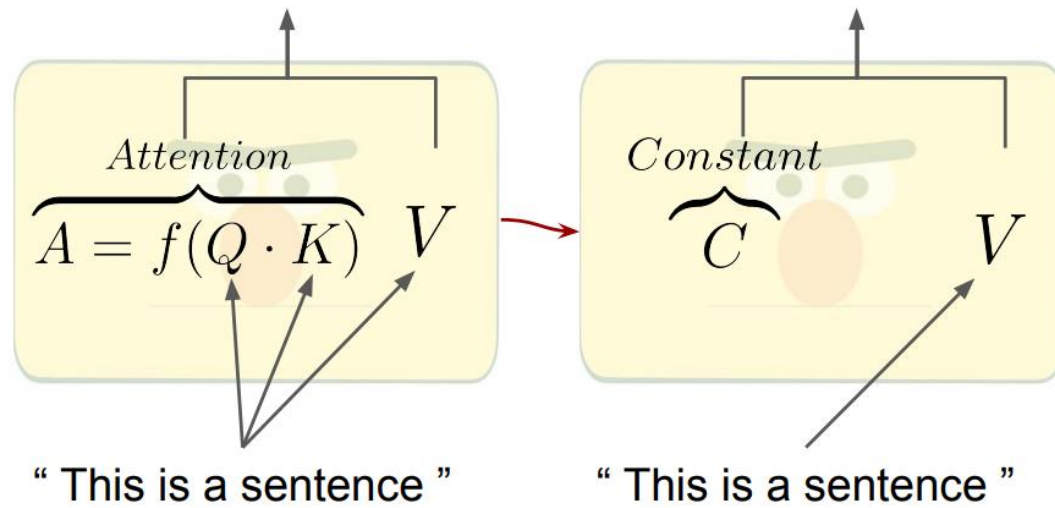


What is Context Mixing?

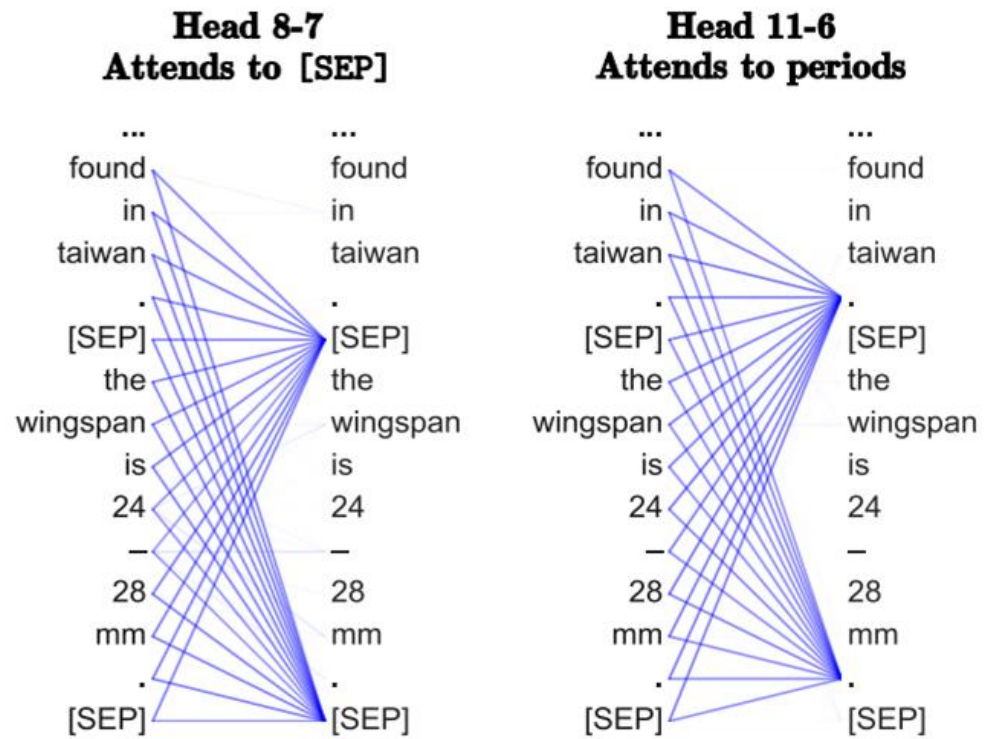
Either you win the game or you <mask> the game.



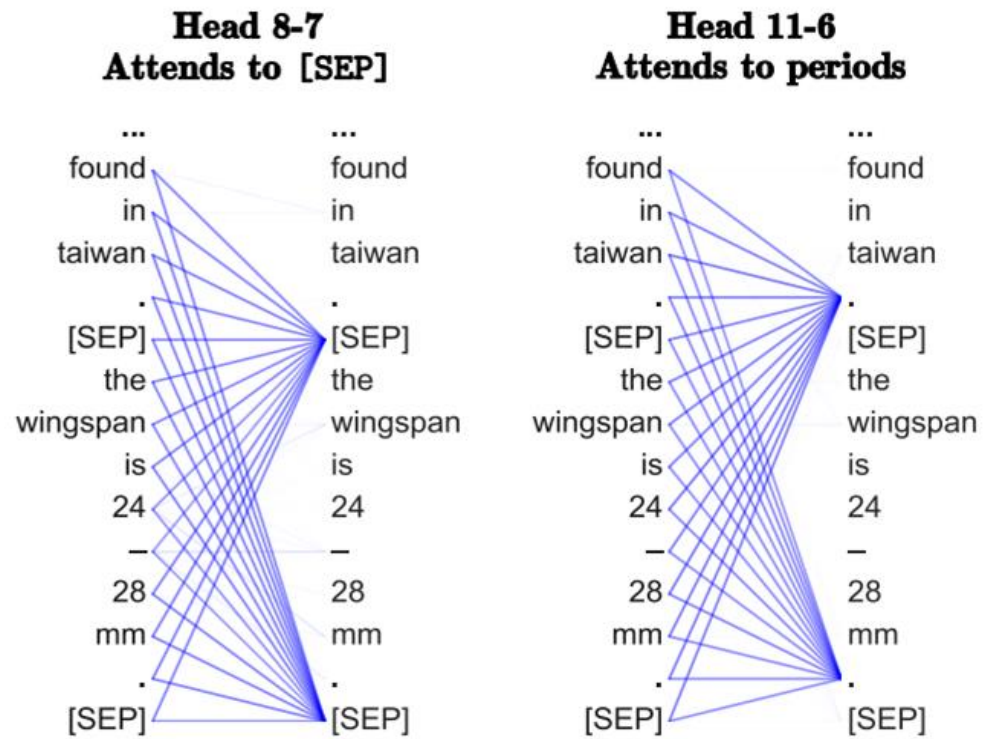
Self-Attention



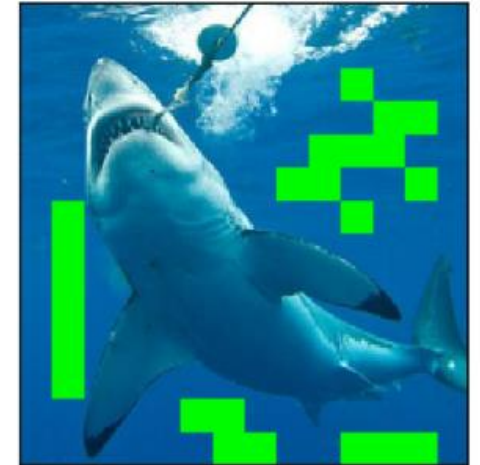
Self-Attention



Self-Attention

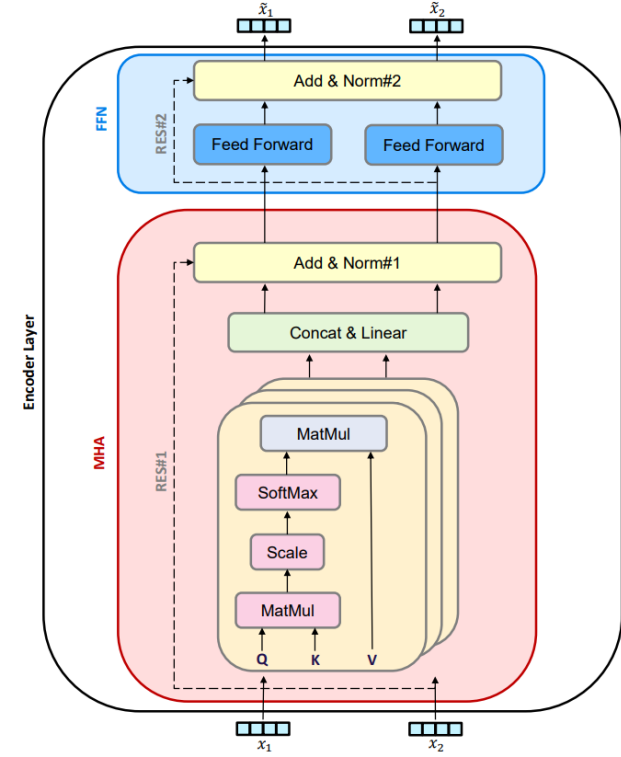
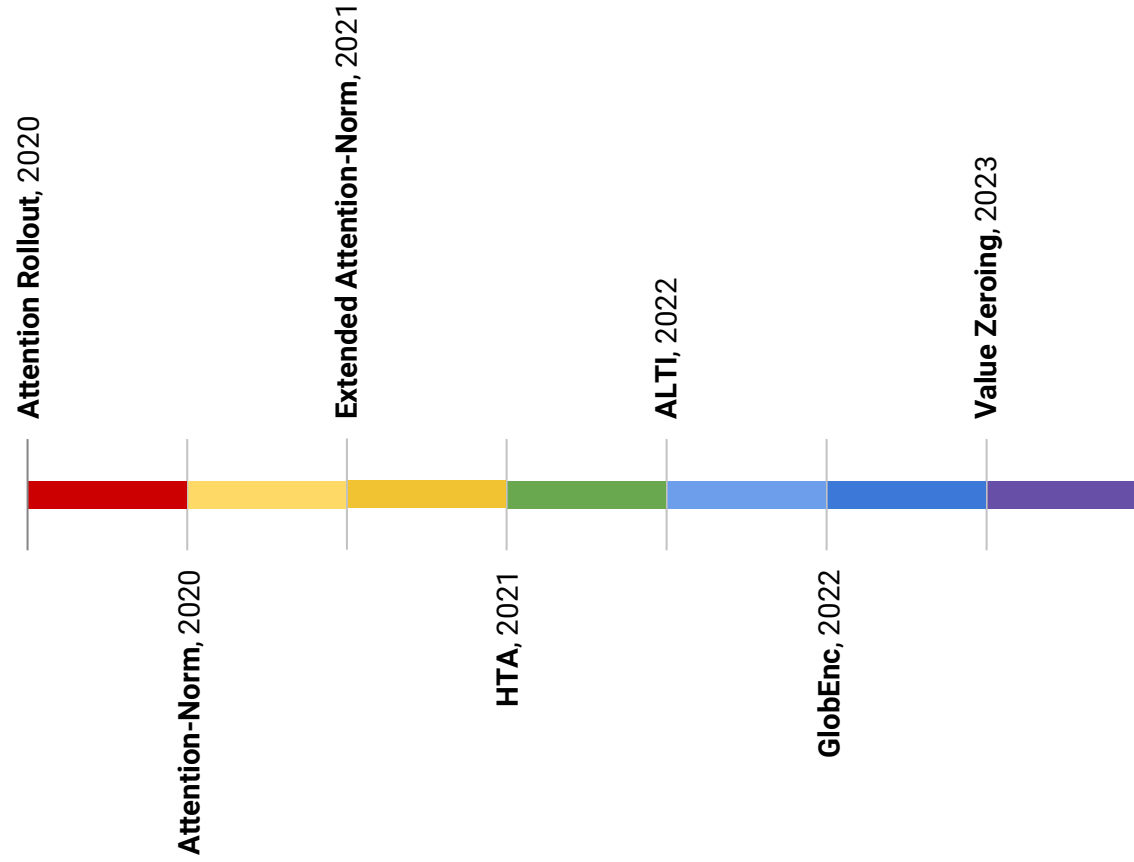


(a)

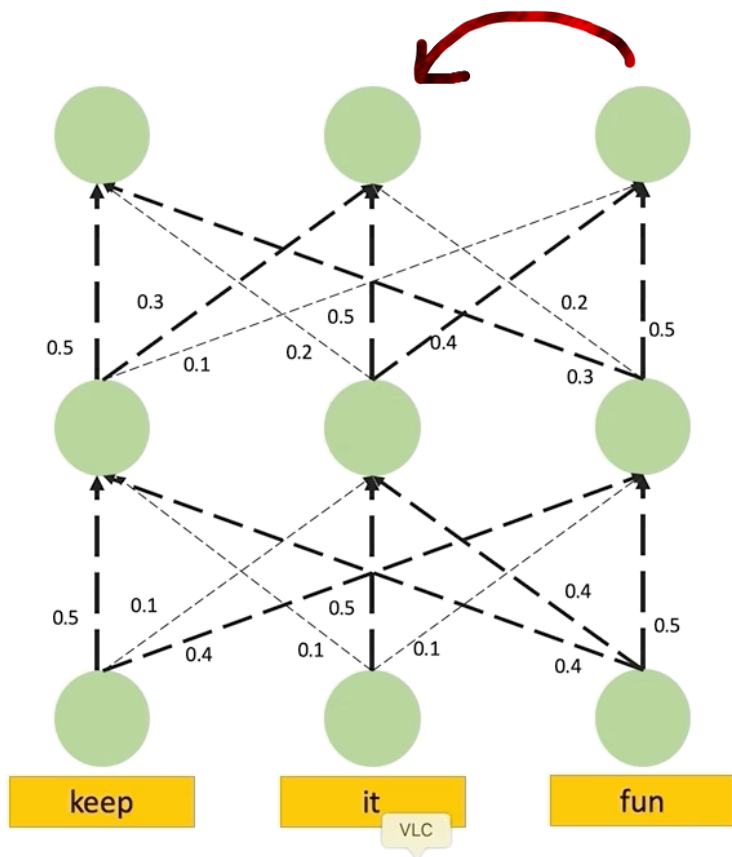


(b)

