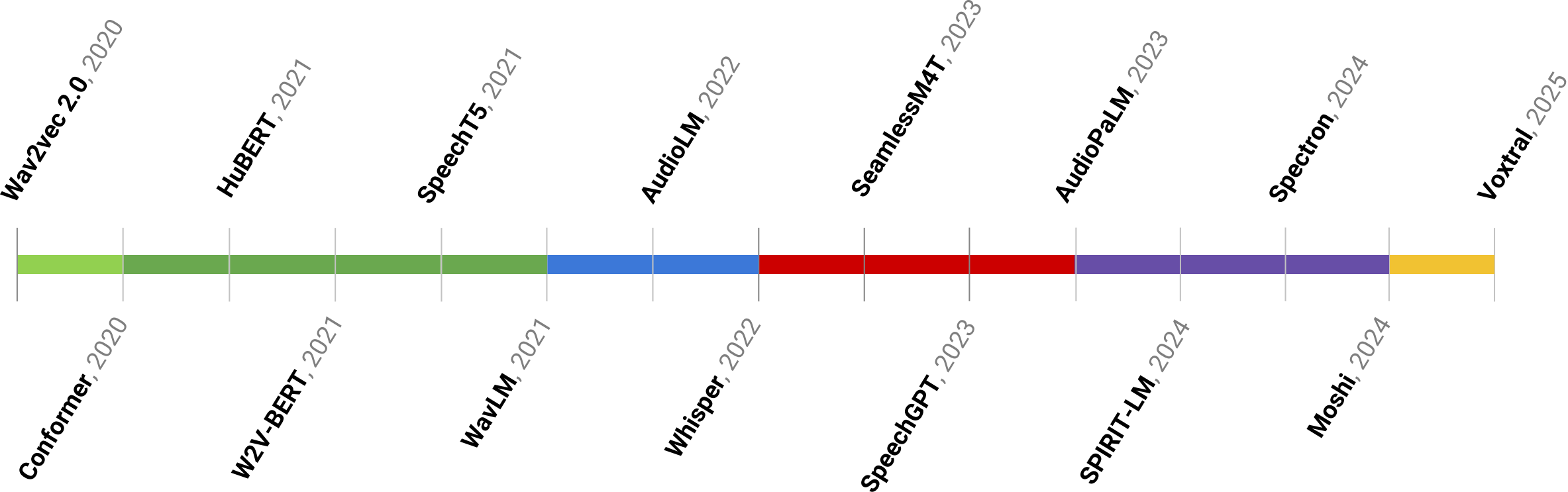# Context-Mixing in Speech Transformers
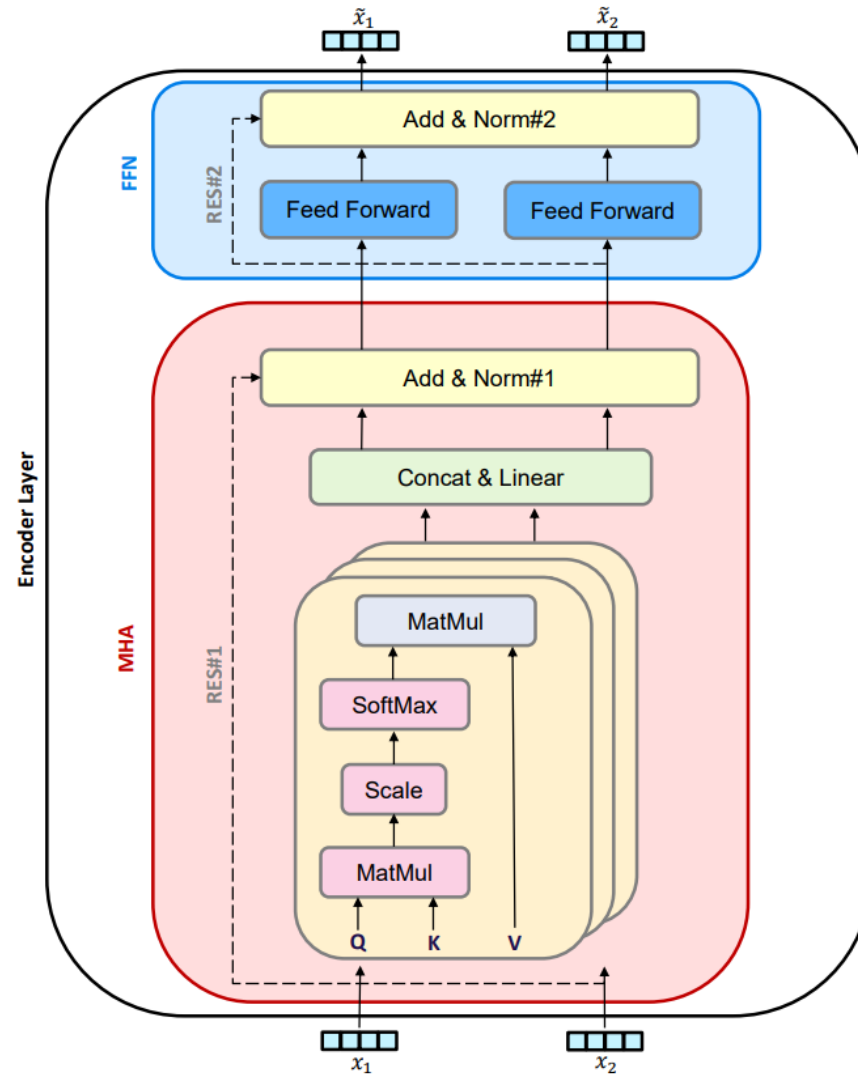
Hosein Mohebbi

