What and How Can Speech Models Hear?

Evaluation and Analysis of Speech SSL Models' Hearing Abilities.

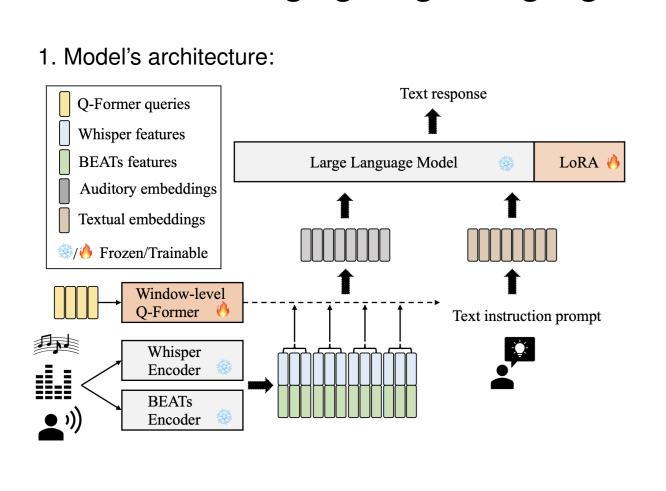


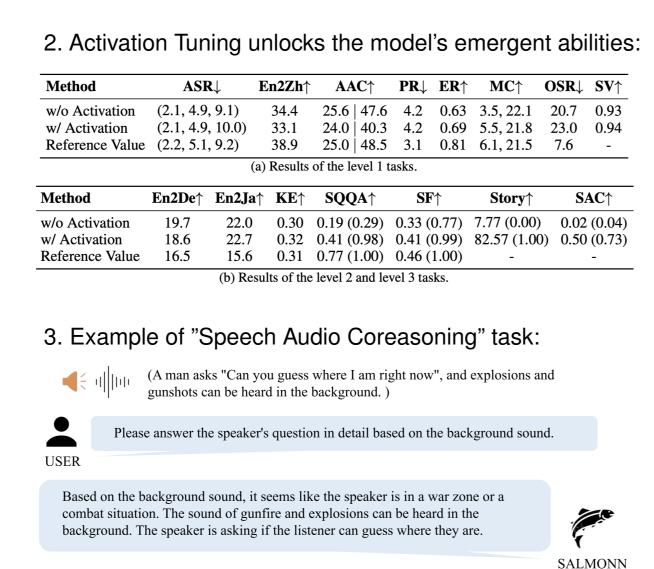
INTRODUCTION & MOTIVATION

Speech self-supervised models have emerged as a central focus in contemporary research. The unprecedented achievements of NLP foundation models serve as an inspiration for advancing speech processing techniques. This study examines the auditory capabilities of speech models to evaluate their comparative advantages and limitations.

SALMONN'S ABILITIES

SALMONN: Bridging Large Language Models and Audio Perception[2]





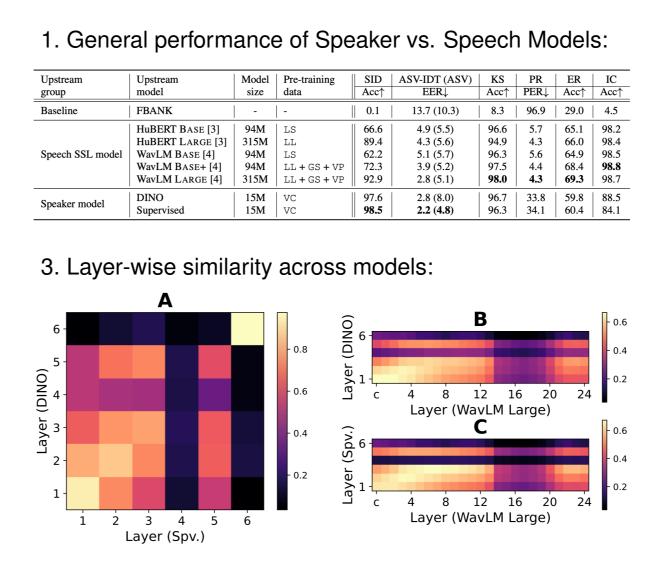
OBJECTIVES

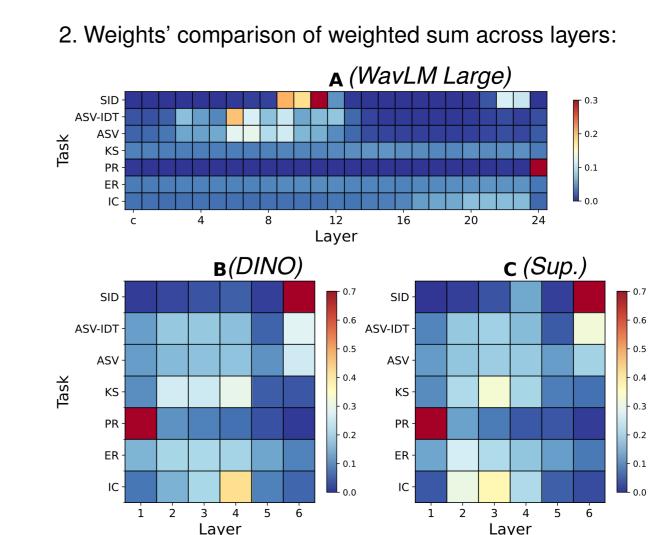
The objectives of this study are:

- To evaluate the auditory processing mechanisms in speech self-supervised models
- To analyze the hierarchical feature representation in speech and speaker models
- To analyze a large-scale evaluation of Speech FMs using (Dynamic-) SUPERB frameworks

SPEAKER VS. SPEECH MODELS

Comparison of models' layers representations of different speech-auditory tasks[1]



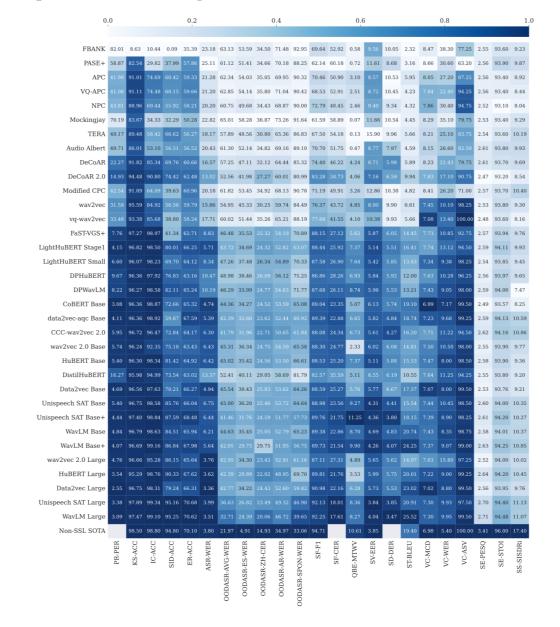


LARGE-SCALE SUPERB BENCHMARKING OF SPEECH MODELS

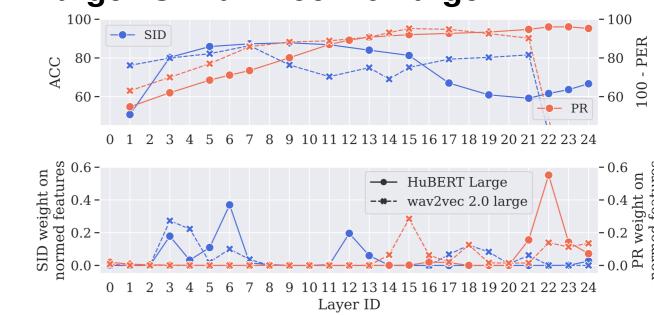
A large scale evaluation of 33 Speech Models across 15 SUPERB tasks[3]

- The study shows that using a pre-trained model with taskspecific heads leads to good generalization across speech tasks.
- The proposed multi-tasking framework proves effective, allowing competitive performance without extensive model finetuning.
- A long-term maintained platform is provided for deterministic evaluation and collaborative benchmarking.
- In-depth analyses reveal how models process information and demonstrate the benefits of few-shot activation tuning.
- The study confirms the reliability of benchmark results through statistical significance and robust evaluation methods.

Heatmap comparison of 33 foundation models across 15 speech processing tasks



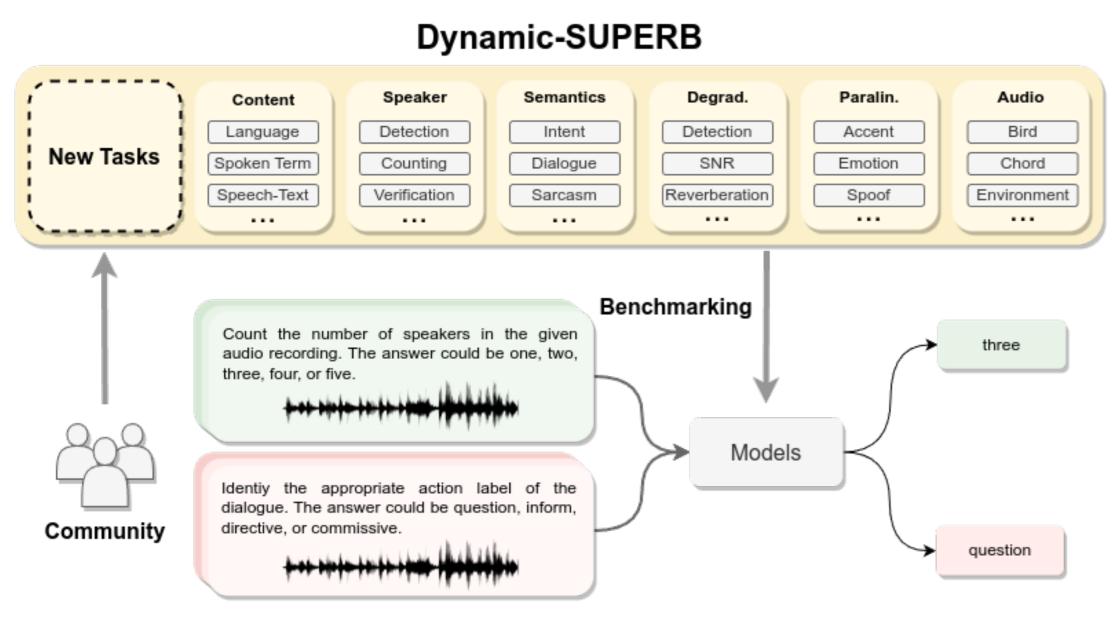
Performance Comparison of Layer Weights by Task and Model: SID and PR in HuBERT Large vs. wav2vec 2.0 Large