Measures of Context Mixing

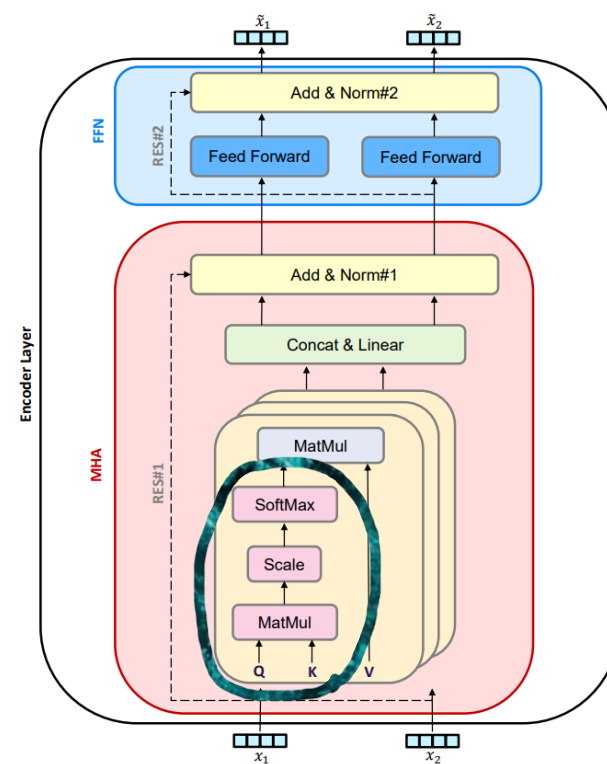


Attention-Rollout

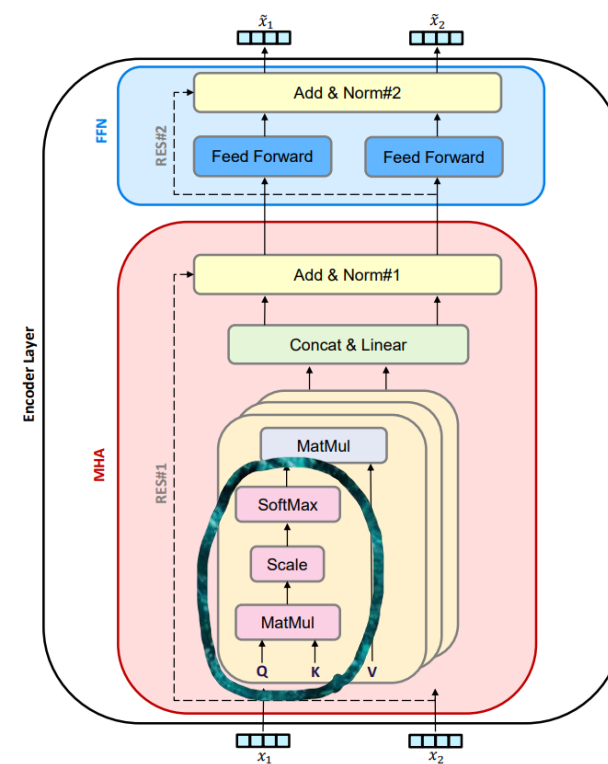
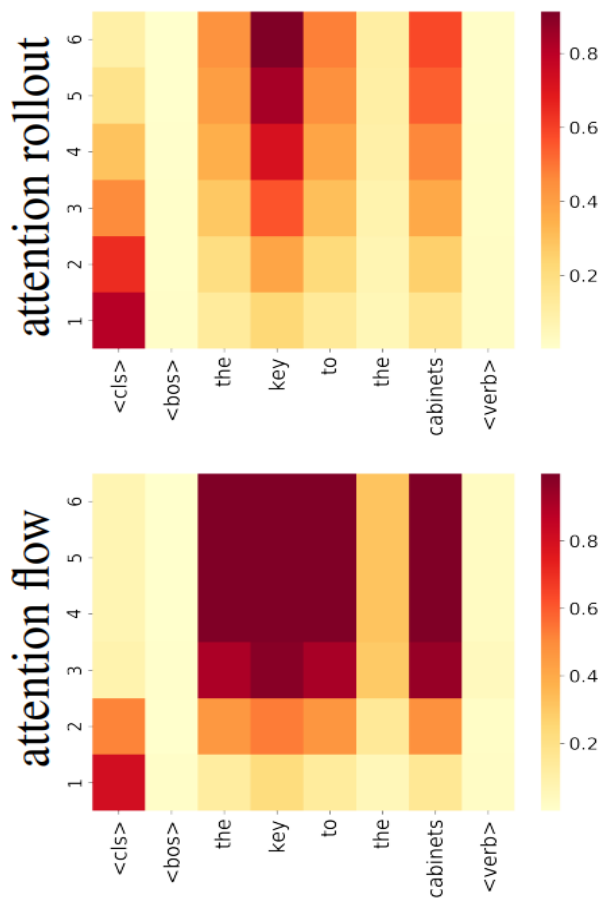
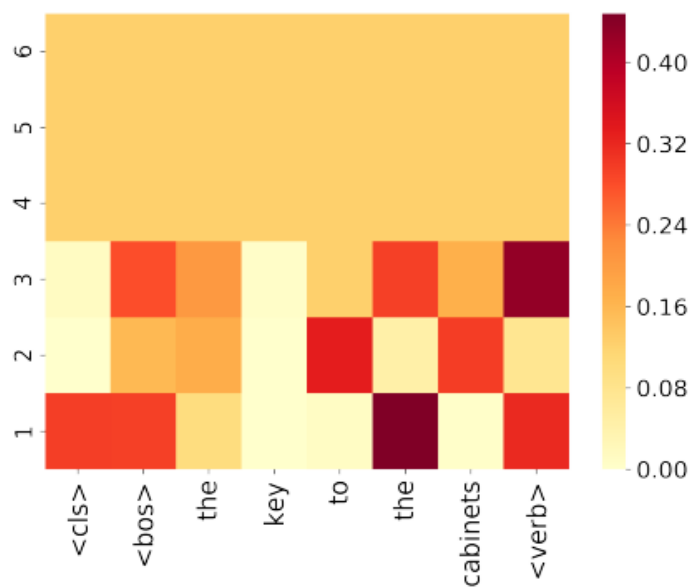


$$\tilde{\mathbf{A}}_{\ell} = \begin{cases} \hat{\mathbf{A}}_{\ell} \tilde{\mathbf{A}}_{\ell-1} & \ell > 1 \\ \hat{\mathbf{A}}_{\ell} & \ell = 1 \end{cases}$$

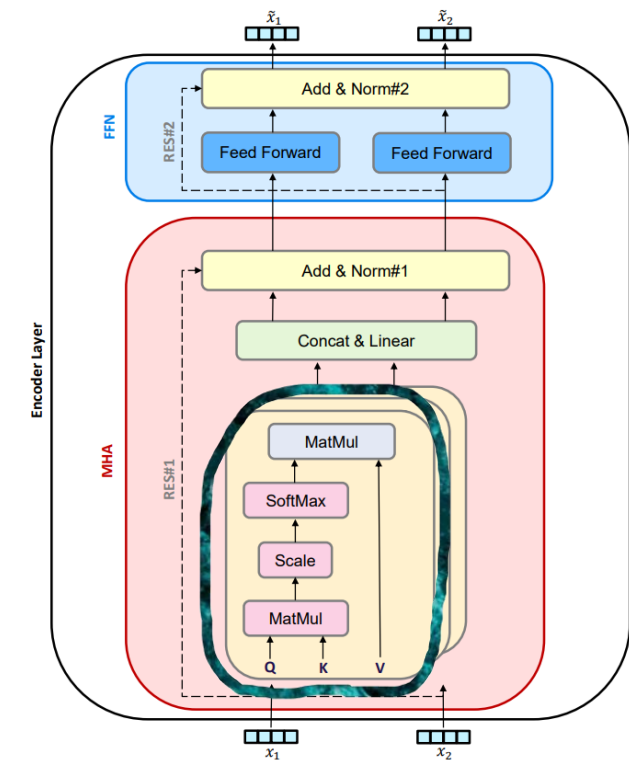
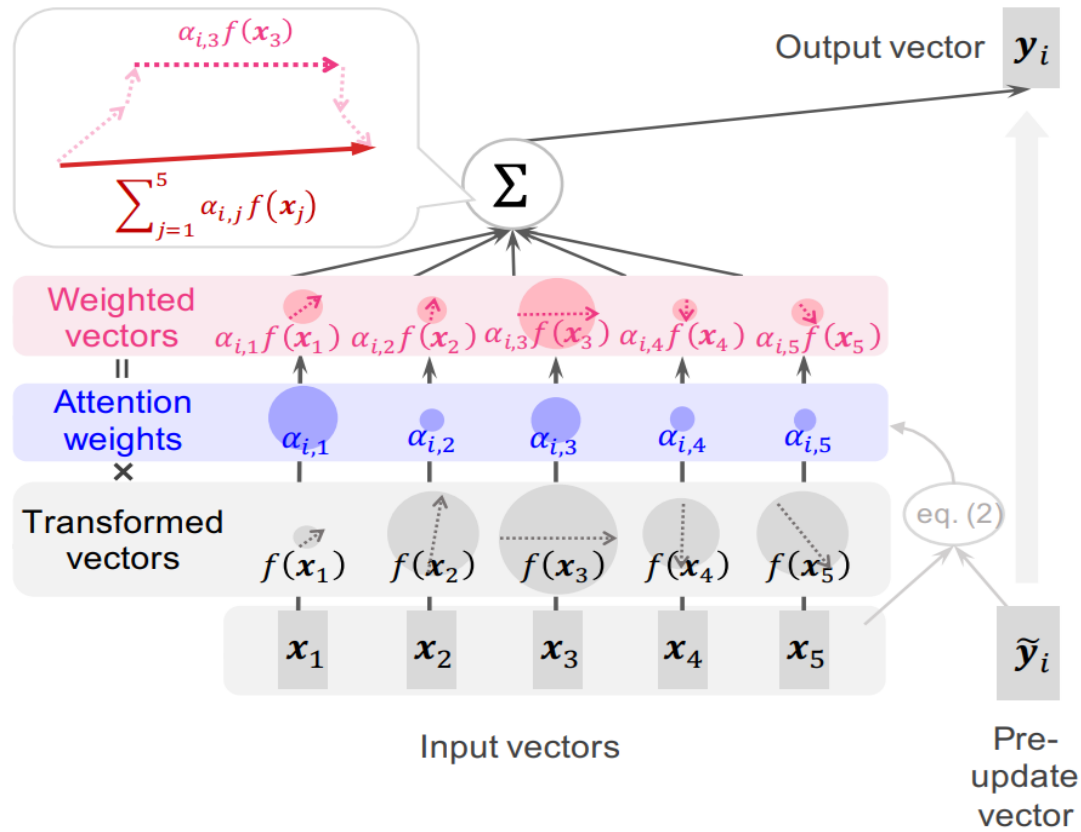
$$\hat{\mathbf{A}}_{\ell} = 0.5 \bar{\mathbf{A}}_{\ell} + 0.5 \mathbf{I}$$