August 17, 2025

Rotterdam

# Transformers for speech processing
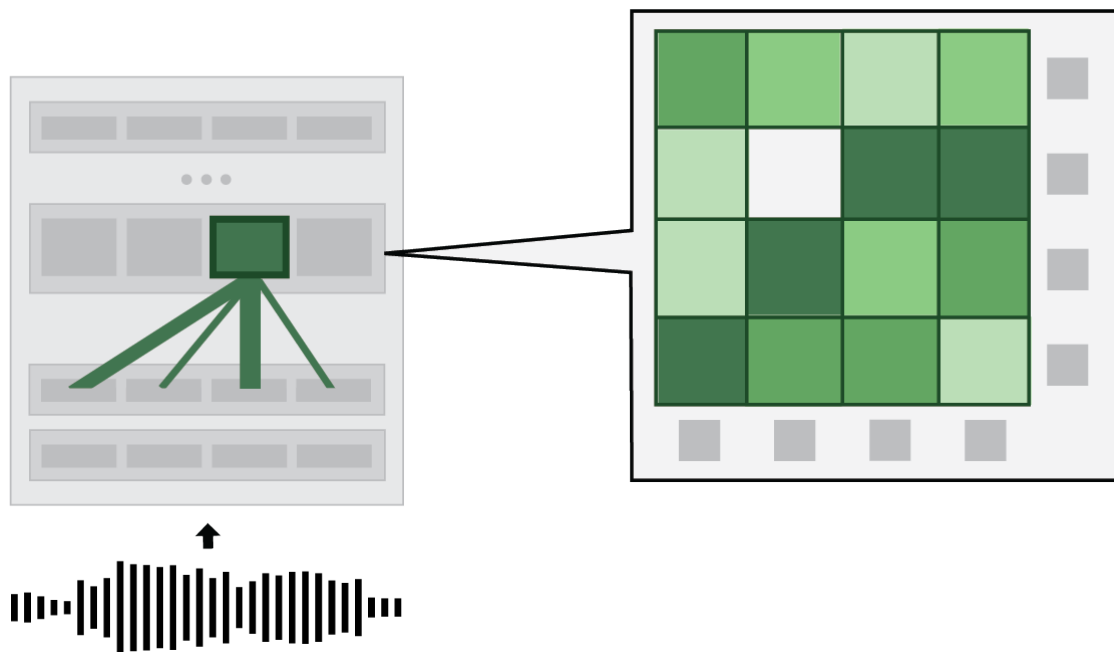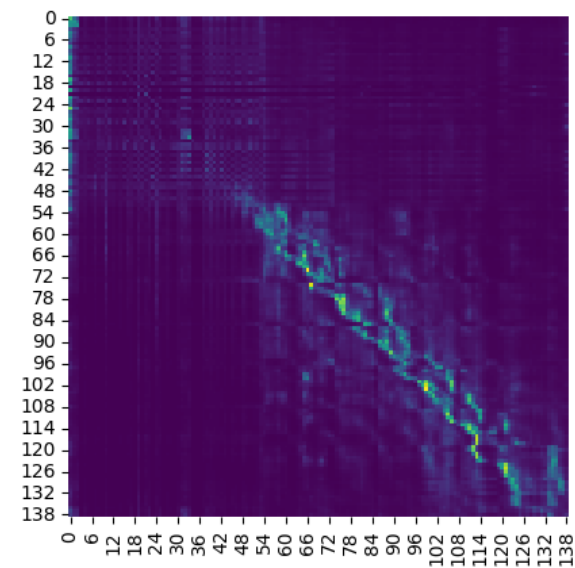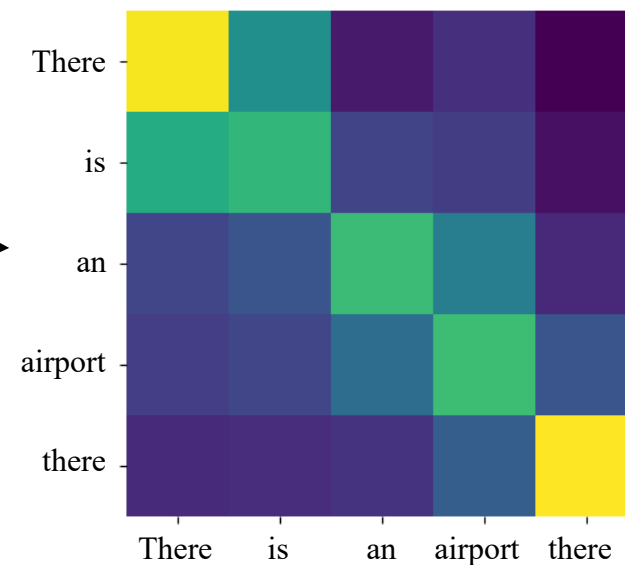
# Transformer



