

# SING: Spatial Context in Large Language Model for Next-Gen Wearables



Ayushi Mishra  
(co-primary)



Yang Bai  
(co-primary)



Priyadarshan  
Narayanasamy



Nakul Garg



Nirupam Roy

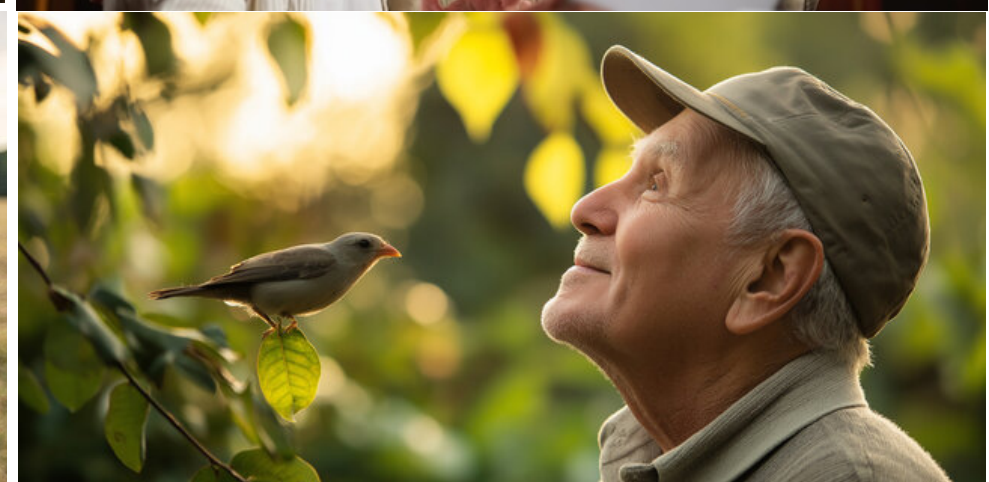
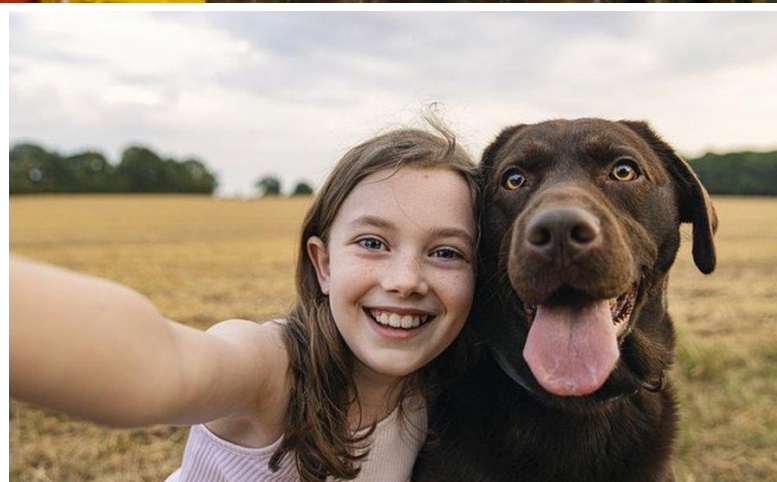