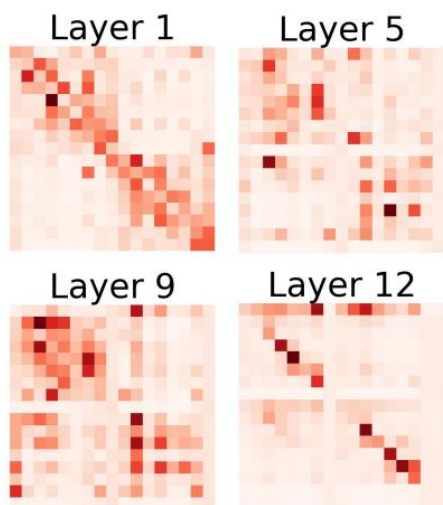
Attention-Rollout



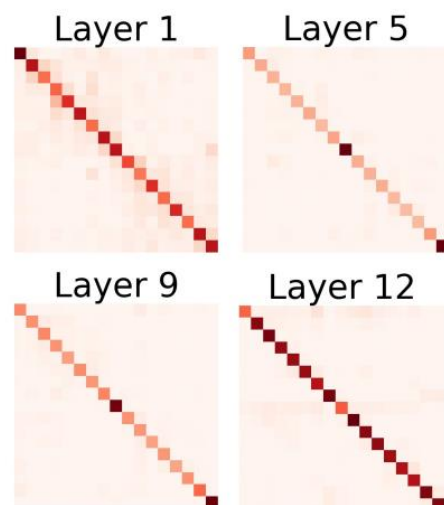
Attention-Norm



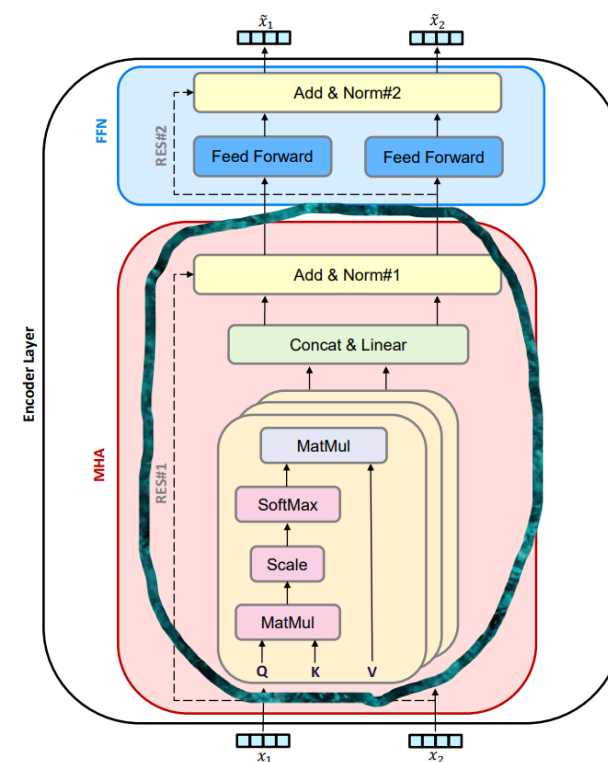
Attention-Norm (extended)



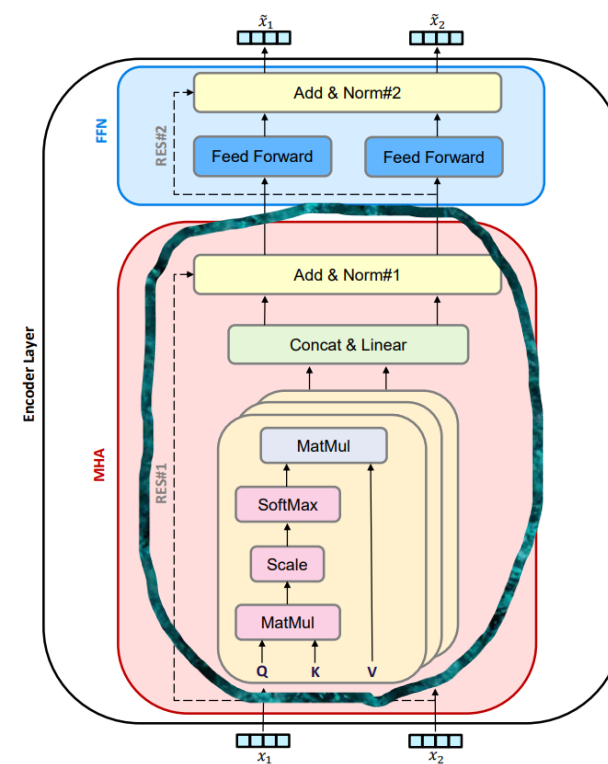
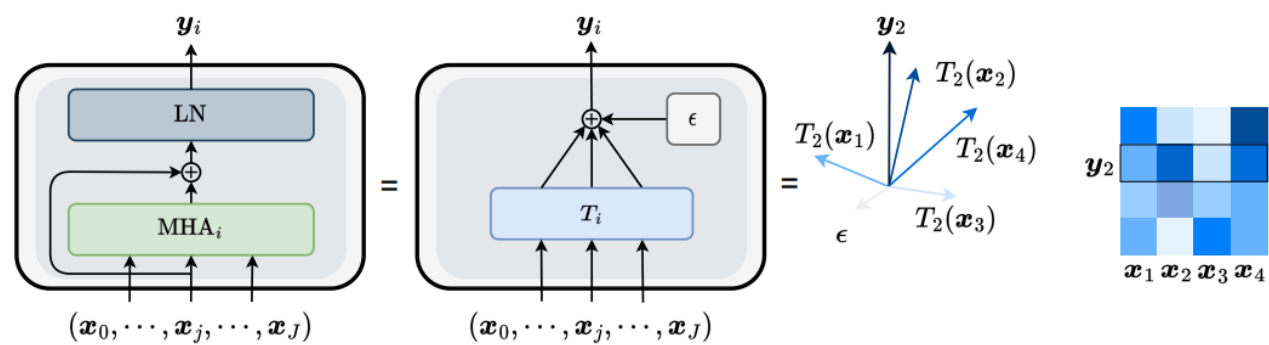
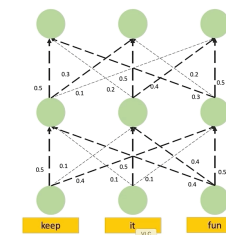
(a) Existing analysis focusing only on the multi-head attention (Kobayashi et al., 2020).



(b) Proposed method incorporating the whole attention block (i.e., multi-head attention, residual connection, and layer normalization) into the analysis.



ALTI



ALTI

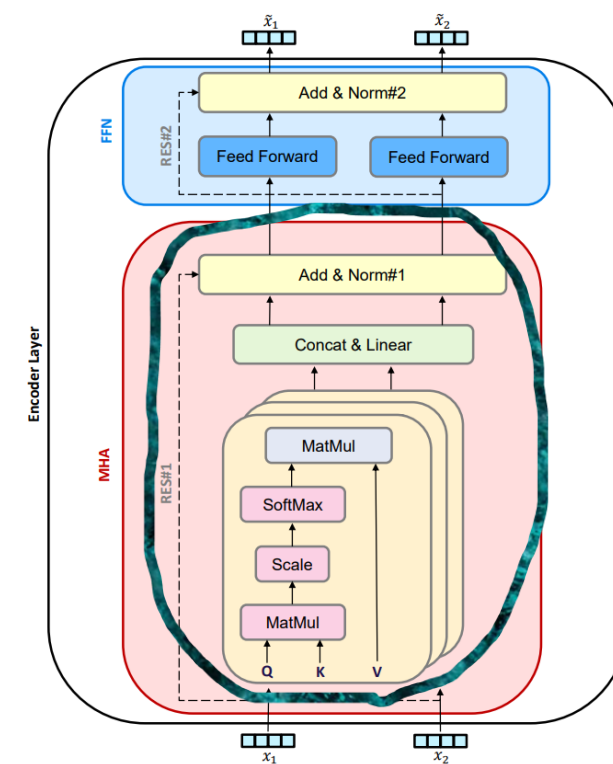
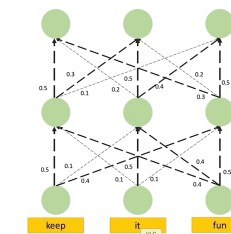
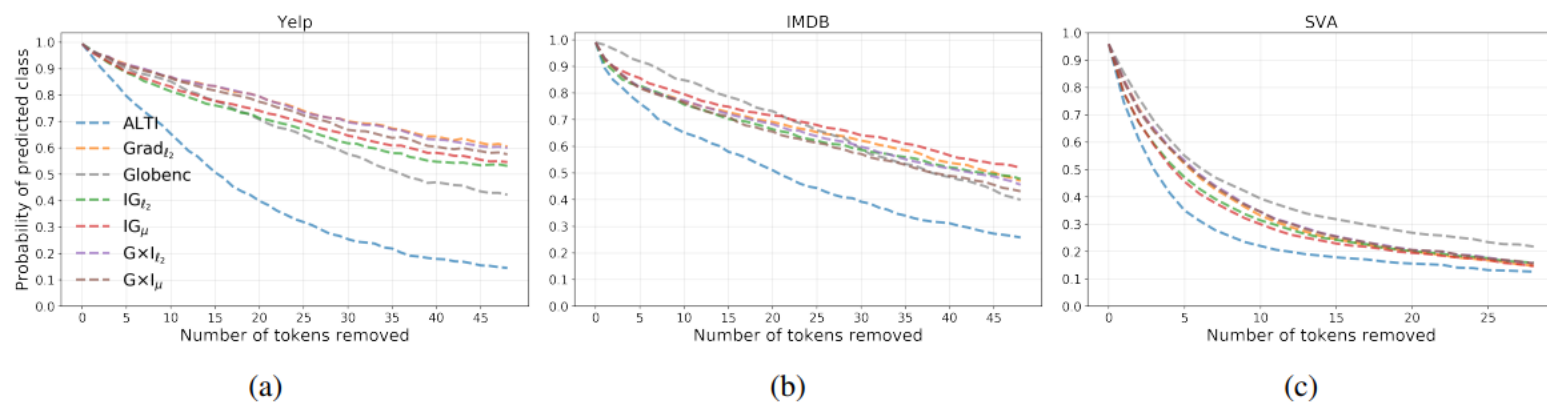


Figure 6: Probability drop in BERT predictions when removing important tokens, obtained by different interpretability methods. We show results on three datasets.

GlobEnc

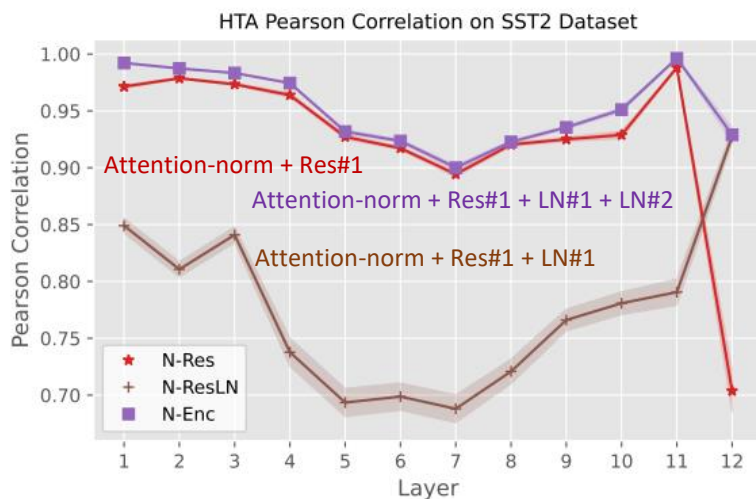


Figure 4: Single layer Pearson correlation of HTA maps with attribution maps. The 99% confidence intervals are shown as shaded areas around each line. $\mathcal{N}_{\text{RESLN}}$ shows considerably less association with HTA.