(Vaswani et al., 2017)

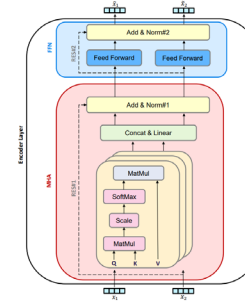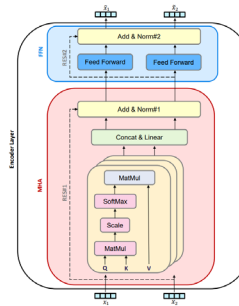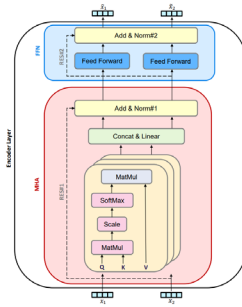# What is Context Mixing?



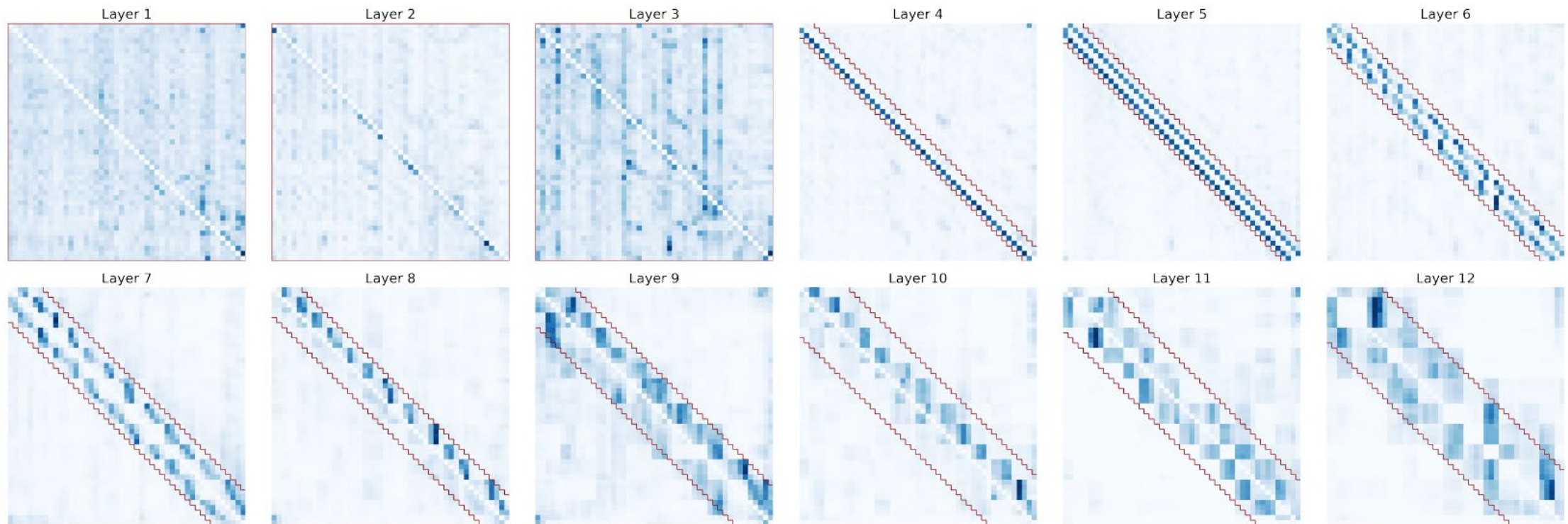Frame-level

Word-level

# Measures of *Context-mixing*

- Attention
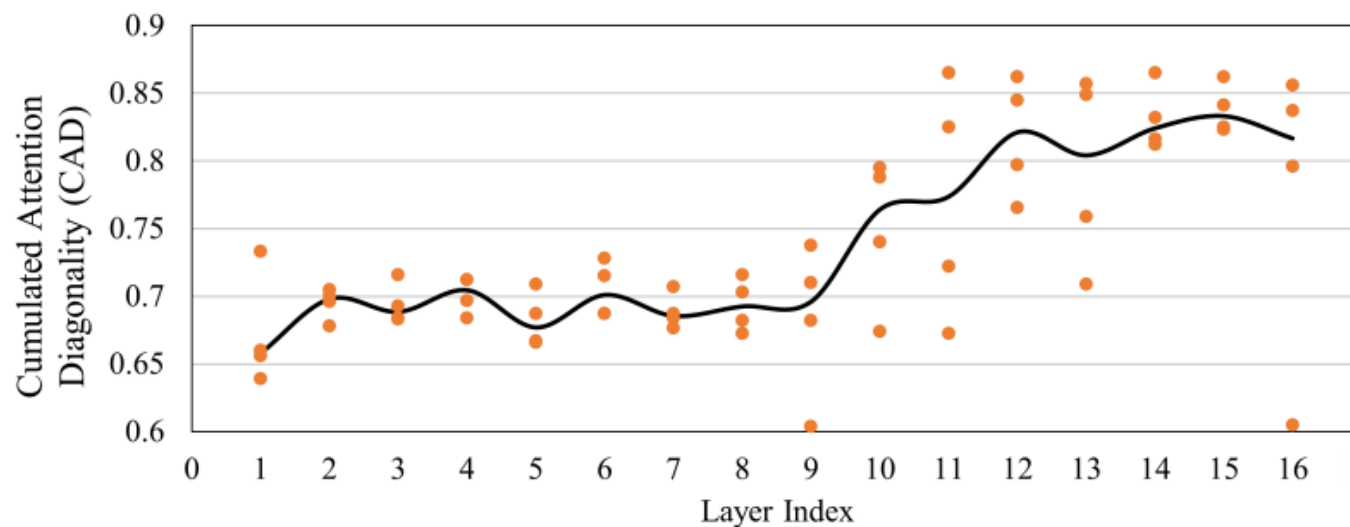
- Attention-Norm

- Value Zeroing

# Diagonality



(Alastruey et al., 2022)

# Two distinct roles

$$\text{CAD}_h = \int_{r=0}^{1} \frac{1}{T} \sum_{i=1}^{T} \left( \sum_{\substack{j=\max(1, \\ i-r(T-1))}}^{\min(T, \\ i+r(T-1))} A_h[i,j] \right) dr = \int_{r=0}^{1} D(r) dr$$





(Shim et al., 2022)

# Diversity Loss

Attention → Highest

Value vectors → Lowest

| Diversity loss | dev | dev-other | test | test-other |
|---|---|---|---|---|
| $d^{\text{A}}(m, n)$ | 6.37 | 6.02 | 6.31 | 6.10 |
| $d^{\text{Q}}(m, n)$ | 0.53 | 0.59 | 0.54 | 0.55 |
| $d^{\text{K}}(m, n)$ | 0.57 | 0.61 | 0.61 | 0.58 |
| $d^{\text{V}}(m, n)$ | 0.13 | 0.14 | 0.13 | 0.14 |

Table 3: *Attention diversity losses summed over all layers of the Conformer acoustic encoder for the baseline full-context Librispeech model.*

(Audhkhasi et al., 2022)

# Attention-Norm



([Kobayashi et al., 2020](#))

# Value Zeroing



(Mohebbi et al., 2023)

# Value Zeroing

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

(Mohebbi et al., 2023)