**ICML** July  
International Conference  
On Machine Learning **2025**

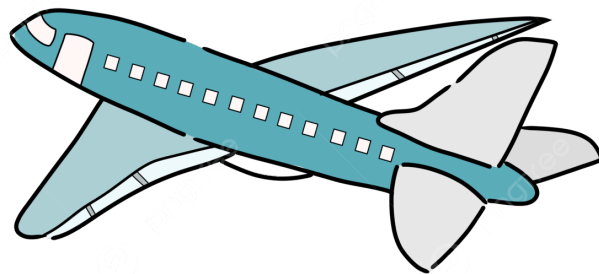




# *A World Alive with Sound*







alexa



Google  
Assistant



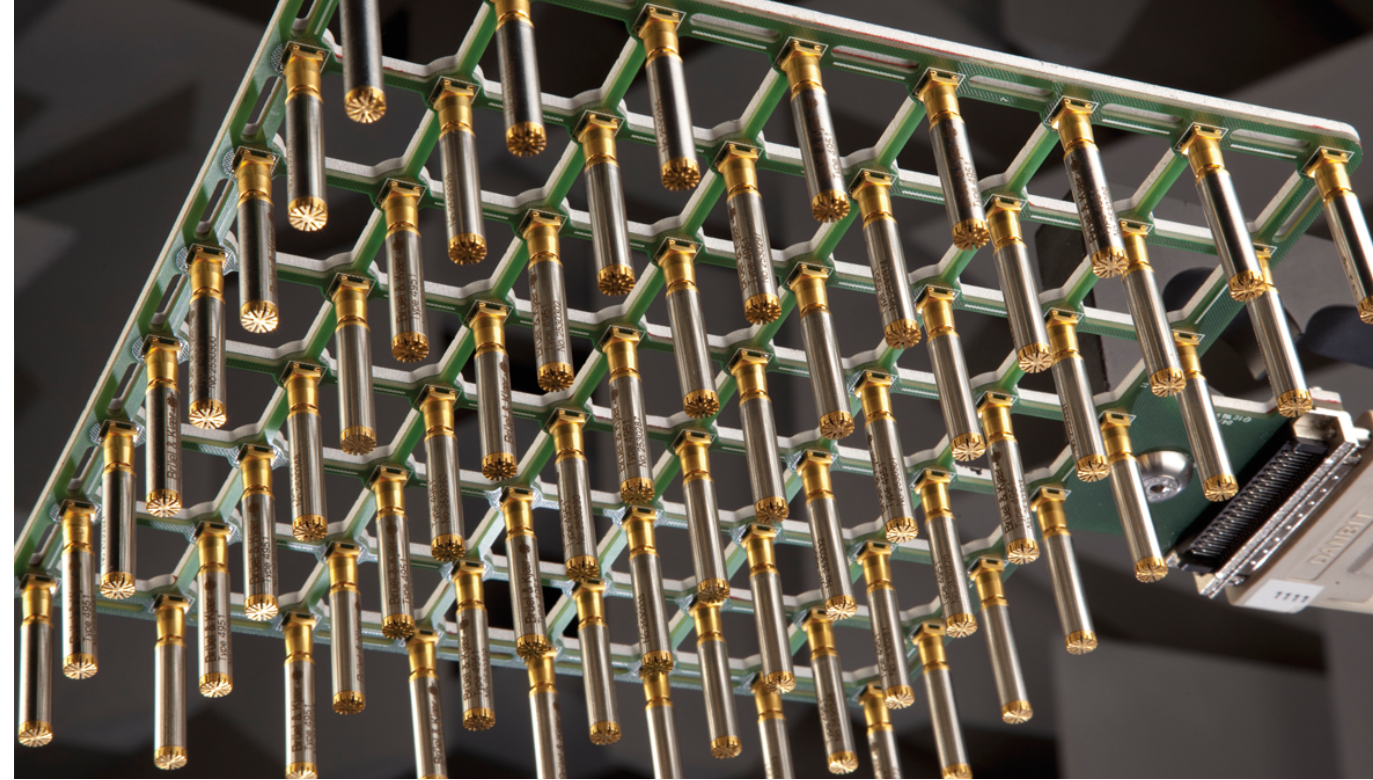
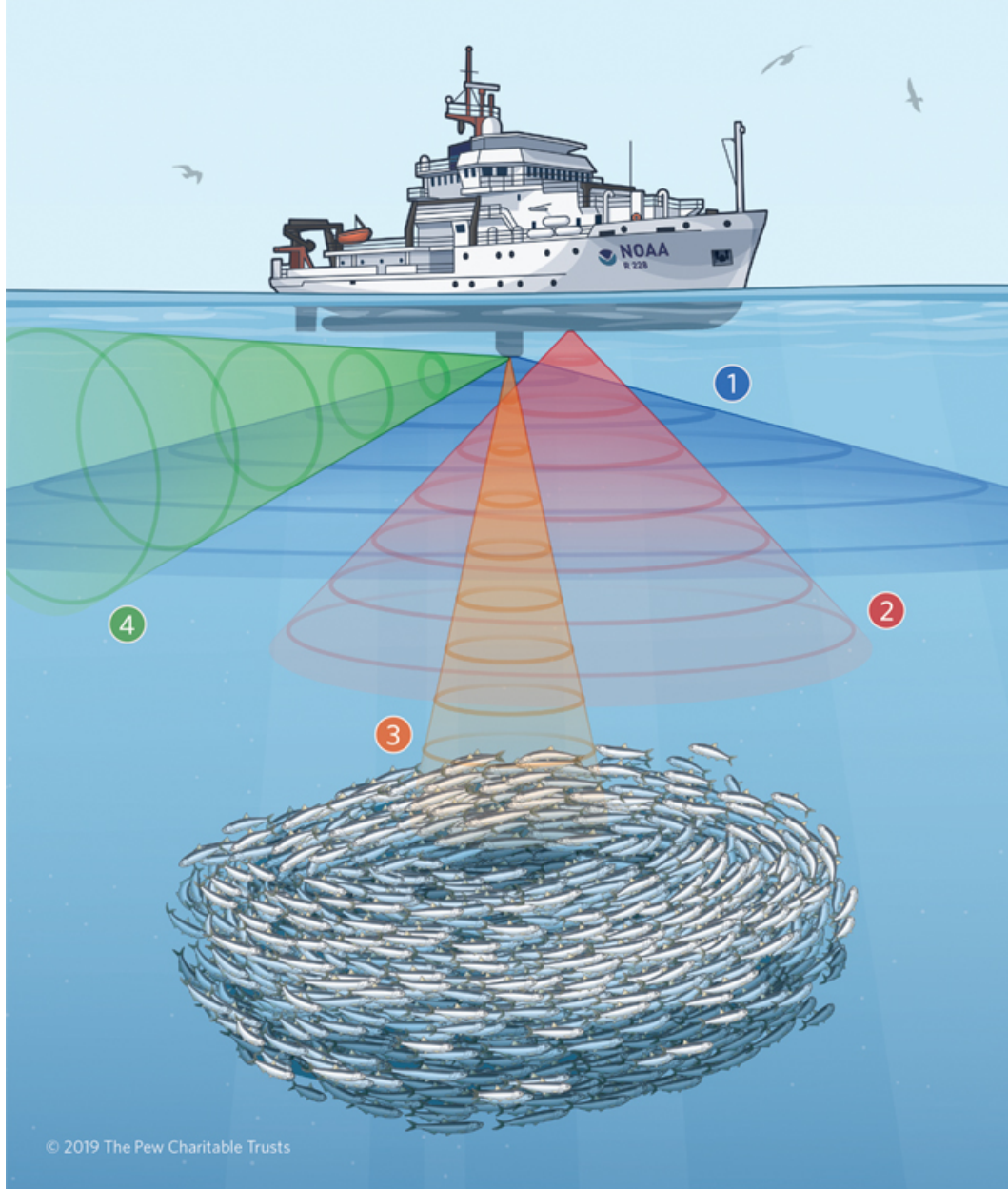
Apple AirPods