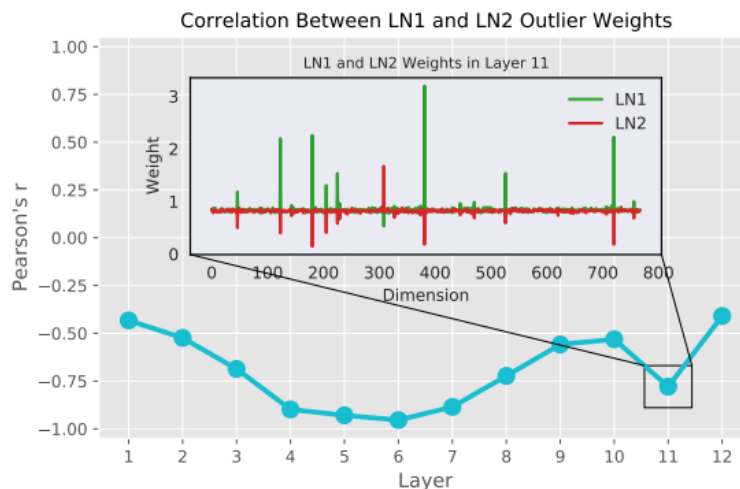
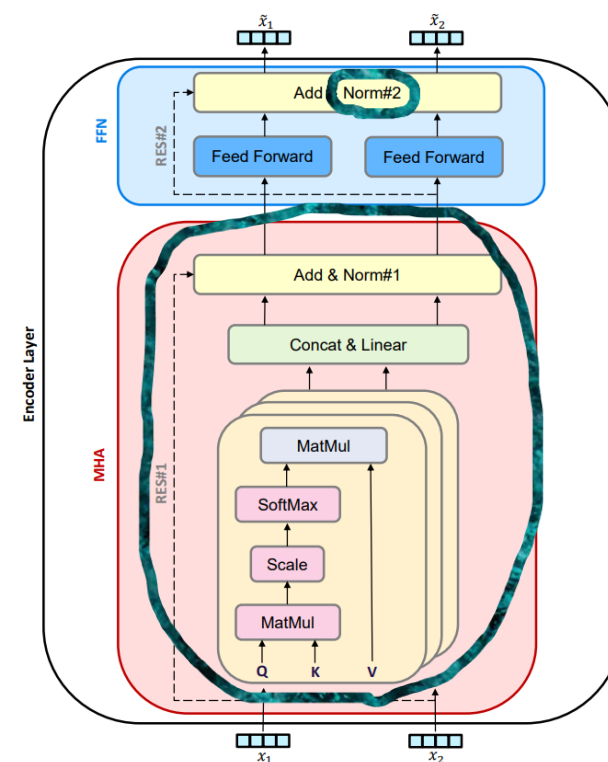
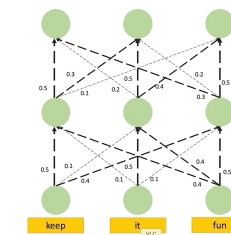
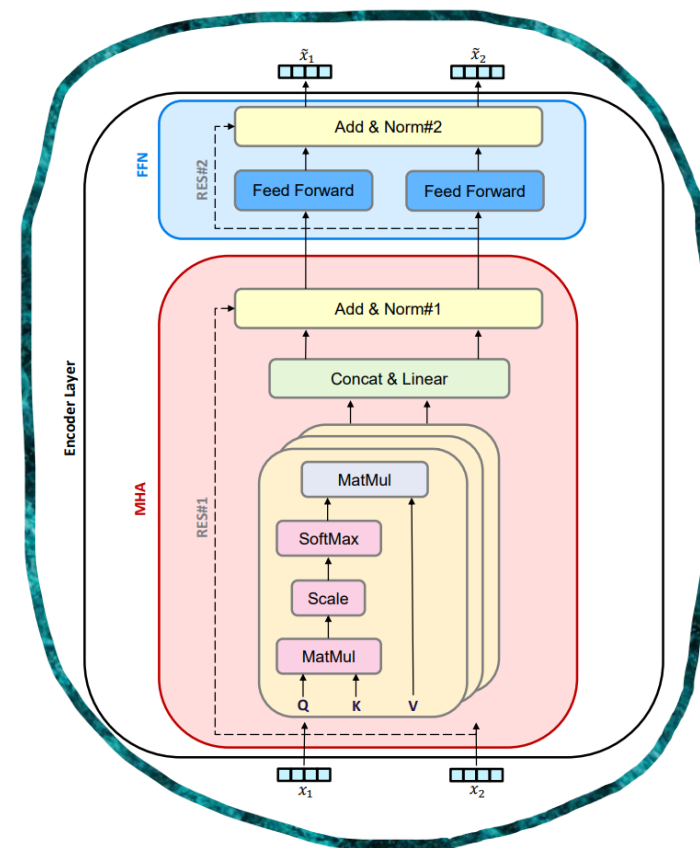
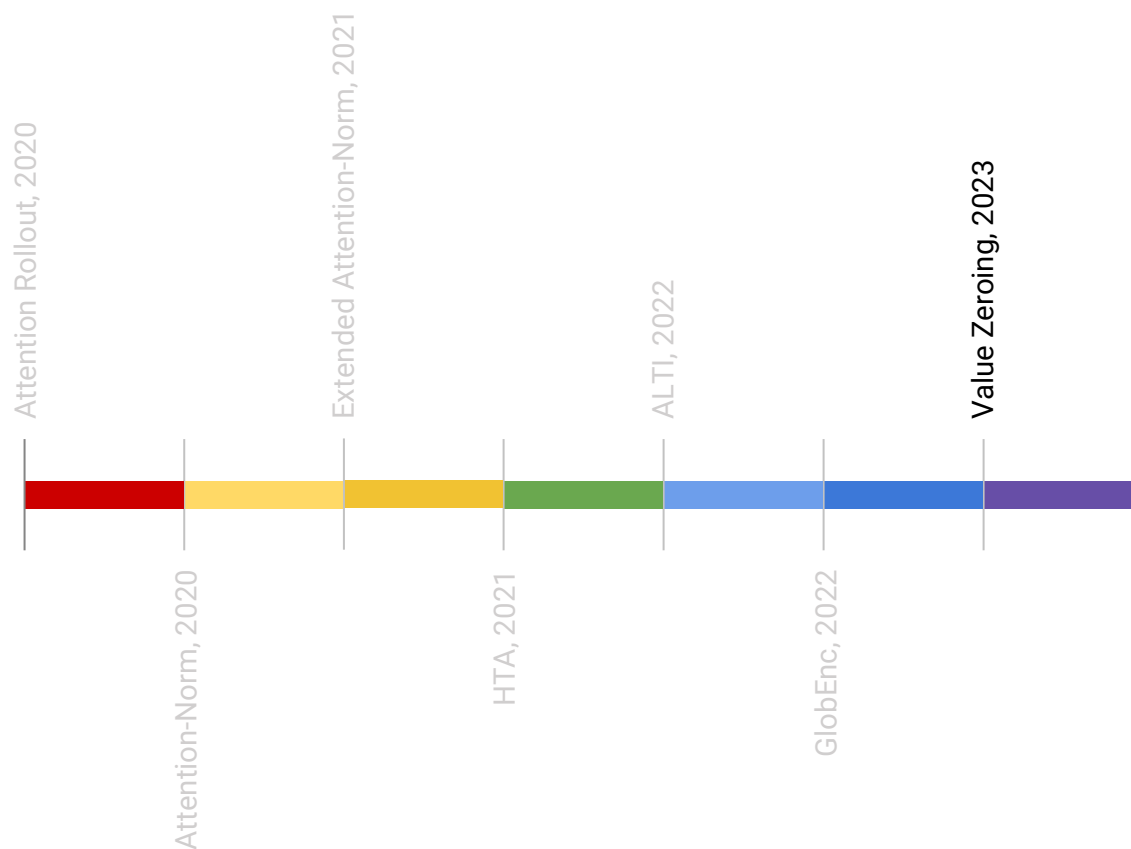


Figure 5: The Pearson correlation between outlier weights of LN#1 and LN#2 across layers. The weight values for layer 11 are shown as well.



Value Zeroing



([Mohebbi et al., 2023](#))

Value Zeroing

$$(\mathbf{x}_1, \dots, \mathbf{x}_n)$$

$$\left. \begin{aligned} \mathbf{q}_i^h &= \mathbf{x}_i \mathbf{W}_Q^h + \mathbf{b}_Q^h \\ \mathbf{k}_i^h &= \mathbf{x}_i \mathbf{W}_K^h + \mathbf{b}_K^h \\ \mathbf{v}_i^h &= \mathbf{x}_i \mathbf{W}_V^h + \mathbf{b}_V^h \end{aligned} \right\} \alpha_{i,j} = \operatorname{softmax}_{\mathbf{x}_j \in \mathcal{X}} \left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d}} \right) \in \mathbb{R}$$

$$\mathbf{z}_i^h = \sum_{j=1}^n \alpha_{i,j}^h \mathbf{v}_j^h$$

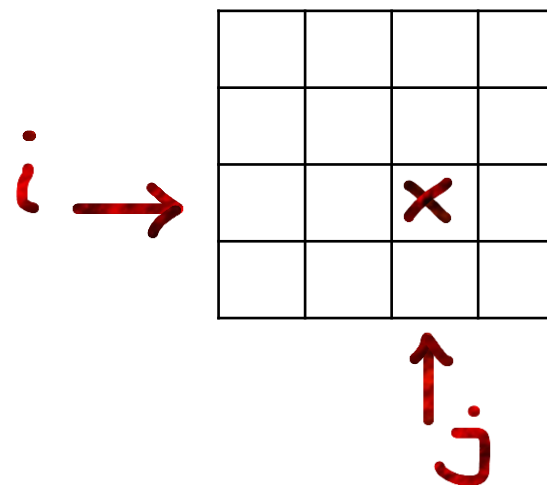
$$\mathbf{z}_i = \operatorname{CONCAT}(\mathbf{z}_i^1, \dots, \mathbf{z}_i^H) \mathbf{W}_O$$

$$\mathbf{z}_i = \operatorname{LN}_{\text{MHA}}(\mathbf{z}_i + \mathbf{x}_i)$$

$$\tilde{\mathbf{x}}_i = \max(0, \mathbf{z}_i \mathbf{W}_1 + \mathbf{b}_1) \mathbf{W}_2 + \mathbf{b}_2$$

$$\tilde{\mathbf{x}}_i = \operatorname{LN}_{\text{FFN}}(\tilde{\mathbf{x}}_i + \mathbf{z}_i)$$

$$C_{i,j} = ?$$



Value Zeroing

$$(\mathbf{x}_1, \dots, \mathbf{x}_n)$$

$$\left. \begin{aligned} \mathbf{q}_i^h &= \mathbf{x}_i \mathbf{W}_Q^h + \mathbf{b}_Q^h \\ \mathbf{k}_i^h &= \mathbf{x}_i \mathbf{W}_K^h + \mathbf{b}_K^h \\ \mathbf{v}_i^h &= \mathbf{x}_i \mathbf{W}_V^h + \mathbf{b}_V^h \end{aligned} \right\} \alpha_{i,j} = \operatorname{softmax}_{\mathbf{x}_j \in \mathcal{X}} \left(\frac{\mathbf{q}_i \mathbf{k}_j^\top}{\sqrt{d}} \right) \in \mathbb{R}$$

$$\mathbf{z}_i^h = \sum_{j=1}^n \alpha_{i,j}^h \mathbf{v}_j^h$$

$$\mathbf{z}_i = \operatorname{CONCAT}(\mathbf{z}_i^1, \dots, \mathbf{z}_i^H) \mathbf{W}_O$$

$$\mathbf{z}_i = \operatorname{LN}_{\text{MHA}}(\mathbf{z}_i + \mathbf{x}_i)$$

$$\tilde{\mathbf{x}}_i = \max(0, \mathbf{z}_i \mathbf{W}_1 + \mathbf{b}_1) \mathbf{W}_2 + \mathbf{b}_2$$

$$\tilde{\mathbf{x}}_i = \operatorname{LN}_{\text{FFN}}(\tilde{\mathbf{x}}_i + \mathbf{z}_i)$$

$$\mathcal{C}_{i,j} = ?$$

$$\mathbf{v}_j^h \leftarrow \mathbf{0}, \forall h \in H$$

$$\mathcal{C}_{i,j} = \tilde{\mathbf{x}}_i^{\neg j} * \tilde{\mathbf{x}}_i$$

Value Zeroing

 (x_1, \dots, x_n) No ablation here!

$$\left. \begin{aligned} q_i^h &= x_i W_Q^h + b_Q^h \\ k_i^h &= x_i W_K^h + b_K^h \\ v_i^h &= x_i W_V^h + b_V^h \end{aligned} \right\} \alpha_{i,j} = \operatorname{softmax}_{x_j \in \mathcal{X}} \left(\frac{q_i k_j^\top}{\sqrt{d}} \right) \in \mathbb{R}$$

$$z_i^h = \sum_{j=1}^n \alpha_{i,j}^h v_j^h$$

$$z_i = \operatorname{CONCAT}(z_i^1, \dots, z_i^H) W_O$$

$$z_i = \operatorname{LN}_{\text{MHA}}(z_i + x_i)$$

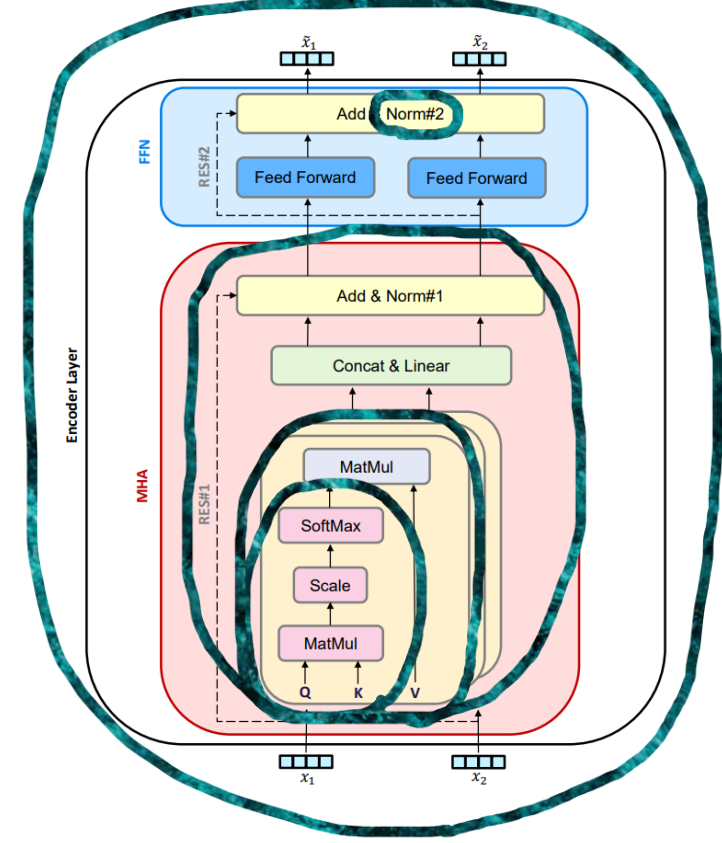
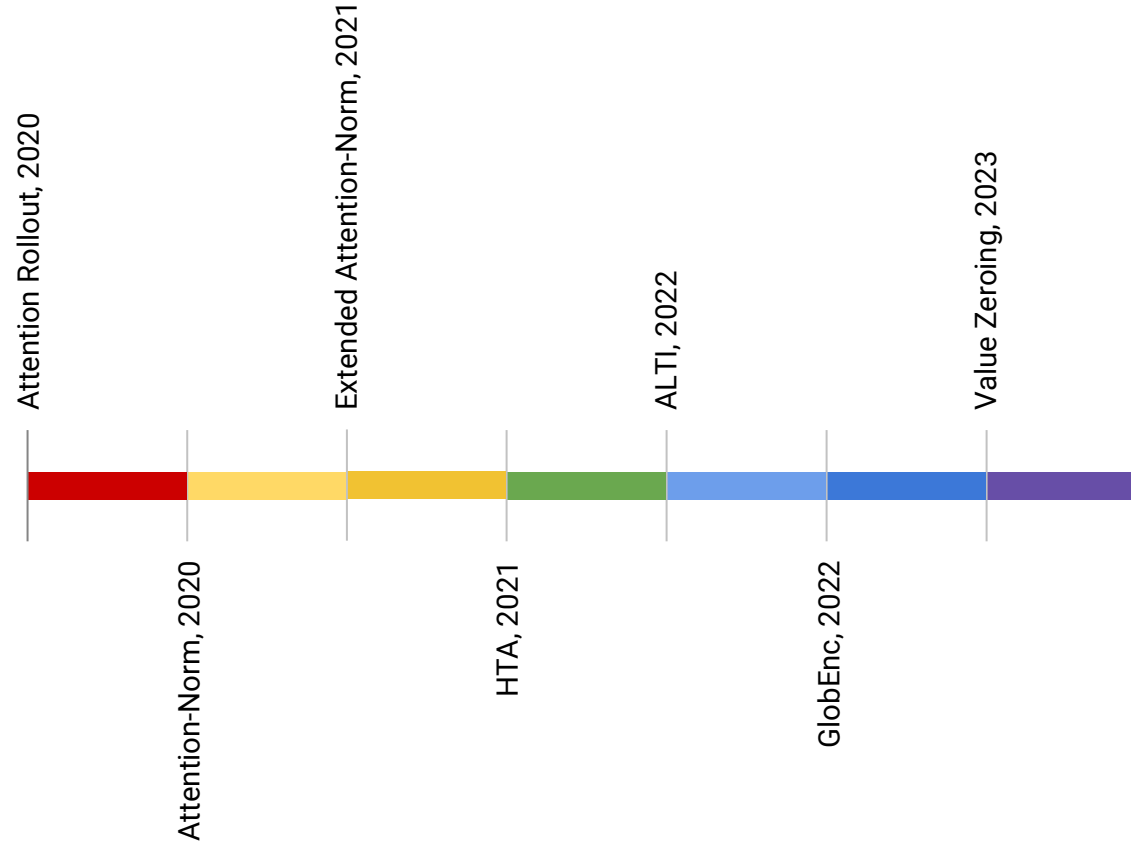
$$\tilde{x}_i = \max(0, z_i W_1 + b_1) W_2 + b_2$$