# Value Zeroing

$$(x_1, ..., x_n)$$

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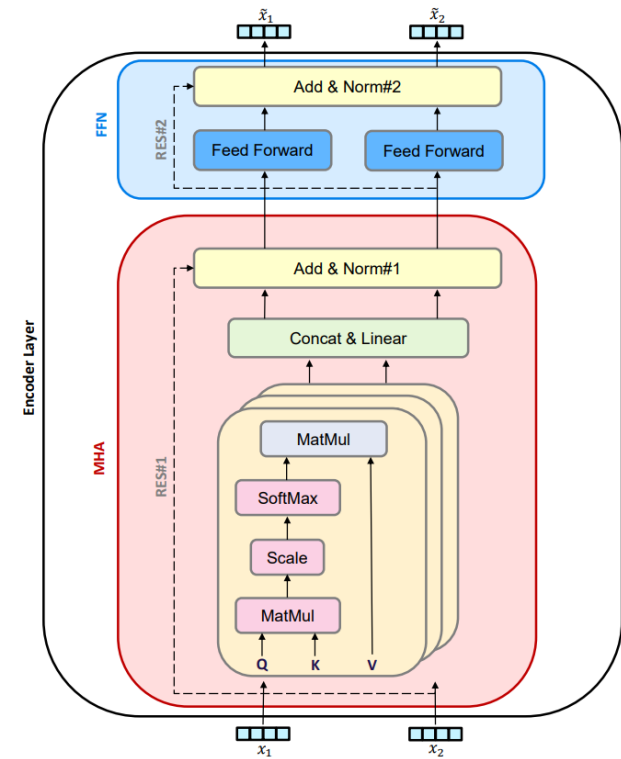\quad \alpha_{i,j} = \underset{x_j \in \mathcal{X}}{\mathrm{softmax}} \left( \frac{q_i k_j^\top}{\sqrt{d}} \right) \in \mathbb{R}
$$

$$v_i^h = x_i W_V^h + b_V^h$$

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

$$z_i = \mathrm{CONCAT}(z_i^1, ..., z_i^H) W_O$$

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

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

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

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

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

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

(Mohebbi et al., 2023)

# Let's Evaluate!

# Evaluation
Controlled task: homophony in French

**Target**

livre (singular)

/livʀ/

livres (plural)

Elle a perdu les **livres**

(She lost the books)

# Evaluation
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<br>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<br><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<br>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 <u>underlined</u> and **bolded**, respectively.