*Sound Becomes Smart*



# Spatial perception with acoustics



Sonar and microphone array system:  
For spatial analysis of sound

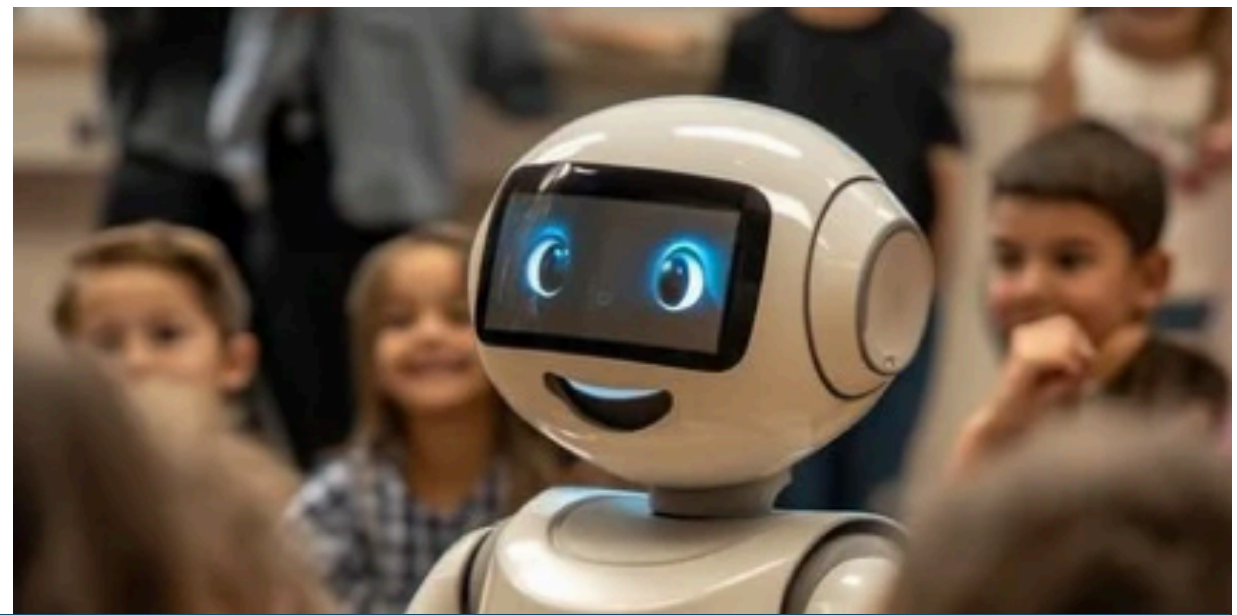


# Spatial acoustic sensing in nature



Animals use body structure for spatial audio.

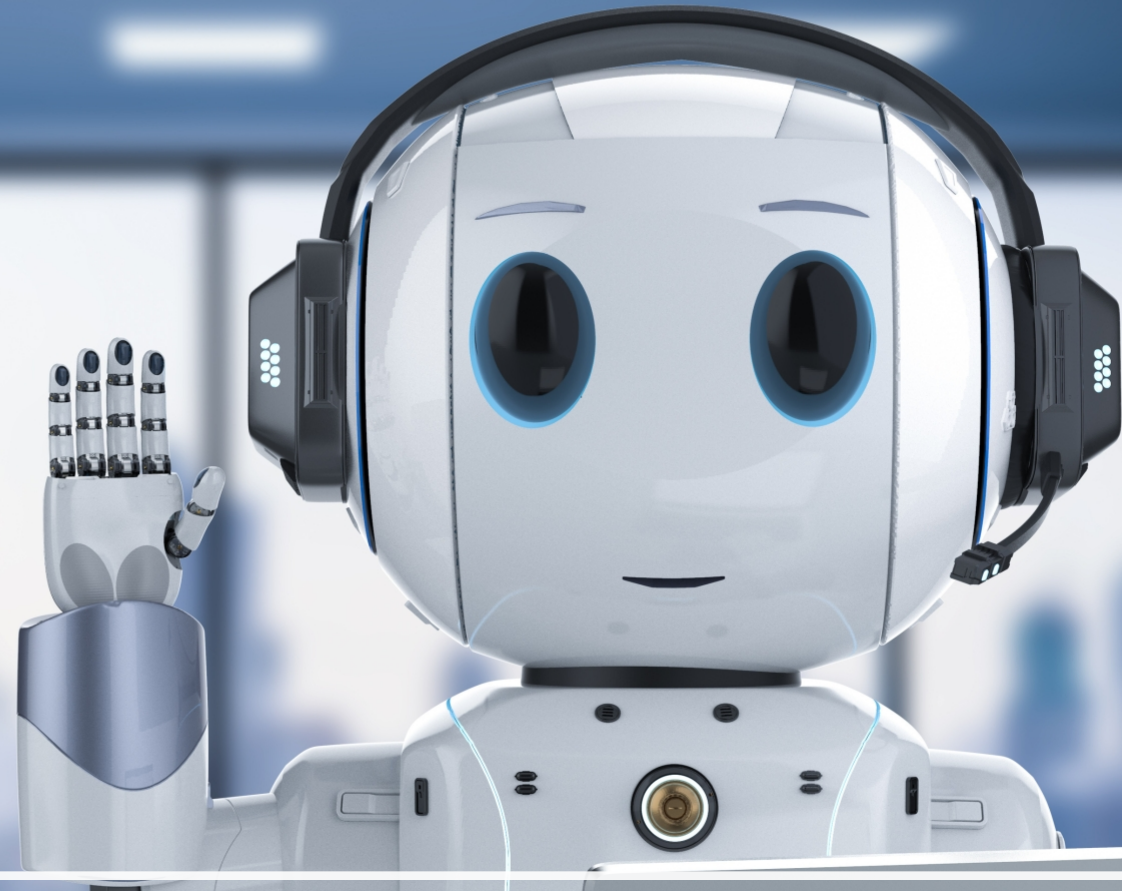




## *Spatial Soundscape*



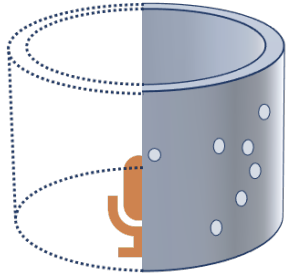




Can Large Language Model understand Spatial Speech?