$$\tilde{x}_i = \operatorname{LN}_{\text{FFN}}(\tilde{x}_i + z_i)$$

$$C_{i,j} = ?$$

$$\begin{aligned} v_j^h &\leftarrow \mathbf{0}, \forall h \in H \\ C_{i,j} &= \tilde{x}_i^{\neg j} * \tilde{x}_i \end{aligned}$$

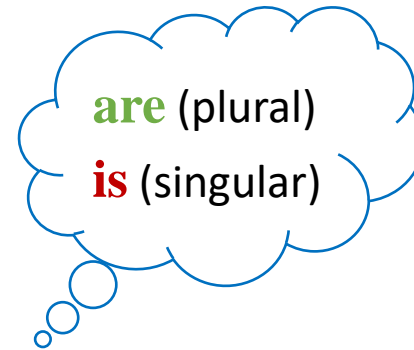
Let's evaluate & compare



Evaluation in Text

Controlled task: grammatical agreements

Target



The books in the library [MASK] read by many.



Evaluation in Text

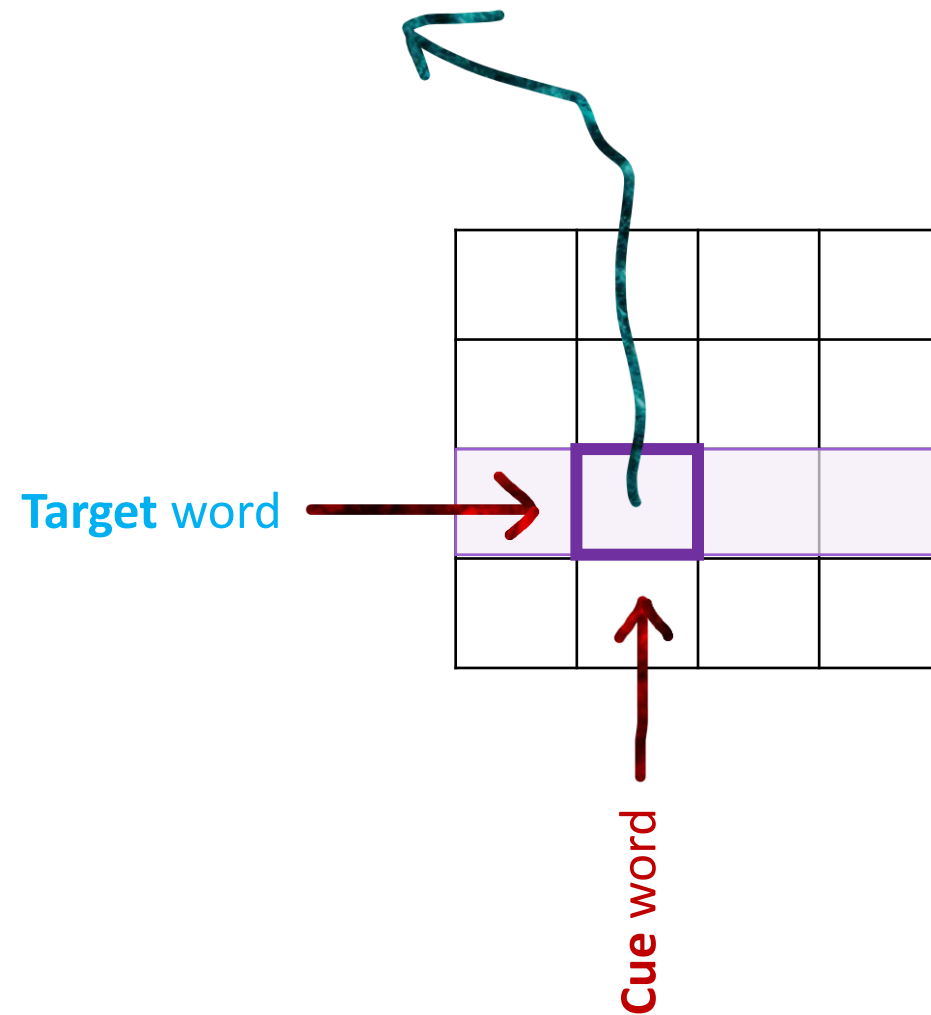
Controlled task: grammatical agreements

Phenomenon	UID	Example	Target word	Foil word
Anaphor Number Agreement	ana	<u>Many teenagers</u> were helping [MASK].	themselves	herself
Determiner-Noun Agreement	dna	Jeffrey has not passed [MASK] <u>museums</u> .	these	this
	dnaa	Sara noticed [MASK] white <u>hospitals</u> .	these	this
Subject-Verb Agreement	darn	The <u>pictures</u> of Martha [MASK] not disgust Anne.	do	does
	rpsv	<u>Kristen</u> [MASK] fixed this chair.	has	have

Table 1: Examples of the selected tasks with our annotations from the BLiMP benchmark (UIDs are unique identifiers used in BLiMP). *Cue* words are underlined.

(Accuracy for BERT-base-uncased is 0.96 in the pre-trained and 0.99 in the fine-tuned setup)

'Cue Contribution' score



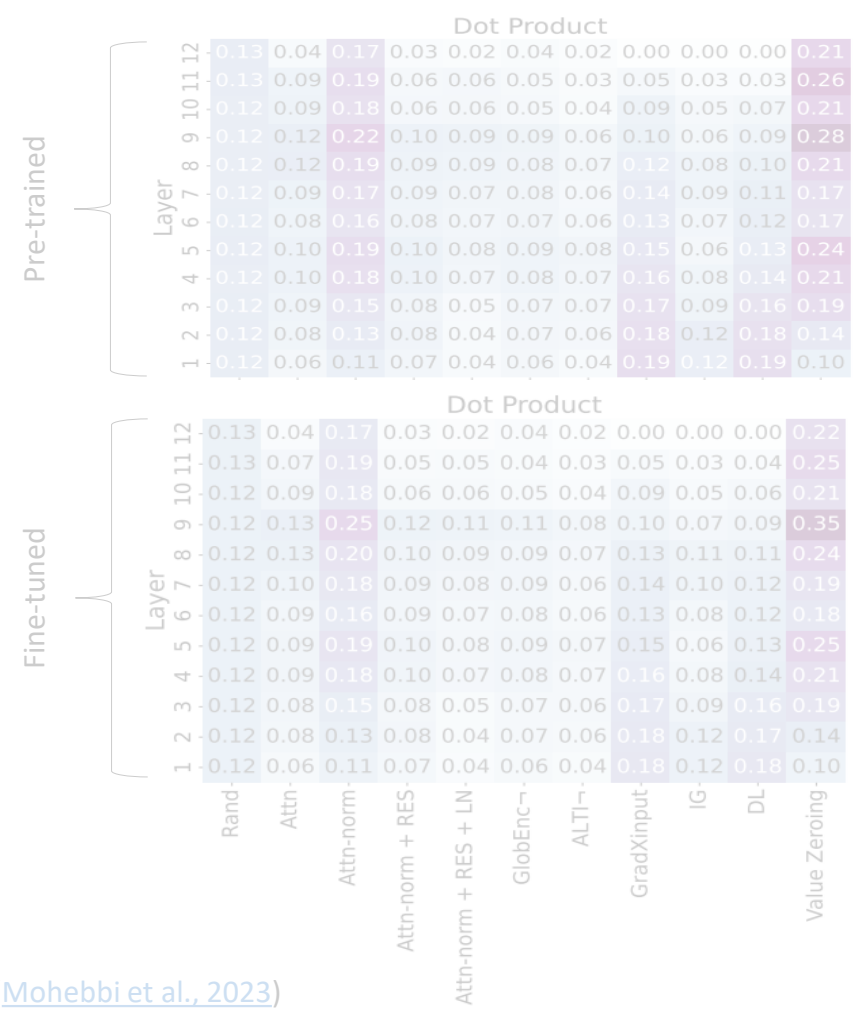
Evaluation in Text

		Dot Product											
Pre-trained	Layer	1	2	3	4	5	6	7	8	9	10	11	12
		0.12	0.08	0.13	0.08	0.04	0.07	0.06	0.18	0.12	0.18	0.14	0.10
		0.12	0.09	0.15	0.08	0.05	0.07	0.07	0.17	0.09	0.16	0.19	0.19
		0.12	0.10	0.18	0.10	0.07	0.08	0.07	0.16	0.08	0.14	0.21	0.21
		0.12	0.10	0.19	0.10	0.08	0.09	0.08	0.15	0.06	0.13	0.24	0.24
		0.12	0.08	0.16	0.08	0.07	0.07	0.06	0.13	0.07	0.12	0.17	0.17
		0.12	0.09	0.17	0.09	0.07	0.08	0.06	0.14	0.09	0.11	0.19	0.19
		0.12	0.12	0.19	0.09	0.09	0.08	0.07	0.12	0.08	0.10	0.28	0.28
		0.12	0.12	0.22	0.10	0.09	0.09	0.06	0.10	0.06	0.09	0.28	0.28
		0.12	0.09	0.18	0.06	0.06	0.05	0.04	0.09	0.05	0.07	0.21	0.21
		0.13	0.09	0.19	0.06	0.06	0.05	0.03	0.05	0.03	0.03	0.26	0.26
		0.13	0.04	0.17	0.03	0.02	0.04	0.02	0.00	0.00	0.00	0.21	0.21
		0.13	0.04	0.17	0.03	0.02	0.04	0.02	0.00	0.00	0.00	0.21	0.21
Fine-tuned	Layer	1	2	3	4	5	6	7	8	9	10	11	12
		0.12	0.08	0.13	0.08	0.04	0.07	0.06	0.18	0.12	0.18	0.14	0.10
		0.12	0.09	0.15	0.08	0.05	0.07	0.07	0.17	0.09	0.16	0.19	0.19
		0.12	0.10	0.18	0.10	0.07	0.08	0.07	0.16	0.08	0.14	0.21	0.21
		0.12	0.09	0.19	0.10	0.08	0.09	0.07	0.15	0.06	0.13	0.25	0.25
		0.12	0.09	0.16	0.09	0.07	0.08	0.06	0.13	0.08	0.12	0.18	0.18
		0.12	0.10	0.18	0.09	0.08	0.09	0.06	0.14	0.10	0.12	0.19	0.19
		0.12	0.13	0.20	0.10	0.09	0.09	0.07	0.13	0.11	0.11	0.24	0.24
		0.12	0.13	0.25	0.12	0.11	0.11	0.08	0.10	0.07	0.09	0.35	0.35
		0.12	0.09	0.18	0.06	0.06	0.05	0.04	0.09	0.05	0.06	0.21	0.21
		0.13	0.07	0.19	0.05	0.05	0.04	0.03	0.05	0.03	0.04	0.25	0.25
		0.13	0.04	0.17	0.03	0.02	0.04	0.02	0.00	0.00	0.00	0.22	0.22
		0.13	0.04	0.17	0.03	0.02	0.04	0.02	0.00	0.00	0.00	0.22	0.22
		0.13	0.04	0.17	0.03	0.02	0.04	0.02	0.00	0.00	0.00	0.22	0.22

