How Redundant is the Transformer Stack in Speech Representation Models?



Teresa Dorszewski*, Albert Kjøller Jacobsen*, Lenka Tětková, Lars Kai Hansen Technical University of Denmark, Section for Cognitive Systems



Motivation	Notation
Transformer-based speech representation models perform well but are computationally demanding, limiting their on-device applications.	The general architecture of a transformer-based speech representation models with a word classification output layer:
Recent studies on LLMs and speech models [1, 2, 3, 4, 5] reveal redundancy of transformer layers.	$egin{aligned} ext{Fine-tuned Transformer model} & h_c \Big(g^{(i)} \Big(f^{(i)}(\mathbf{x}) \Big) \Big) = \hat{y} \ & egin{aligned} extbf{A} \ & egin{aligne} extbf{A} \ & $
> Reducing computational requirements allows for efficient inference	

in resource-constrained environments.

Our contribution

Showing redundancy in speech representation models through similarity analyses, pruning and mimicking selected layers.





Representations of the speech transformer at layers *i* and *j* are denoted:

$$A=f^{(i)}(X)\in \mathbb{R}^{n imes d} \qquad B=f^{(j)}(X)\in \mathbb{R}^{n imes d}$$

Layerwise similarity analysis

Similarity analysis of transformer-stack representations using three metrics:

$$egin{aligned} & \mathrm{S}_{cos}(i,j) = rac{1}{n} \sum_{l=1}^n rac{A_{l,\cdot}^T B_{l,\cdot}}{||A_{l,\cdot}|| \cdot ||B_{l,\cdot}||} \ & \mathrm{S}_{CKA}(i,j) = rac{||B^T A||_F^2}{||A^T A||_F ||B^T B||_F} \ & \mathrm{S}_{kNN}(i,j) = rac{1}{n} \sum_{l=1}^n igg(rac{1}{k} |\mathcal{N}_k(A_{l,\cdot}) \cap \mathcal{N}_k(B_{l,\cdot})|) \end{aligned}$$



Pruning the transformer by heuristics





Cosine and kNN similarities are used for block-influence pruning heuristic [1, 6]:

$$\mathrm{BI}(i) = 1 - \mathrm{S}_{cos}(i-1,i)$$

Knowledge distillation by mimicking selected layers



Mimicking networks learn to reproduce selected latent representations.

- > **Step 1: mimicking phase** using MSE loss
- > Step 2: adaption phase using fine-tuning (NLL) loss. Non-mimicker models only learn with NLL loss.



Conclusions

References & sponsors

- We find:
 - 1. **block-like similarity structure** suggesting two main processing steps.
 - that **pruning up to 45% of the layers** of transformer-based speech 2. representation models can be done **with low performance drop**.
 - importance of **both similarity blocks** to maintain performance. 3.
 - that **mimicking layers maintain 95% of original performance while** 4. reducing inference time by up to 87%.
- Xin Men et al., "Shortgpt: Layers in large language models are more redundant than you expect", *arXiv* preprint 2024
- Andrey Gromov et al., "The unreasonable ineffectiveness of the deeper layers", arXiv preprint 2024.
- Anton Razzhigaev et al., "Your transformer is secretly linear", arXiv preprint 2024.
- Ankita Pasad et al, "Comparative layer-wise analysis of self-supervised speech models", ICASSP IEEE 2023
- Teresa Dorszewski et al., "Convexity-based pruning of speech representation models", MLSP IEEE 2024
- Minyoung Huh et al., "The platonic representation hypothesis", ICML 2024 6.

PIONEER CENTRE FOR ARTIFICIAL INTELLIGENCE







International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2025, Hyderabad, India