# Cue Contribution score
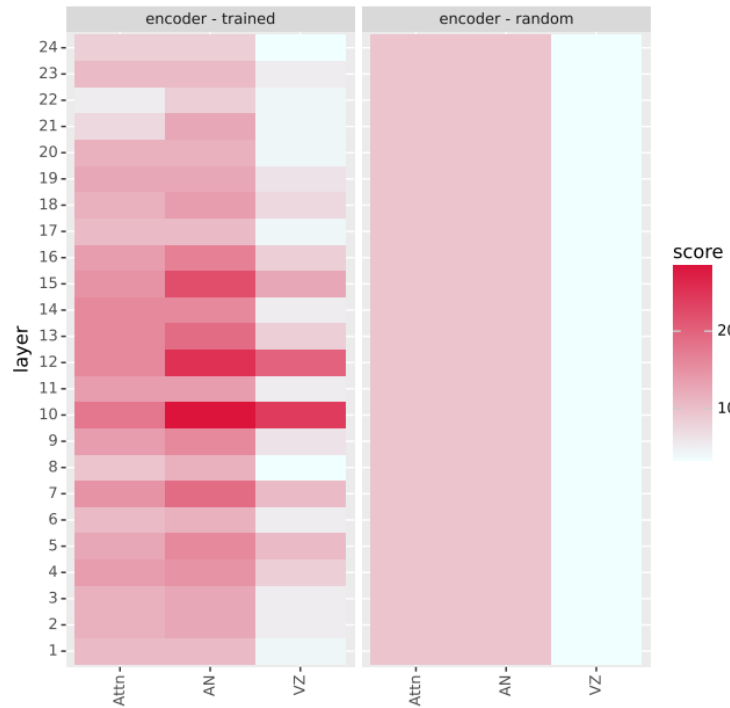


Target word

Cue word

# Cue Contribution



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

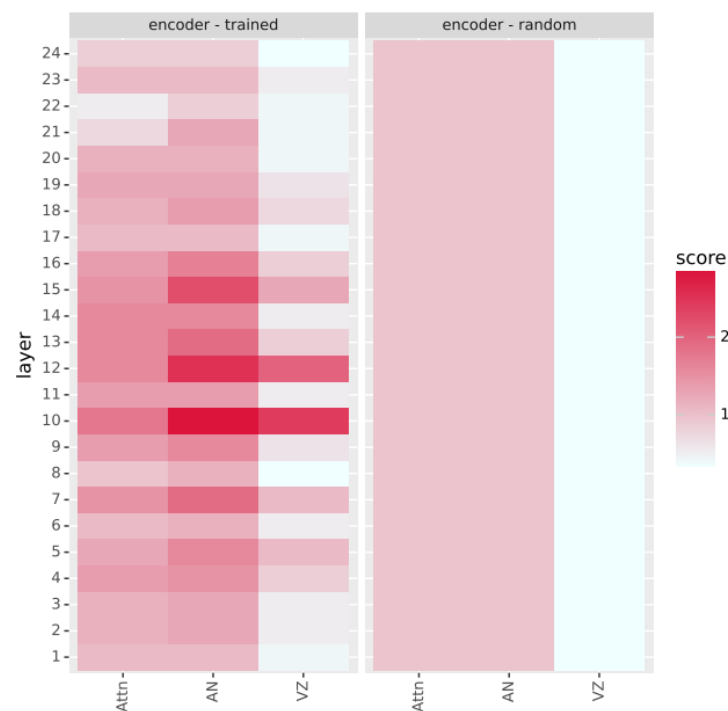(Mohebbi et al., 2023)

# Cue Contribution



Figure 1: Layer-wise cue contribution according to different analysis methods averaged over all examples for XLSR-53, trained (left) vs. randomly initialized (right).