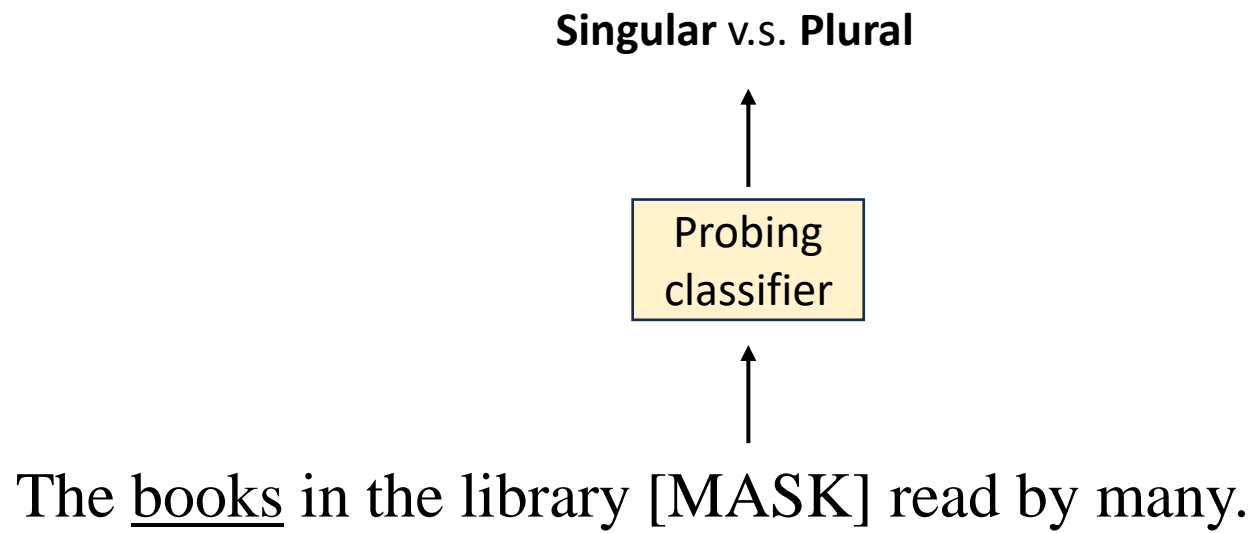
Mohebbi et al., 2023)

(Mohebbi et al., 2023)

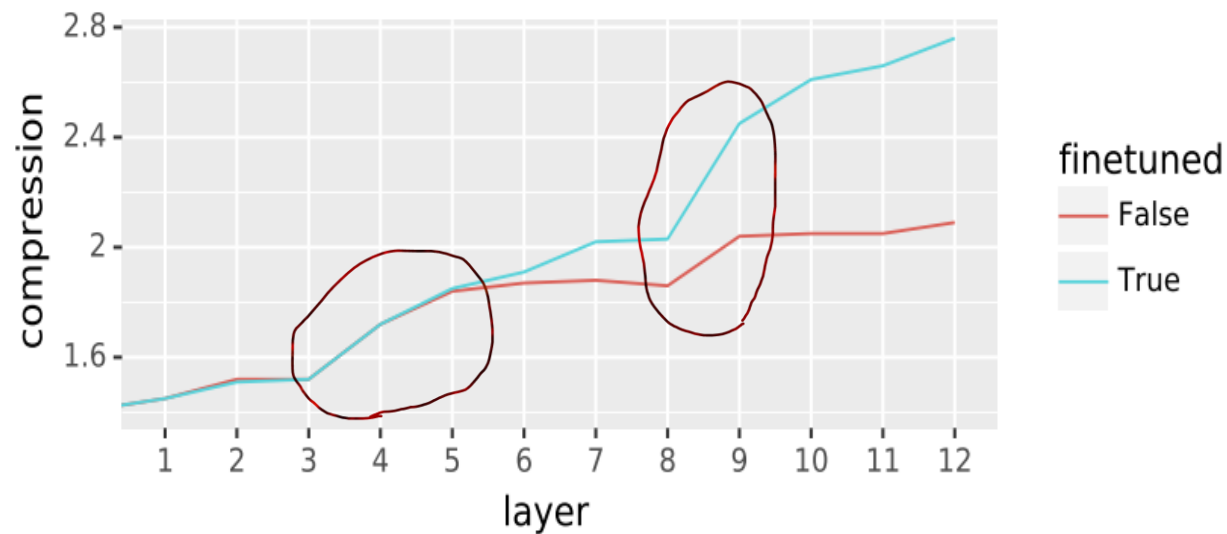
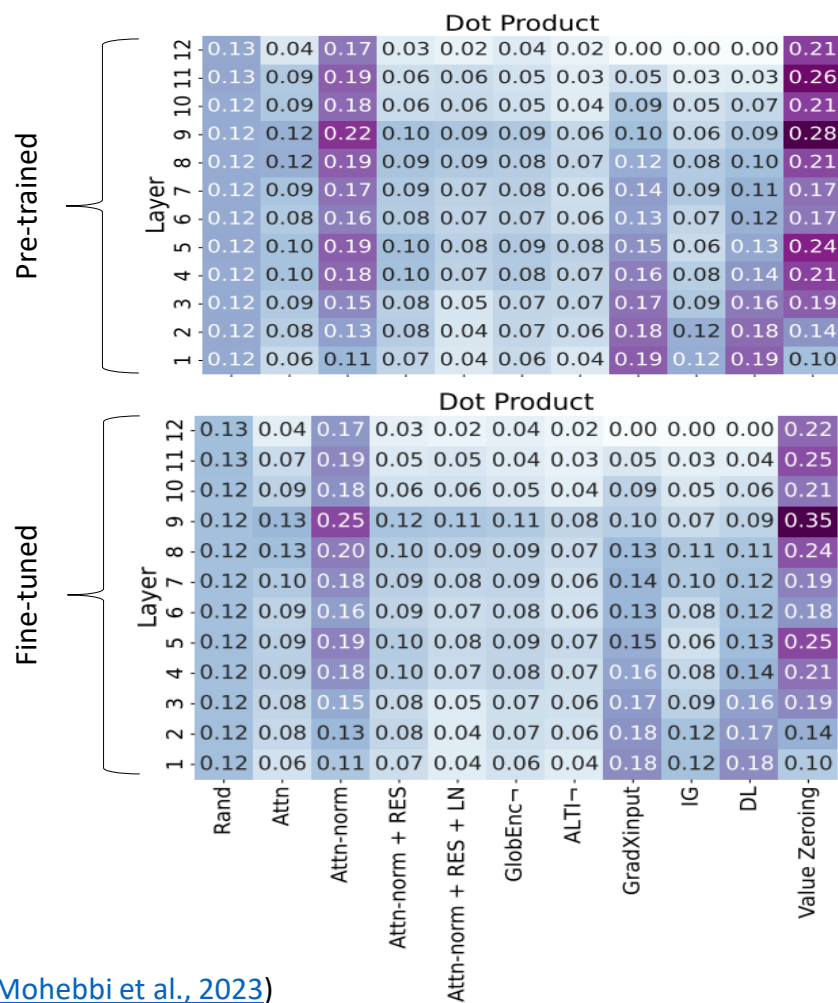
Evaluation in Text



(Mohebbi et al., 2023)



Cue Contribution v.s. Number encoding probing



Aggregated scores for a SVA example

The pictures of some hat [MASK] scaring Marcus.

Aggregated scores for a SVA example

The pictures of some hat [MASK] scaring Marcus.

Attn:

Attn-norm:

Attn-norm+RES:

Attn-norm+RES+LN:

GlobEnc:

ALTI:

GradXinput:

IG:

DL:

Value Zeroing:

Aggregated scores for a SVA example

The pictures of some hat [MASK] scaring Marcus.

Attn:	[CLS] the pictures of some hat [MASK] scar ##ing marcus . [SEP]
Attn-norm:	[CLS] the pictures of some hat [MASK] scar ##ing marcus . [SEP]
Attn-norm+RES:	[CLS] the pictures of some hat [MASK] scar ##ing marcus . [SEP]
Attn-norm+RES+LN:	[CLS] the pictures of some hat [MASK] scar ##ing marcus . [SEP]
GlobEnc:	[CLS] the pictures of some hat [MASK] scar ##ing marcus . [SEP]
ALTI:	[CLS] the pictures of some hat [MASK] scar ##ing marcus . [SEP]
GradXinput:	[CLS] the pictures of some hat [MASK] scar ##ing marcus . [SEP]
IG:	[CLS] the pictures of some hat [MASK] scar ##ing marcus . [SEP]
DL:	[CLS] the pictures of some hat [MASK] scar ##ing marcus . [SEP]
Value Zeroing:	[CLS] the pictures of some hat [MASK] scar ##ing marcus . [SEP]

Evaluation in Speech

Controlled task: homophony in French

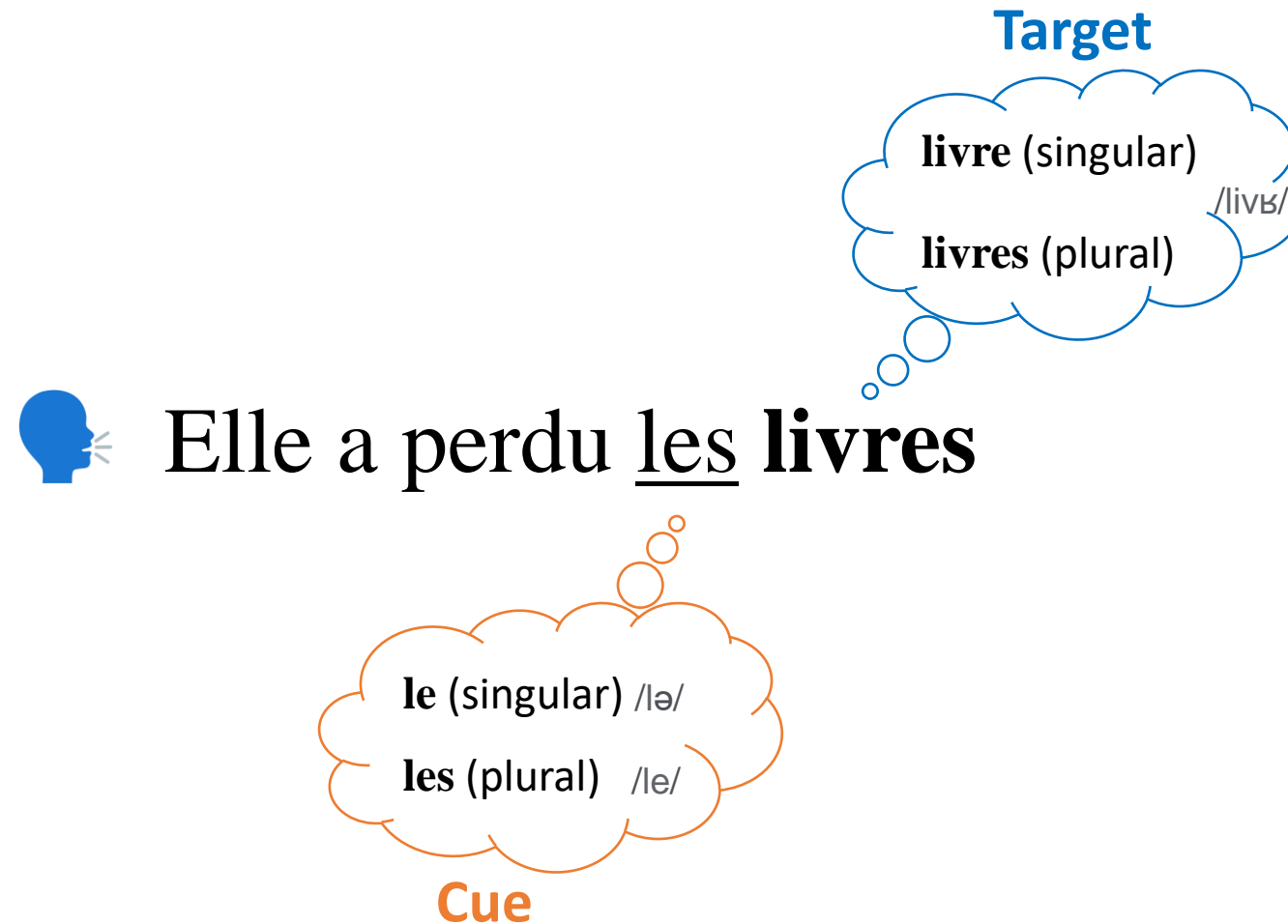


Elle a perdu les livres

(She lost the books)

Evaluation in Speech

Controlled task: homophony in French



Defined Templates

Pattern	Examples of transcription	#
Det_Noun	C'est <u>le</u> septième titre de champion de Syrie de l'histoire du club Il y mène <u>une</u> vie d'études et de recherches	720
Pronoun_Verb	Chaque jour, leurs concurrents les voient sortir de pistes dont <u>ils</u> ignorent l'existence <u>On</u> y trouve une plage naturiste	257
Det_Noun_Verb	Peu après cette élimination, <u>le</u> club et Alexander se séparent à l'amiable À la fin, <u>les</u> enfants se révoltent et détruisent l'école.	23

Table 1: Examples of the extracted audios from the Common Voice corpus based on defined patterns. Last column shows the number of examples obtained. Cue and Target words are underlined and **bolded**, respectively.