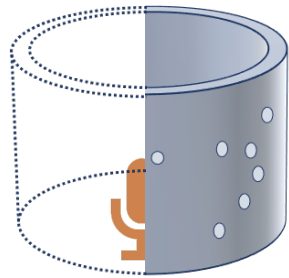


# Core idea: Structure-assisted spatial sensing

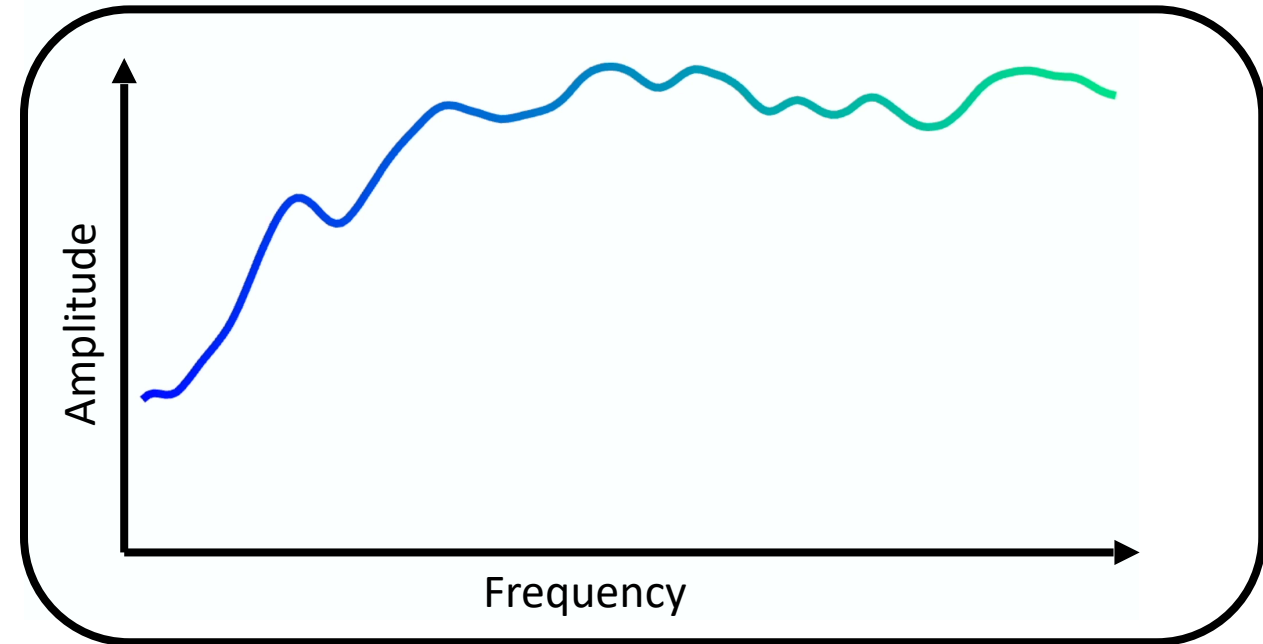


**Metamaterial**

# Core idea: Structure-assisted spatial sensing



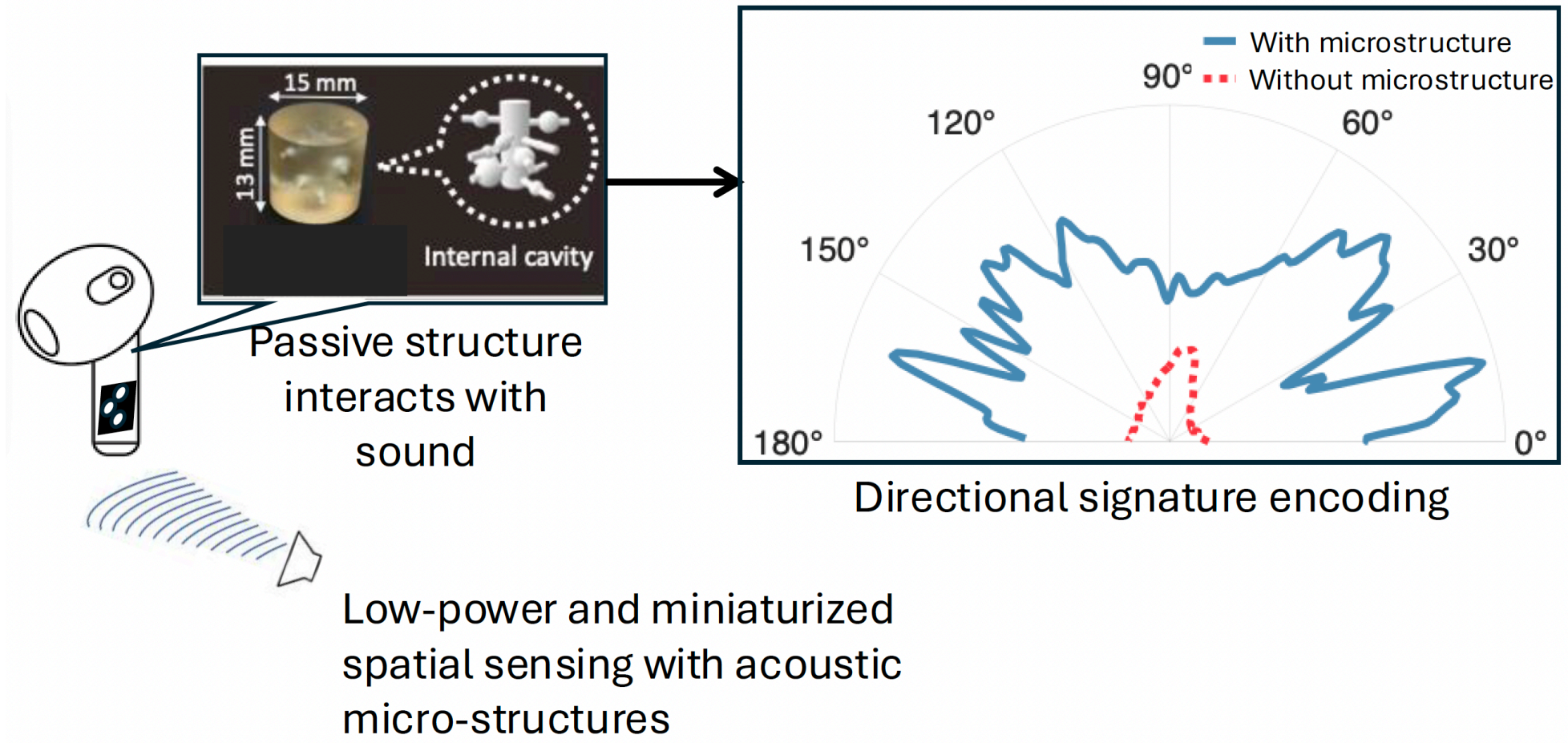
