

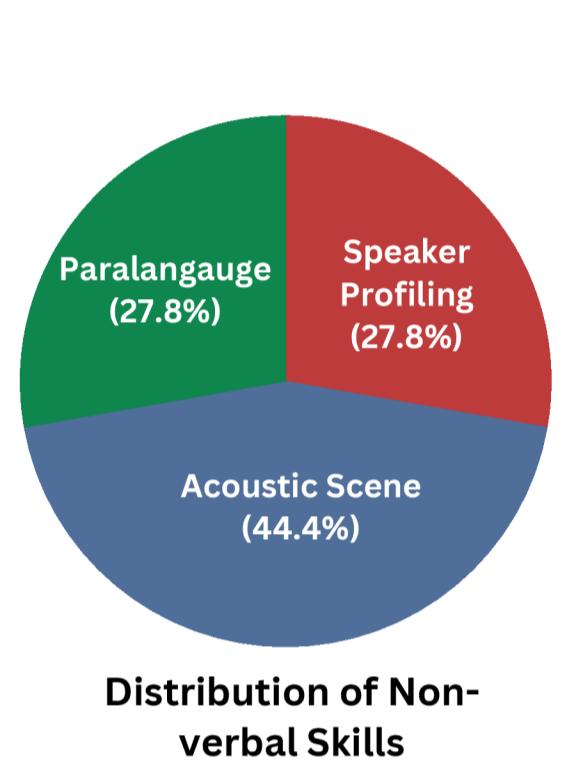
MultiVox: A Benchmark for Evaluating Voice Assistants for Multimodal Interactions



Ramaneswaran Selvakumar, Ashish Seth, Nishit Anand
Utkarsh Tyagi, Sonal Kumar, Sreyan Ghosh, Dinesh Manocha
University of Maryland, College Park

Omni-modal Language Models (OLMs) are increasingly being used as voice assistants capable of understanding both visual and spoken inputs

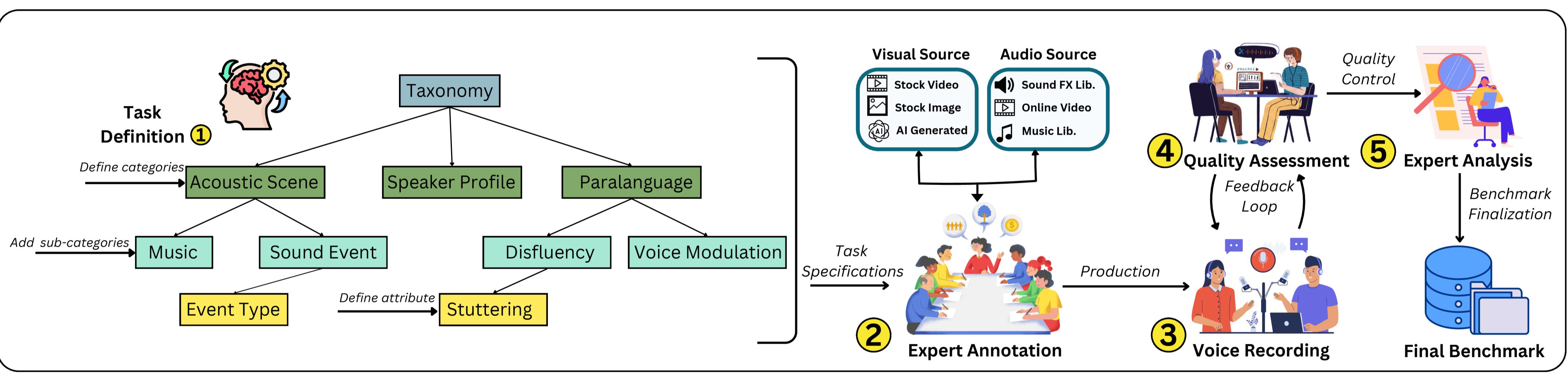
High Level Categorization in MultiVox



We introduce MultiVox, a multi-modal benchmark with 1000 human-spoken utterances to evaluate OLMs as assistants



Benchmark Construction Process



Empirical Results

Model	Acoustic Scene	Paralanguage	Speaker Profile	Avg. CA
Human	4.37	4.33	4.36	4.35
Mini-Omni	1.53	1.79	2.01	1.74
VITA 1.5	2.60	2.56	3.01	2.69
VideoLlama2	1.52	1.59	1.38	1.50
Baichuan-Omni	1.90	2.25	2.01	2.02
Mini CPM	2.87	2.35	2.90	2.69
Intern Omni	2.54	1.94	2.64	2.35
Phi4-MM	2.26	2.63	2.48	2.44
Qwen 2.5 Omni	3.19	2.98	3.06	3.08
Gemini 2.5 Flash	3.55	3.19	3.64	3.44
Gemini 2.5 Pro	3.65	3.32	3.74	3.56

Confounder Analysis