Evaluation in Speech

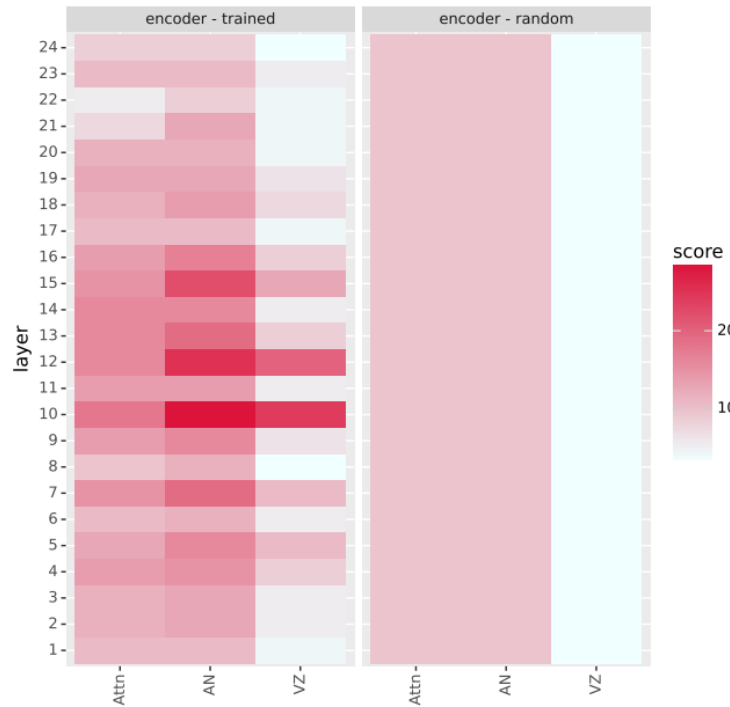


Figure 1: Layer-wise cue contribution according to different analysis methods averaged over all examples for XLSR-53, trained (left) vs. randomly initialized (right).

Evaluation in Speech

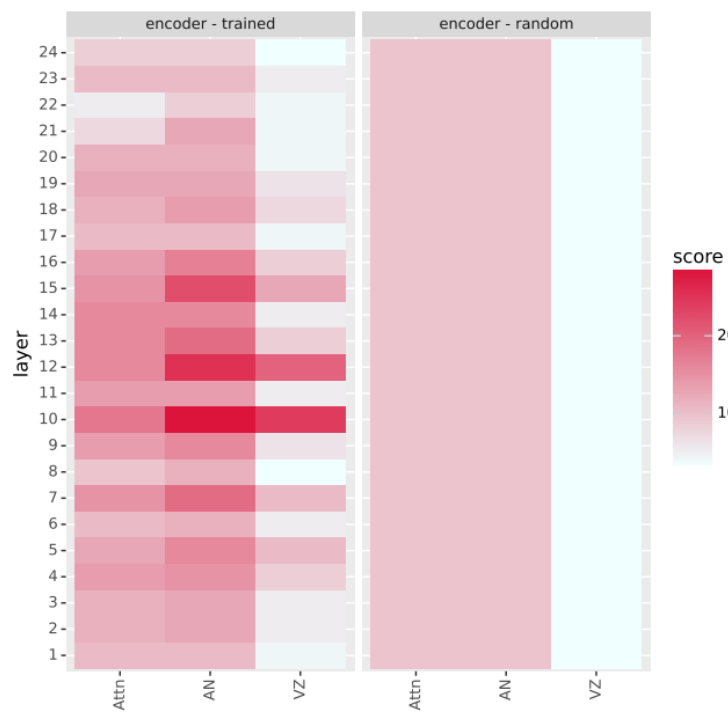


Figure 1: Layer-wise cue contribution according to different analysis methods averaged over all examples for XLSR-53, trained (left) vs. randomly initialized (right).

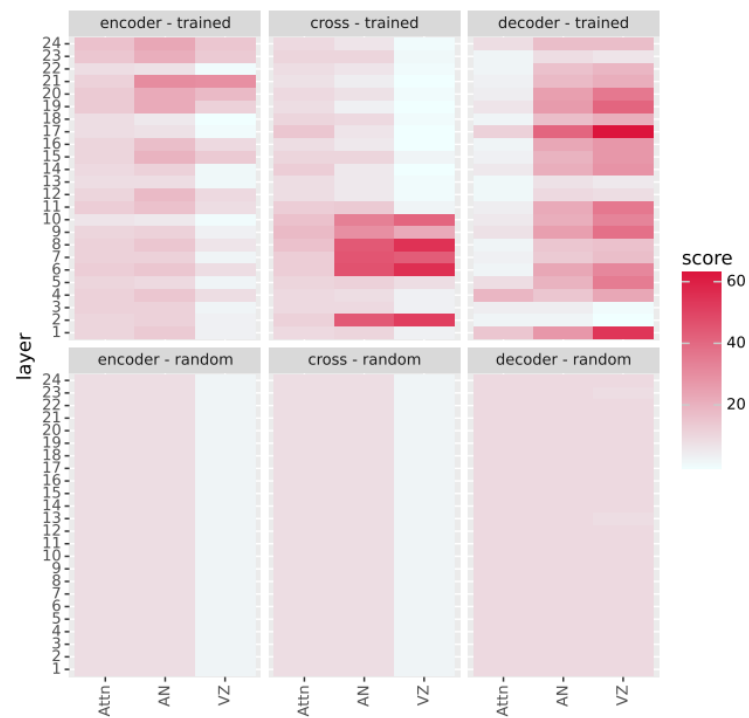


Figure 2: Layer-wise cue contribution according to different analysis methods averaged over all examples for Whisper-medium, trained (top) vs. randomly initialized (bottom).

Cue Contribution v.s. Number encoding probing

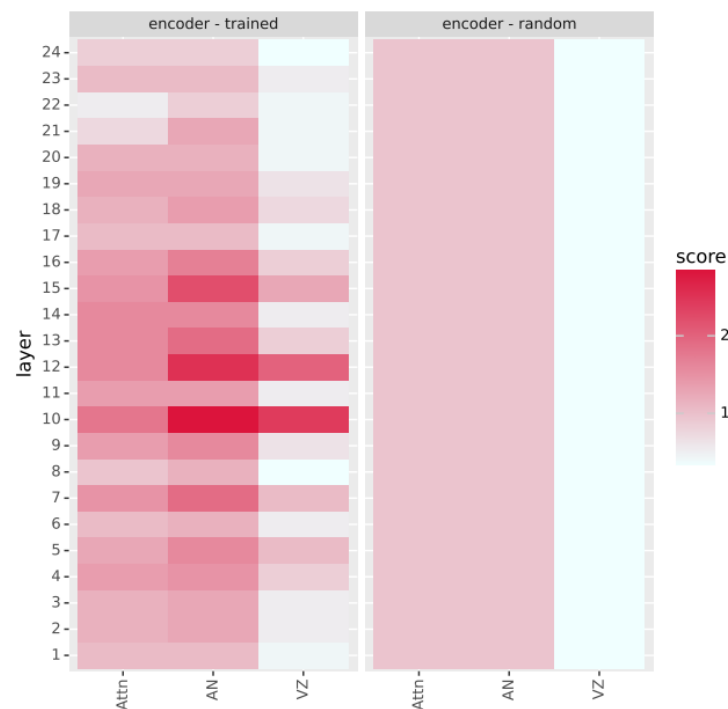


Figure 1: Layer-wise cue contribution according to different analysis methods averaged over all examples for XLSR-53, trained (left) vs. randomly initialized (right).

([Mohebbi et al., 2023](#))

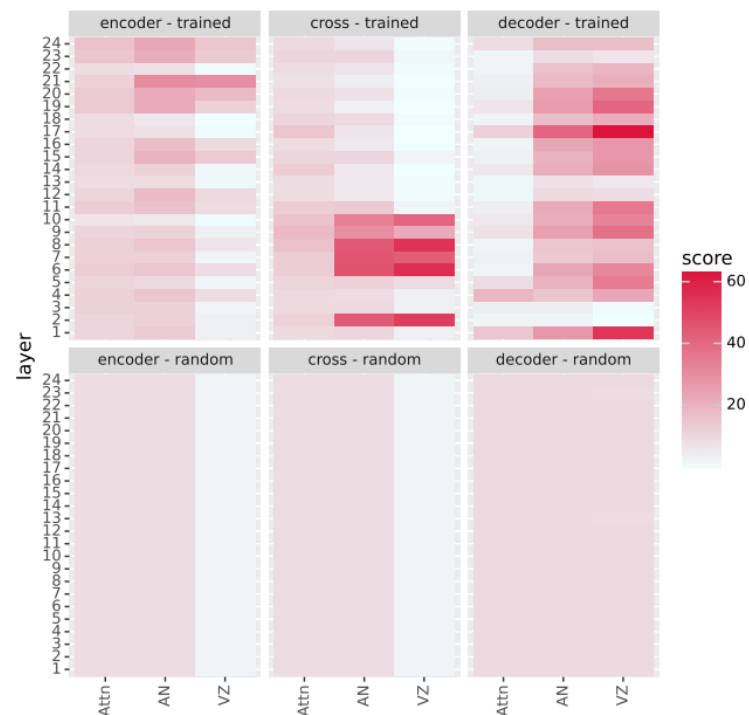


Figure 2: Layer-wise cue contribution according to different analysis methods averaged over all examples for Whisper-medium, trained (top) vs. randomly initialized (bottom).

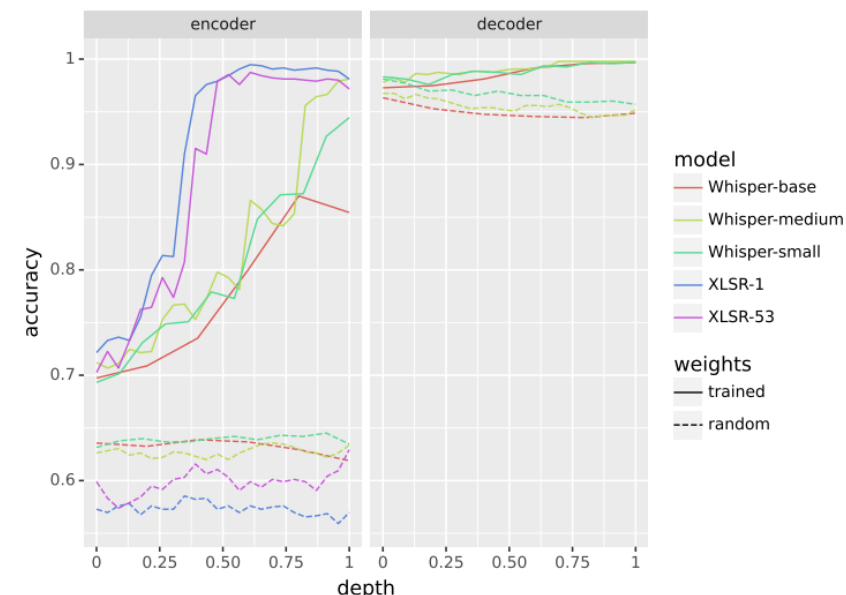


Figure 4: Accuracy of probing classifiers trained on frozen target representations obtained from various ASR models. The depth of Whisper-base (6) and Whisper-small (12), has been normalized to 1 to facilitate comparisons.

Logit-based methods

- LRP-based Attention
- ALTI-Logit
- DecompX

Final thoughts

- Transformers are shared among different modalities; time to converge in our analysis methods to have methods that are:
 - tailored to the **model architecture** (Transformer)
 - irrespective to input **data type** (text, audio, music, image, etc.)
 - irrespective to **training objective** (language modeling, contrastive learning, denoising, etc.)
- We are interested in quantifying context mixing, something that Transformers are made for!
- Attention is **Not** enough for the purpose of context mixing measurement!