**Metamaterial**



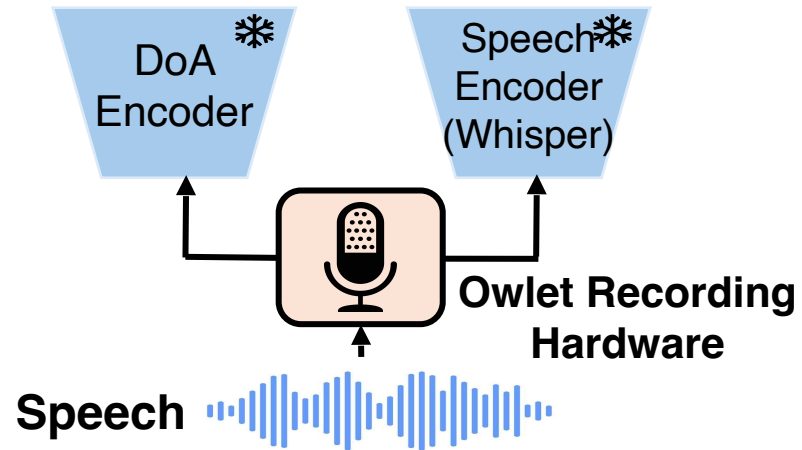
**Received sound signature**



# Spatial Speech

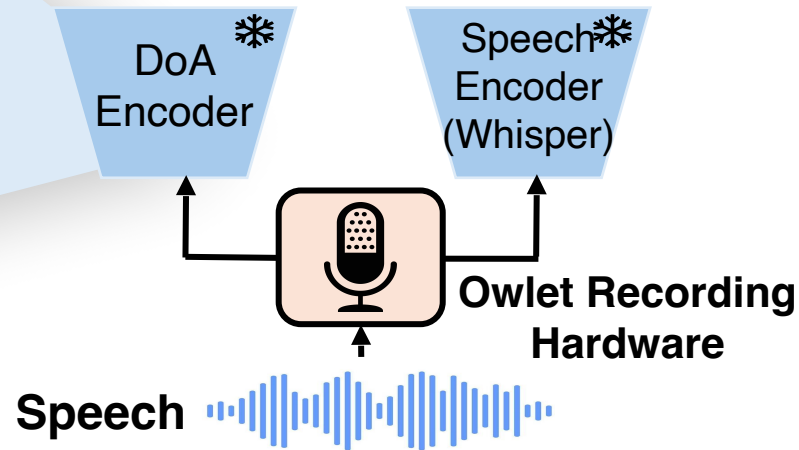
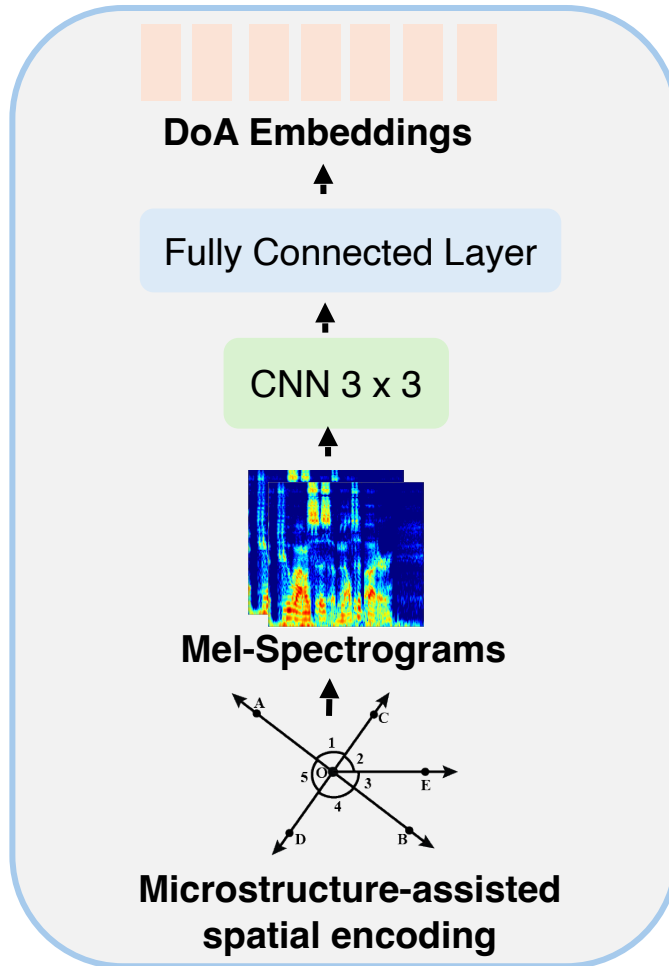