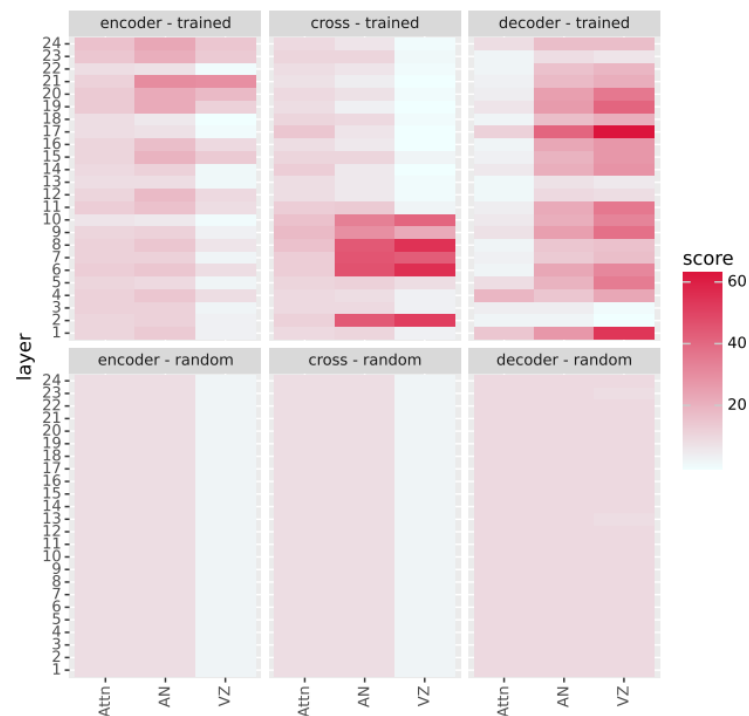
(Mohebbi et al., 2023)

Figure 2: Layer-wise cue contribution according to different analysis methods averaged over all examples for Whisper-medium, trained (top) vs. randomly initialized (bottom).
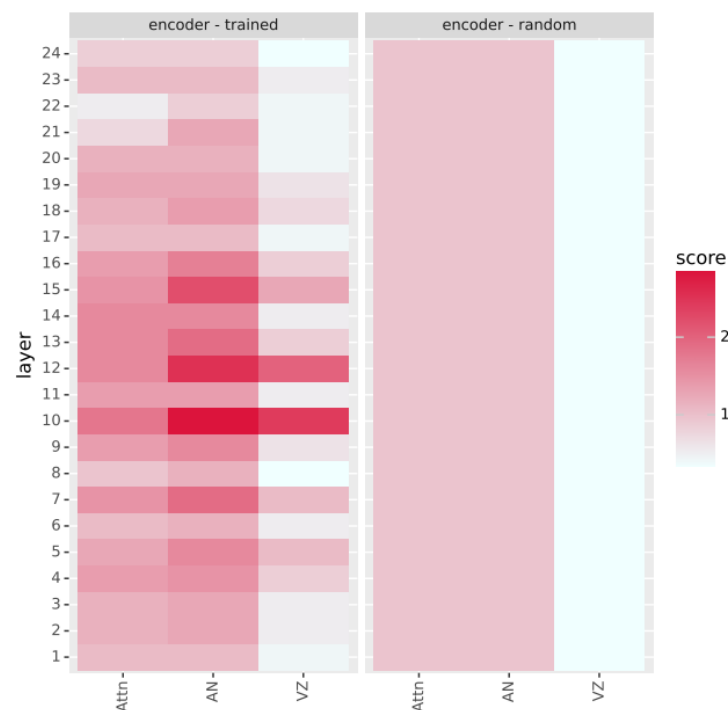
# Cue Contribution v.s. 'Number encoding' probe