Value Zeroing in use



Inseq

Interpretability for Sequence Generation Models 🔍

([Inseq](#))

```
import inseq

model = inseq.load_model("Helsinki-NLP/opus-mt-en-fr", "value_zeroing")
out = model.attribute("A generative language models interpretability tool.")
out.show()
```

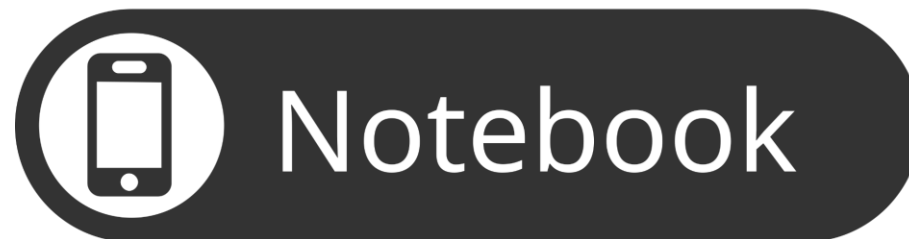
Source Saliency Heatmap

x: Generated tokens, y: Attributed tokens

	__Un	__outil	__d	'	interprétation	__de	__modèles	__de	__langage	__général	atifs	.	</s>
__A	0.784	0.001	0.0	0.0	0.0	0.001	0.0	0.0	0.0	0.0	0.0	0.0	0.0
__	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.001	0.0	0.0	0.0
gener	0.003	0.083	0.349	0.002	0.002	0.237	0.056	0.322	0.01	0.525	0.008	0.0	0.0
ative	0.0	0.0	0.032	0.0	0.0	0.016	0.0	0.036	0.0	0.232	0.945	0.0	0.0
__language	0.003	0.026	0.031	0.044	0.006	0.013	0.026	0.61	0.988	0.001	0.0	0.0	0.0
__models	0.033	0.029	0.027	0.0	0.012	0.72	0.916	0.012	0.0	0.001	0.0	0.001	0.0
__interpret	0.0	0.001	0.464	0.399	0.667	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
ability	0.003	0.001	0.091	0.331	0.31	0.006	0.001	0.0	0.0	0.0	0.0	0.007	0.0
__tool	0.119	0.859	0.001	0.002	0.001	0.001	0.0	0.0	0.0	0.0	0.0	0.002	0.0
.	0.048	0.0	0.002	0.02	0.0	0.003	0.0	0.018	0.0	0.239	0.0	0.979	0.588
</s>	0.006	0.001	0.002	0.201	0.001	0.002	0.001	0.001	0.002	0.001	0.046	0.01	0.412

Thank you! 😊

<https://projects.illc.uva.nl/indeep/tutorial/>

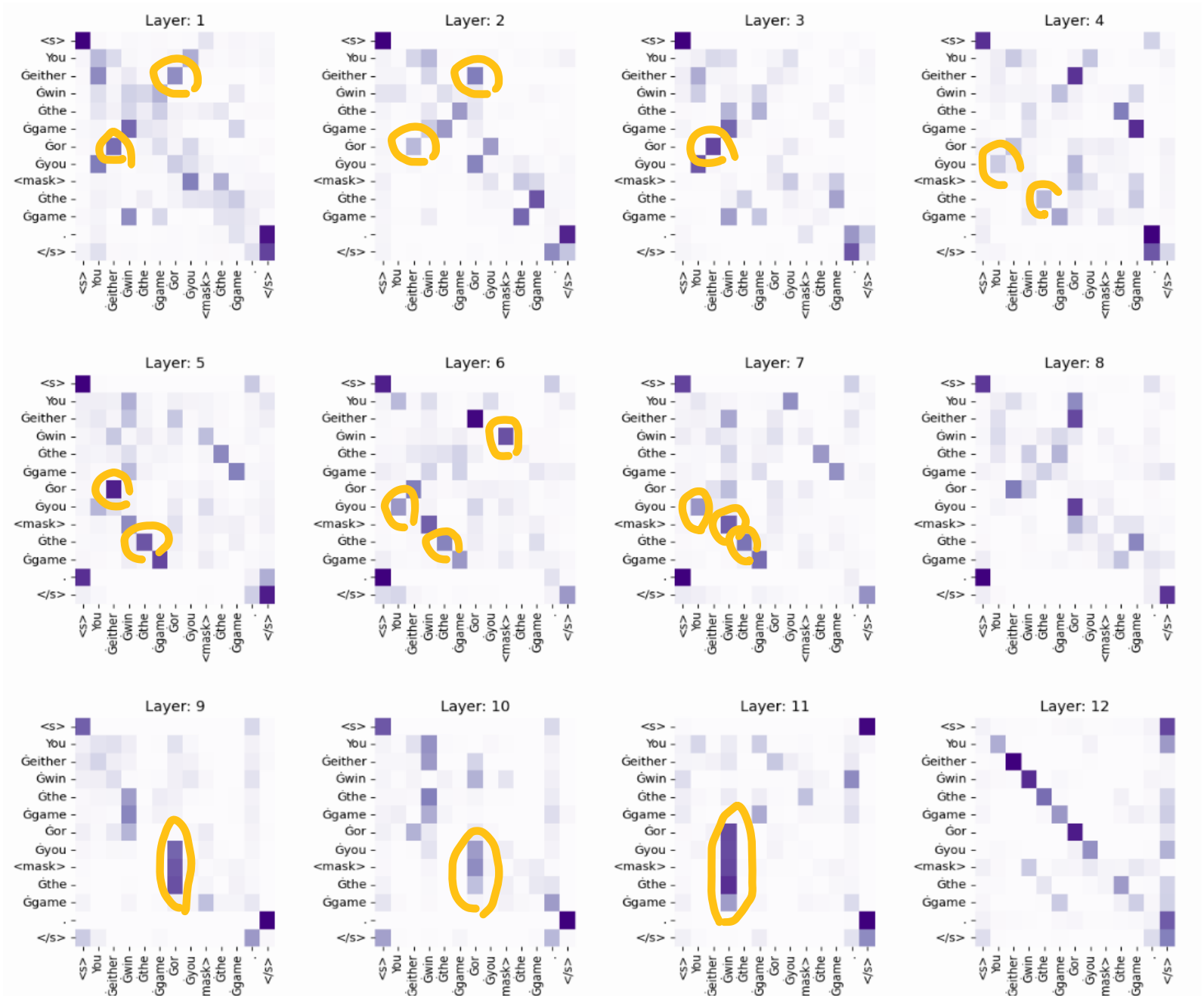


Let's explore!

Either you win the game or you <mask> the game.

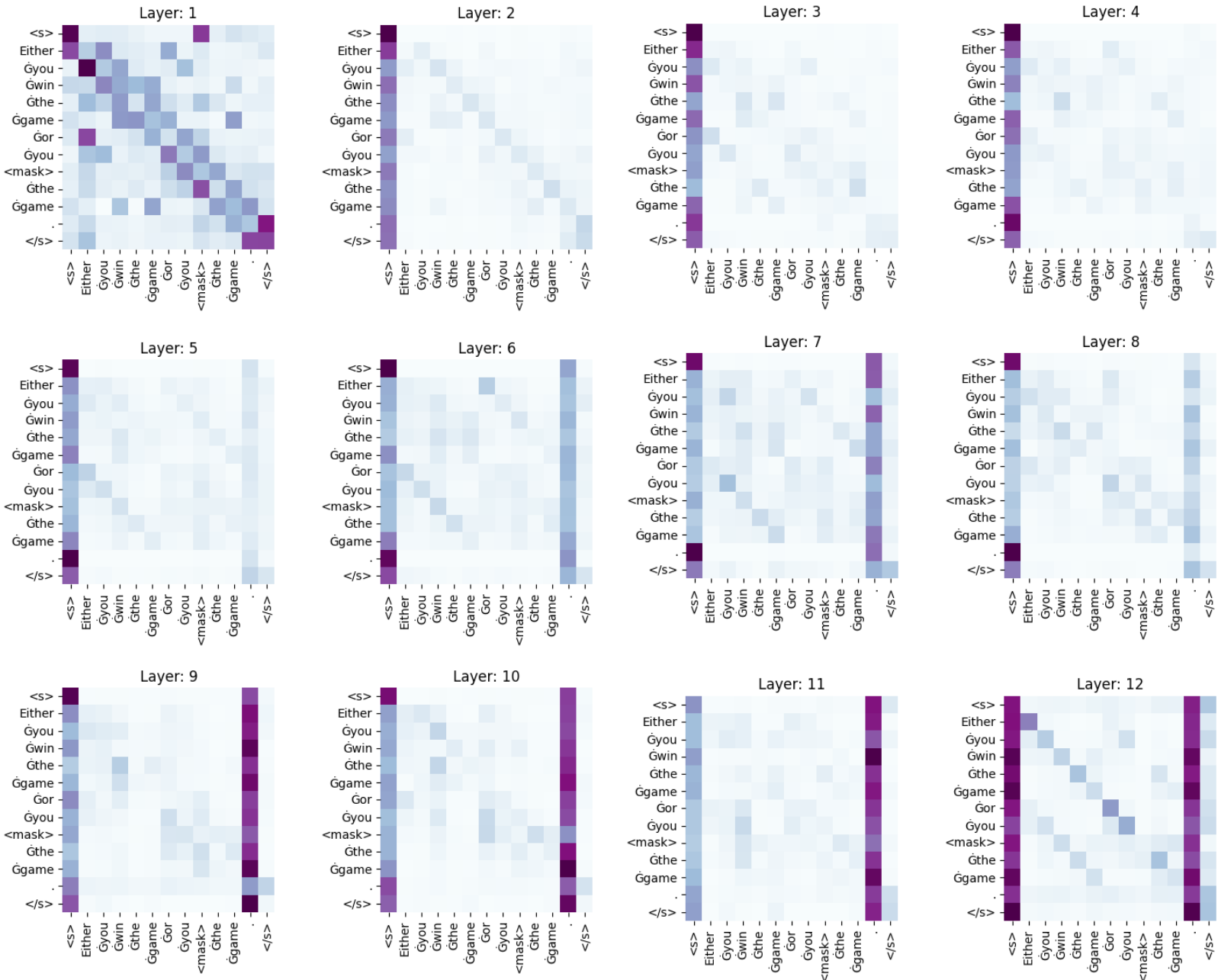
Value Zeroing

Either you win the game or you <mask> the game.



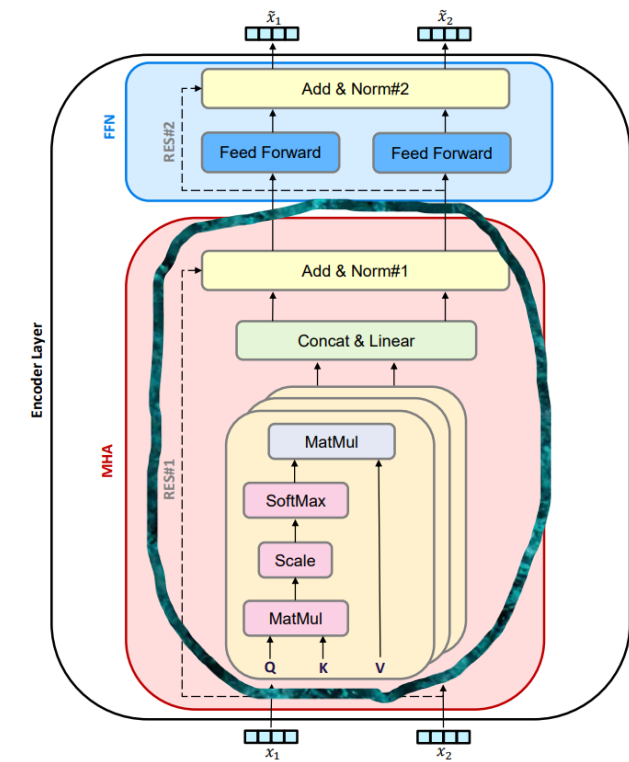
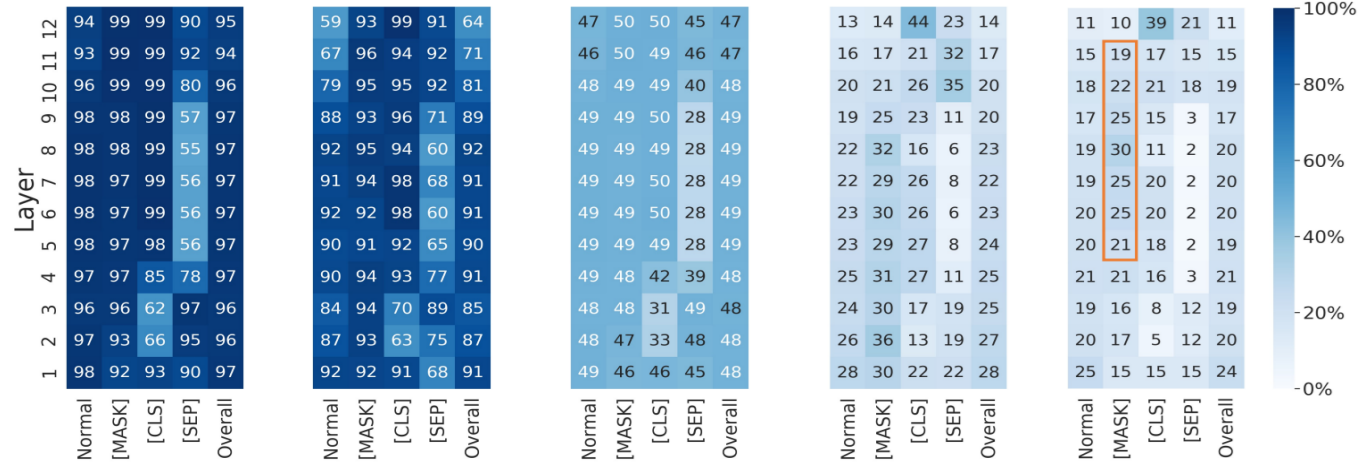
Attention

Either you win the game or you <mask> the game.



Attention-Norm (extended)

$$r_i = \frac{\|\tilde{\mathbf{x}}_{i \leftarrow \text{context}}\|}{\|\tilde{\mathbf{x}}_{i \leftarrow \text{context}}\| + \|\tilde{\mathbf{x}}_{i \leftarrow i}\|}.$$



HTA

$$c_{i,j}^l = \frac{||\nabla_{i,j}^l||_2}{\sum_{k=0}^{d_s} ||\nabla_{k,j}^l||_2}$$

with $\nabla_{i,j}^l = \frac{\delta \mathbf{e}_j^l}{\delta \mathbf{x}_i}$