# Extract speech feature

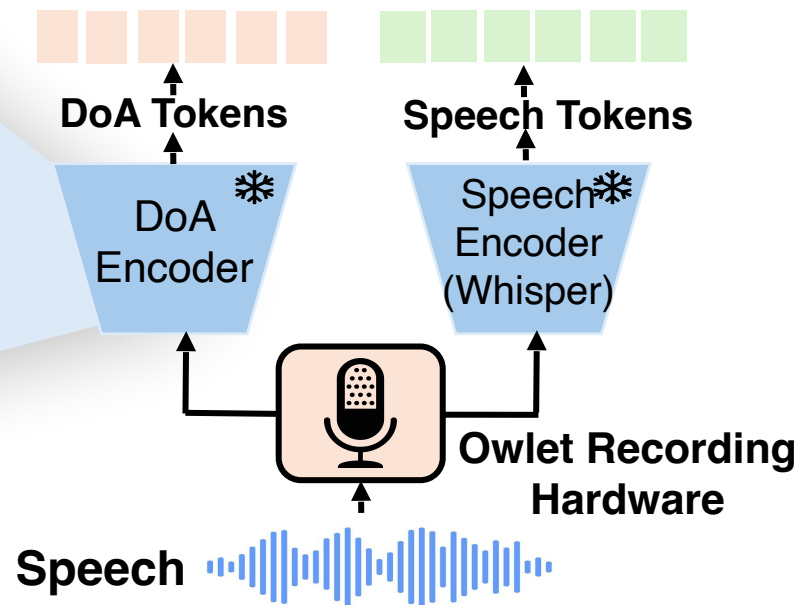
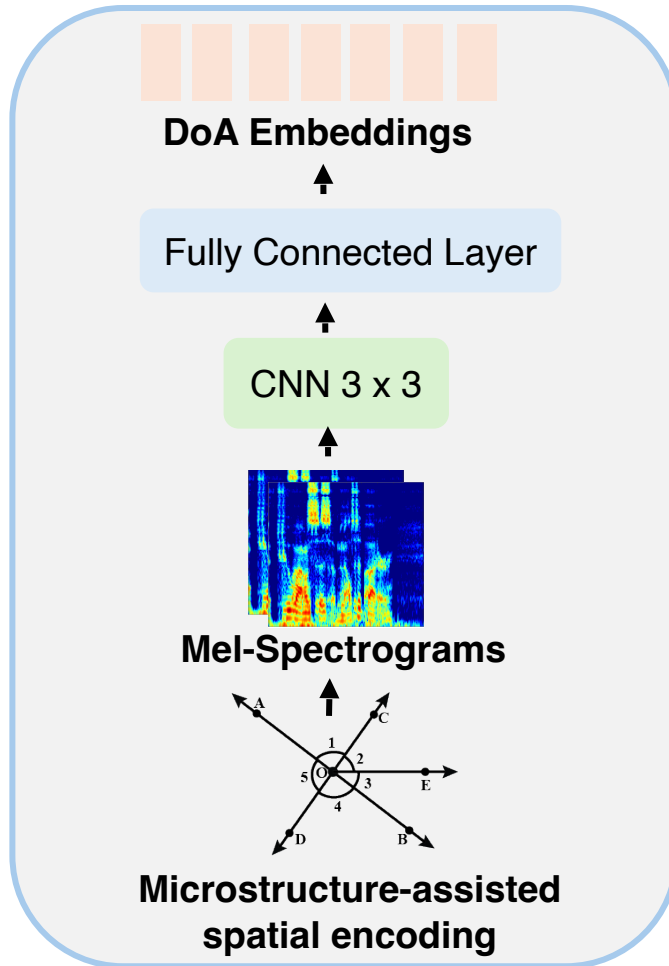