Figure 1: Layer-wise cue contribution according to different analysis methods averaged over all examples for XLSR-53, trained (left) vs. randomly initialized (right).
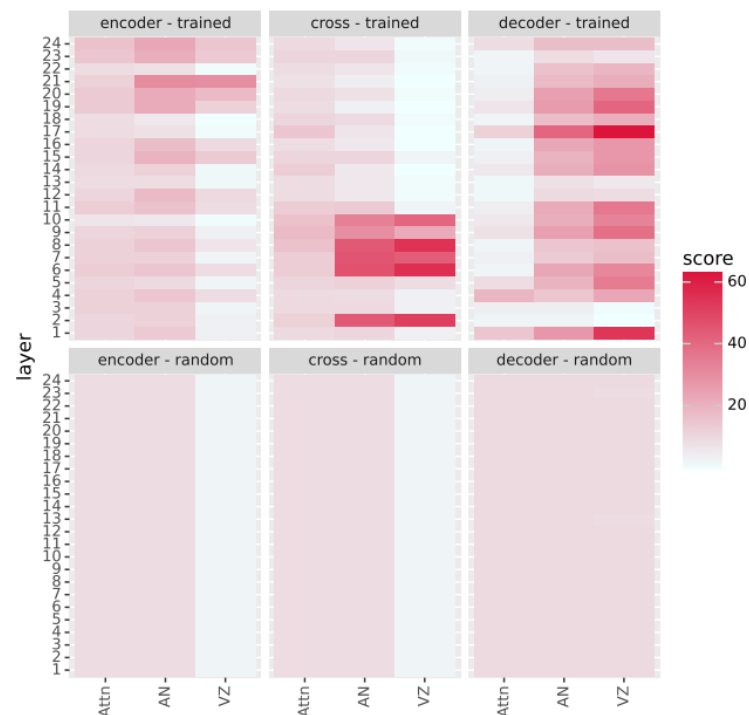
(Mohebbi et al., 2023)



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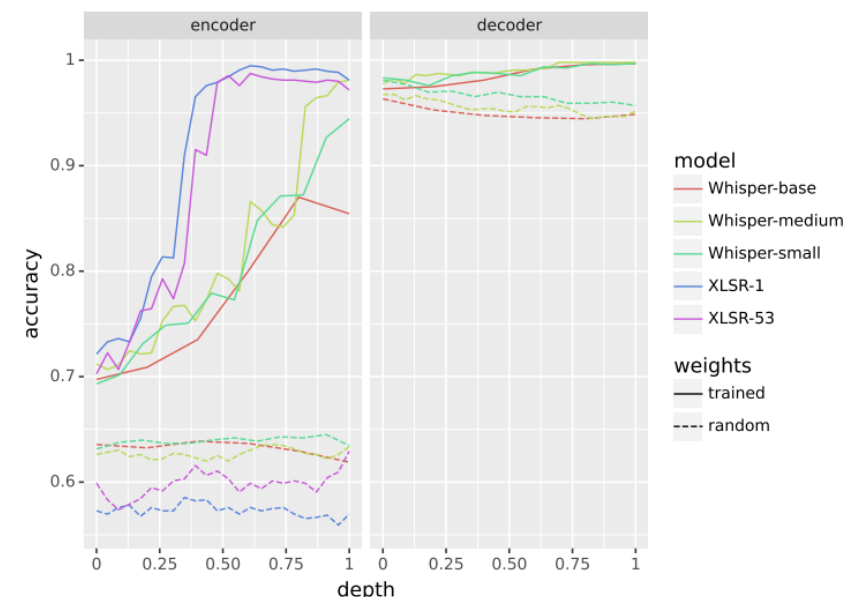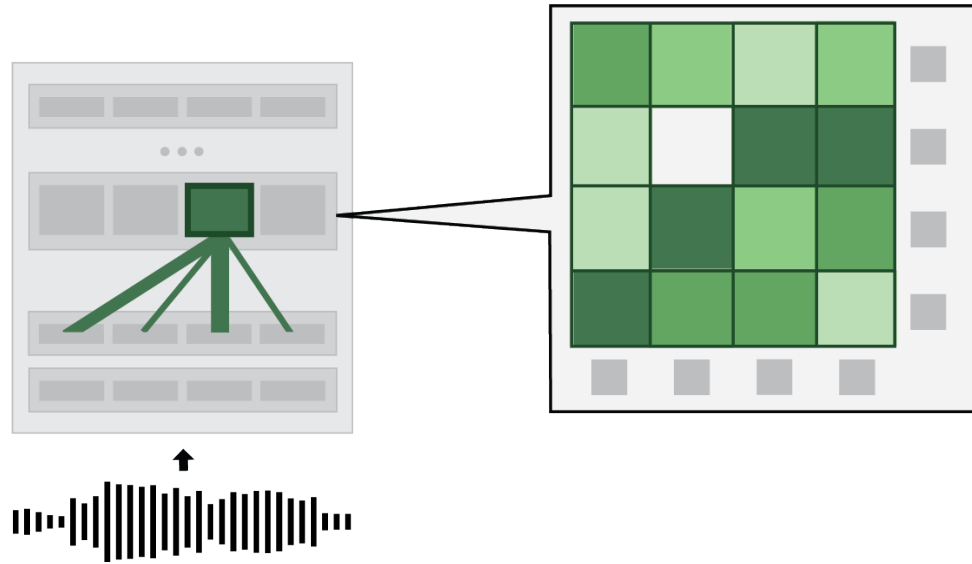