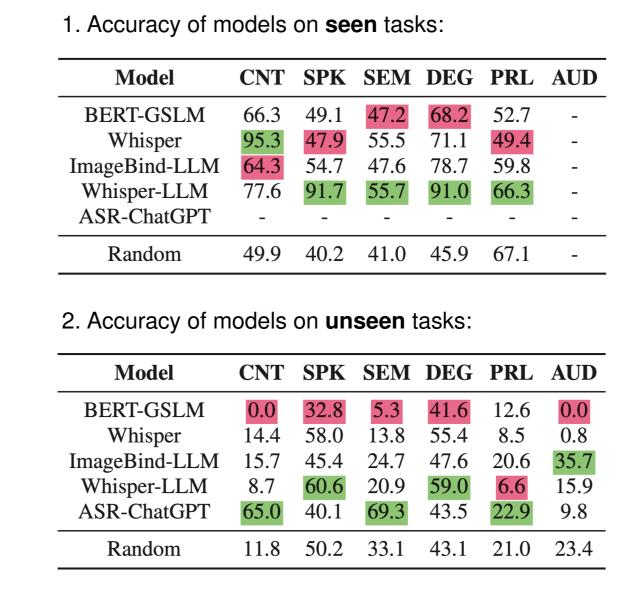
- → The multi-layer approach shows the models' ability to capture a variety of auditory features.
- → Models exhibit different "hearing abilities" for Speech content, Semantics, and Speaker traits.

NEW APPROACH: DYNAMIC SUPERB BENCHMARKING

Introducing the First Dynamic, Collaborative Benchmark for Speech Instruction Tuning: Covering 33 Tasks and 55 Evaluation Instances[4]



Performance results of multiple speech models on seen tasks (with seen/unseen instructions) and on unseen tasks:



- \rightarrow Models often rely on pattern recognition. rather than true semantic comprehension.
- → The need for a more robust instruction understanding and generalization. capabilities.

REFERENCES

- [1] T. Ashihara, M. Delcroix, T. Moriya, K. Matsuura, T. Asami, and Y. Ijima. What do self-supervised speech and speaker models learn? new findings from a cross model layer-wise analysis, 2024.
- [2] C. Tang, W. Yu, G. Sun, X. Chen, T. Tan, W. Li, L. Lu, Z. Ma, and C. Zhang. Salmonn: Towards generic hearing abilities for large language models, 2024.
- [3] S. wen Yang, H.-J. Chang, Z. Huang, A. T. Liu, C.-I. Lai, H. Wu, J. Shi, X. Chang, H.-S. Tsai, W.-C. Huang, T. hsun Feng, P.-H. Chi, Y. Y. Lin, Y.-S. Chuang, T.-H. Huang, W.-C. Tseng, K. Lakhotia, S.-W. Li, A. Mohamed, S. Watanabe, and H. yi Lee. A large-scale evaluation of speech foundation models, 2024.
- [4] C. yu Huang, K.-H. Lu, S.-H. Wang, C.-Y. Hsiao, C.-Y. Kuan, H. Wu, S. Arora, K.-W. Chang, J. Shi, Y. Peng, R. Sharma, S. Watanabe, B. Ramakrishnan, S. Shehata, and H. yi Lee. Dynamic-superb: Towards a dynamic, collaborative, and comprehensive instruction-tuning benchmark for speech, 2024.

CONCLUSION

- Models differ not only in their embeddings of auditory features, but also across different tasks.
- Activation tuning is important for harnessing the emerging abilities of speech models.
- Research could benefit from training on multi-modal data.

MORE DETAILS