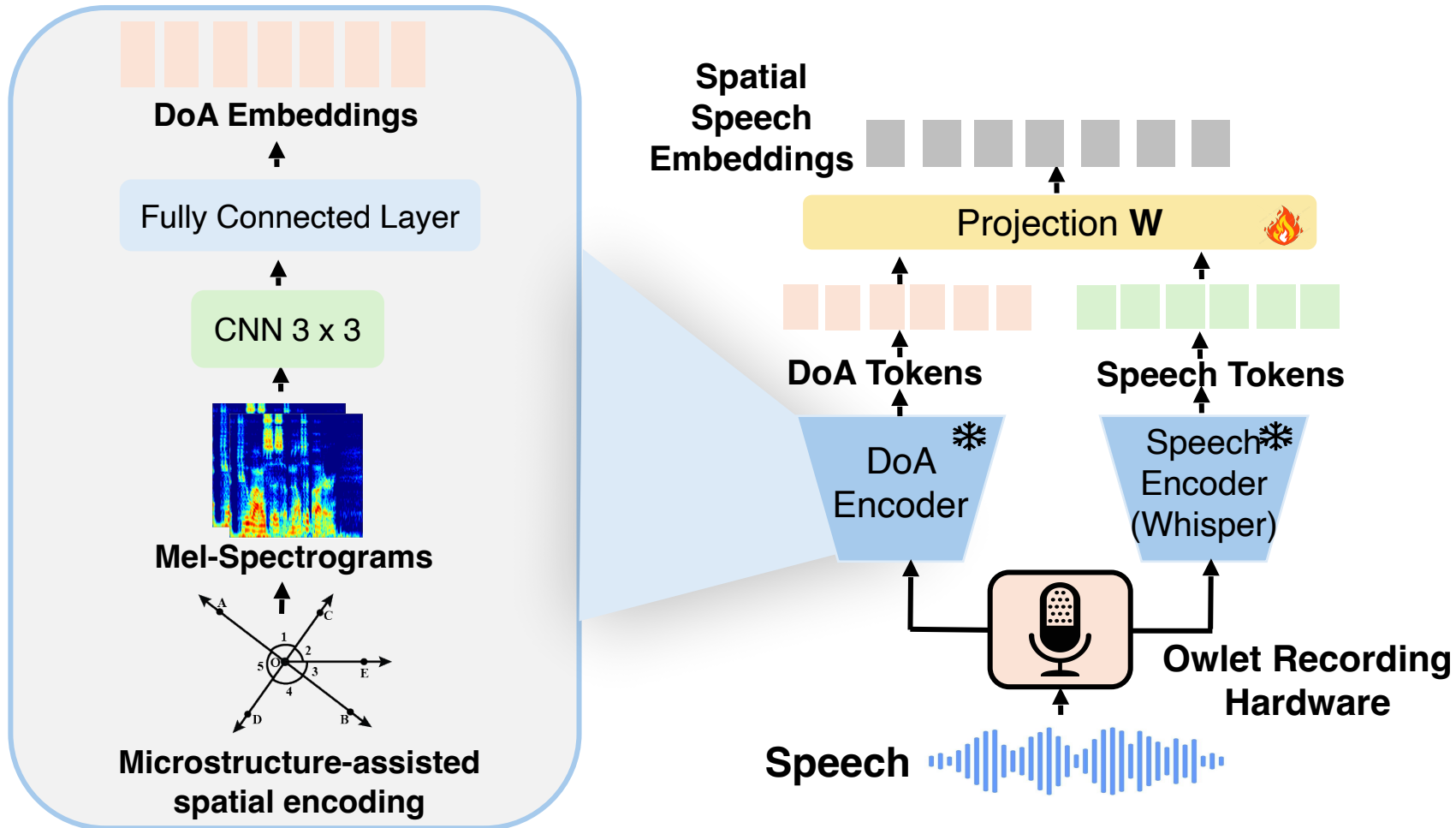
# Extract speech feature



# Extract speech feature

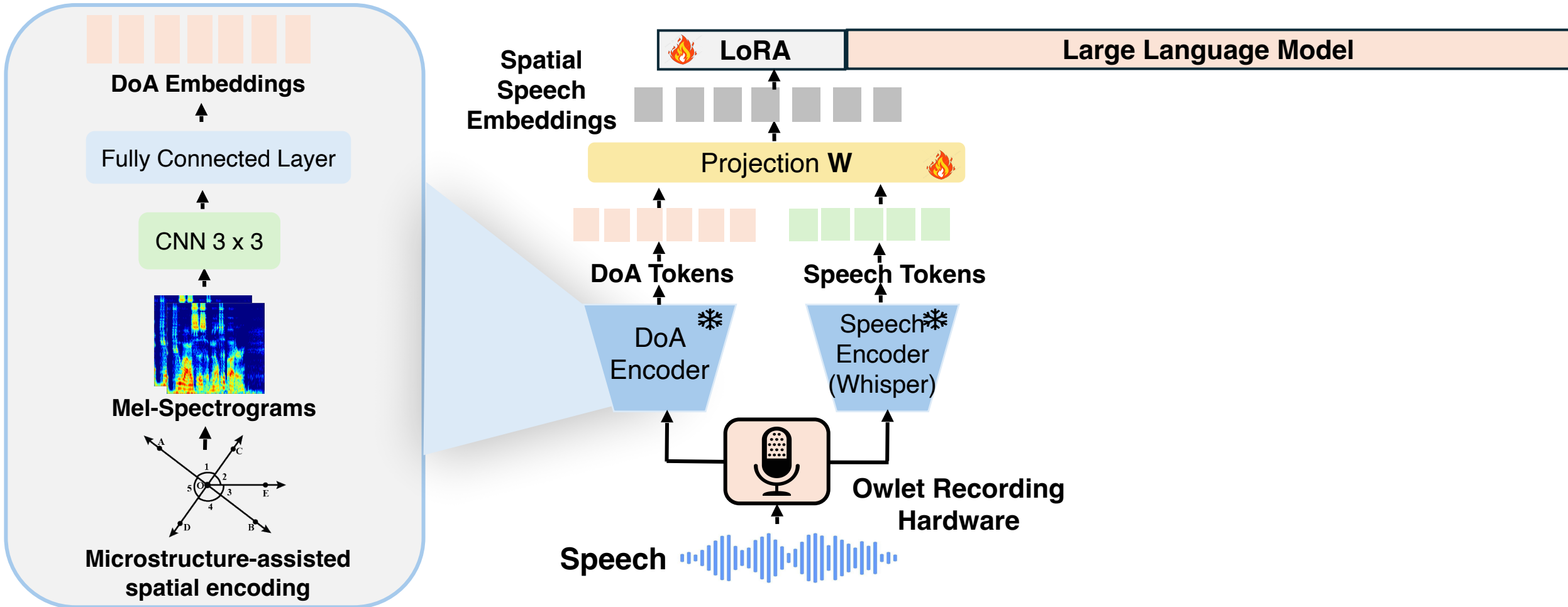


# Alignment with LLM space

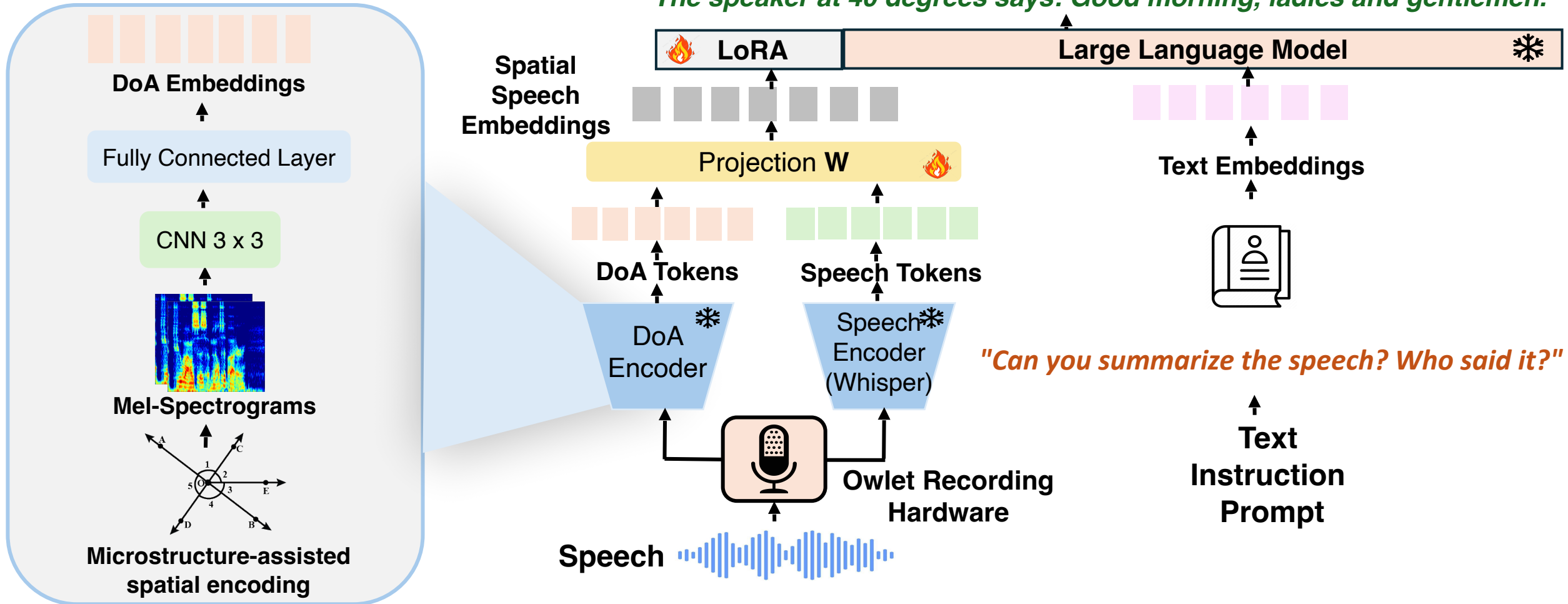




# Fine-tuning LLM for spatial speech



# Fine-tuning LLM for spatial speech





Model	Supported Task	Metric / Performance
BAT	DoA Estimation + Audio Source Recognition (No Speech)	MAE (°) ↓: 88.52
SING [our work]	DoA + Speech ASR (Spatial Awareness)	MAE (°) ↓: <b>25.72</b>
SALMONN	Speech ASR and LLM Functions (No Spatial Awareness)	WER (%) ↓: 2.2
SING [our work]	Speech ASR (without DoA)	WER (%) ↓: <b>1.8</b>
SING [our work]	DoA + Speech ASR (Spatial Awareness)	WER (%) ↓: 5.3

Model	Metric	1 Source	2 Sources	3 Sources	4 Sources	5 Sources
SELDNet	MAE (↓)	90.03	✗	✗	✗	✗