Figure 4: Accuracy of probing classifiers trained on frozen target representations obtained from various ASR models. The depth of Whisper-base (6) and Whisper-small (12), has been normalized to 1 to facilitate comparisons.
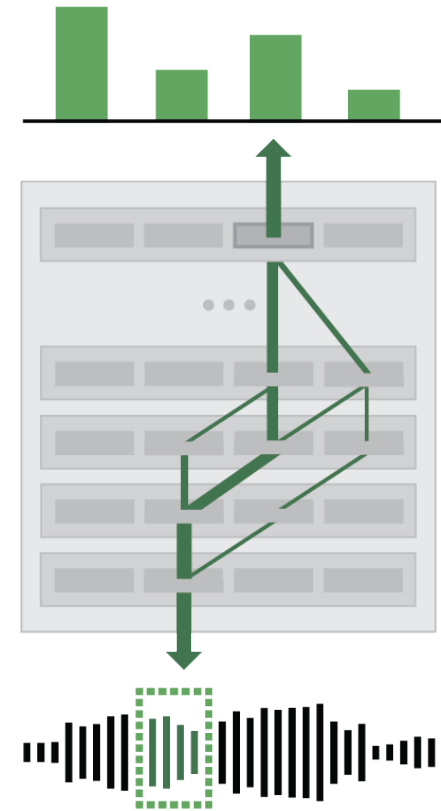
# Wrapping up

- Analyzed the pattern of attentions in speech Transformers (e.g., diagonality, diversity)

- Pointed out the limitations of attention as a measure of context-mixing

- Analyzed context-mixing beyond attention (using e.g., Attention-Norm, Value Zeroing)

- Context-mixing vs. Feature attribution

# Context-Mixing vs. Feature Attribution



(illustration by Marianne)

# Thanks!
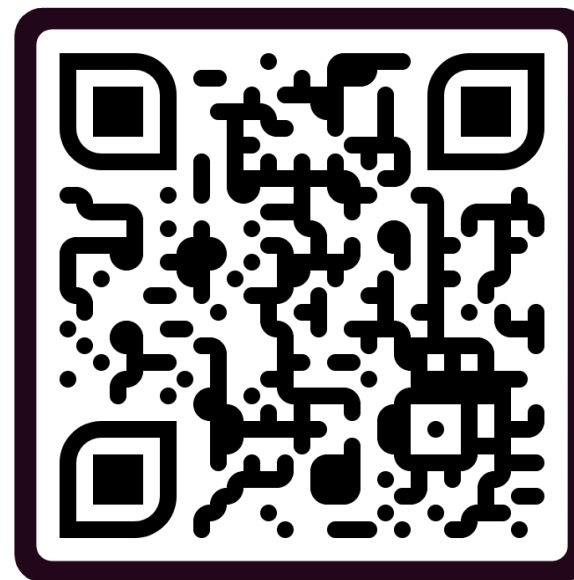# Question?

**Feel free to reach out for any questions:**

**Email:** h.mohebbi@tilburguniversity.edu

**X:** @hmohebbi75

Thank you!

Website

Notebook