	MEEM (↓)	360.12	✗	✗	✗	✗
	Median Error (↓)	90.14	✗	✗	✗	✗
AudioMAE	MAE (↓)	43.79	✗	✗	✗	✗
	MEEM (↓)	43.79	✗	✗	✗	✗
	Median Error (↓)	27.79	✗	✗	✗	✗
SING (Ours)	MAE (↓)	<b>25.72</b>	<b>24.16</b>	<b>28.11</b>	<b>23.31</b>	<b>17.08</b>
	MEEM (↓)	<b>25.72</b>	<b>24.16</b>	<b>28.11</b>	<b>23.31</b>	<b>17.08</b>
	Median Error (↓)	<b>13.00</b>	<b>13.00</b>	<b>20.00</b>	<b>18.00</b>	<b>13.00</b>

# Thank you!

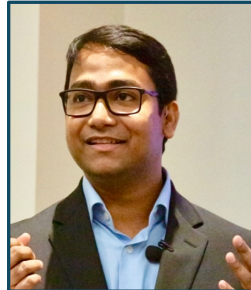
## SING: Spatial Context in Large Language Model for Next-Gen Wearables

[http://icosmos.cs.umd.edu/images/2\\_publication/papers/SING\\_ICML\\_2025\\_icosmos.pdf](http://icosmos.cs.umd.edu/images/2_publication/papers/SING_ICML_2025_icosmos.pdf)

### Contacts:



**Ayushi Mishra**  
amishr13@umd.edu



**Nirupam Roy**  
niruroy@umd.edu



<http://icosmos.cs.umd.edu>