

Dynamic SUPERB



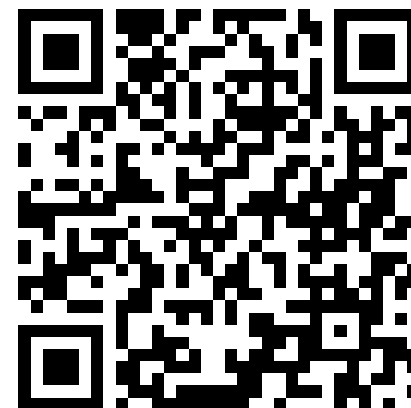
DYNAMIC-SUPERB: TOWARDS A DYNAMIC, COLLABORATIVE, AND COMPREHENSIVE INSTRUCTION-TUNING BENCHMARK FOR SPEECH

*Chien-yu Huang¹, Ke-Han Lu^{*1}, Shih-Heng Wang^{*1}, Chi-Yuan Hsiao^{†1}, Chun-Yi Kuan^{†1}, Haibin Wu^{†1}
Siddhant Arora^{§2}, Kai-Wei Chang^{§1}, Jiatong Shi², Yifan Peng², Roshan Sharma², Shinji Watanabe²
Bhiksha Ramakrishnan^{2,3}, Shady Shehata³, Hung-yi Lee¹*

¹National Taiwan University, Taiwan, ²Carnegie Mellon University, USA

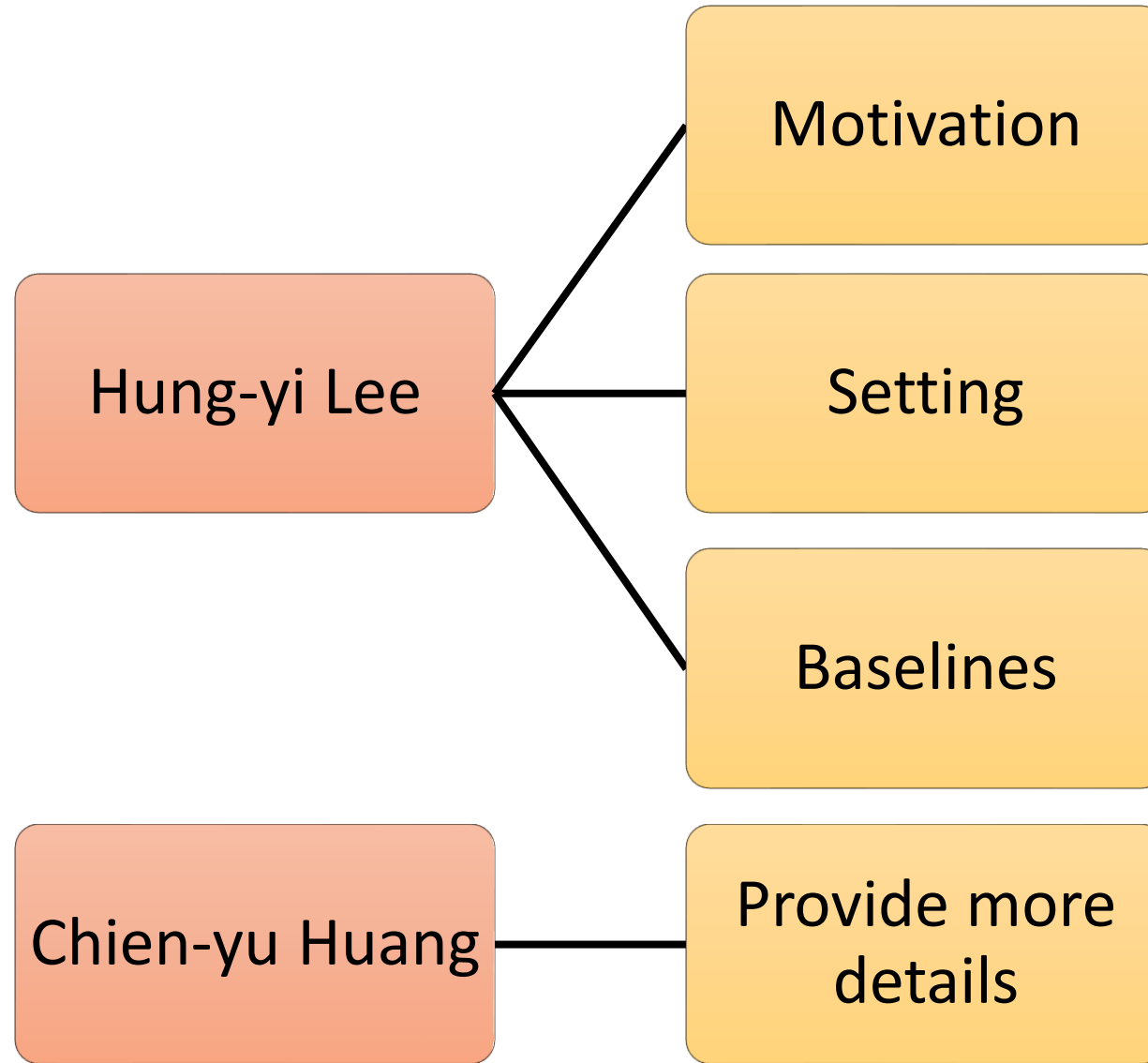
³Mohamed bin Zayed University of Artificial Intelligence, United Arab Emirates

<https://arxiv.org/abs/2309.09510>



Project page: <https://github.com/dynamic-superb/dynamic-superb>

Outline



Chien-yu Huang

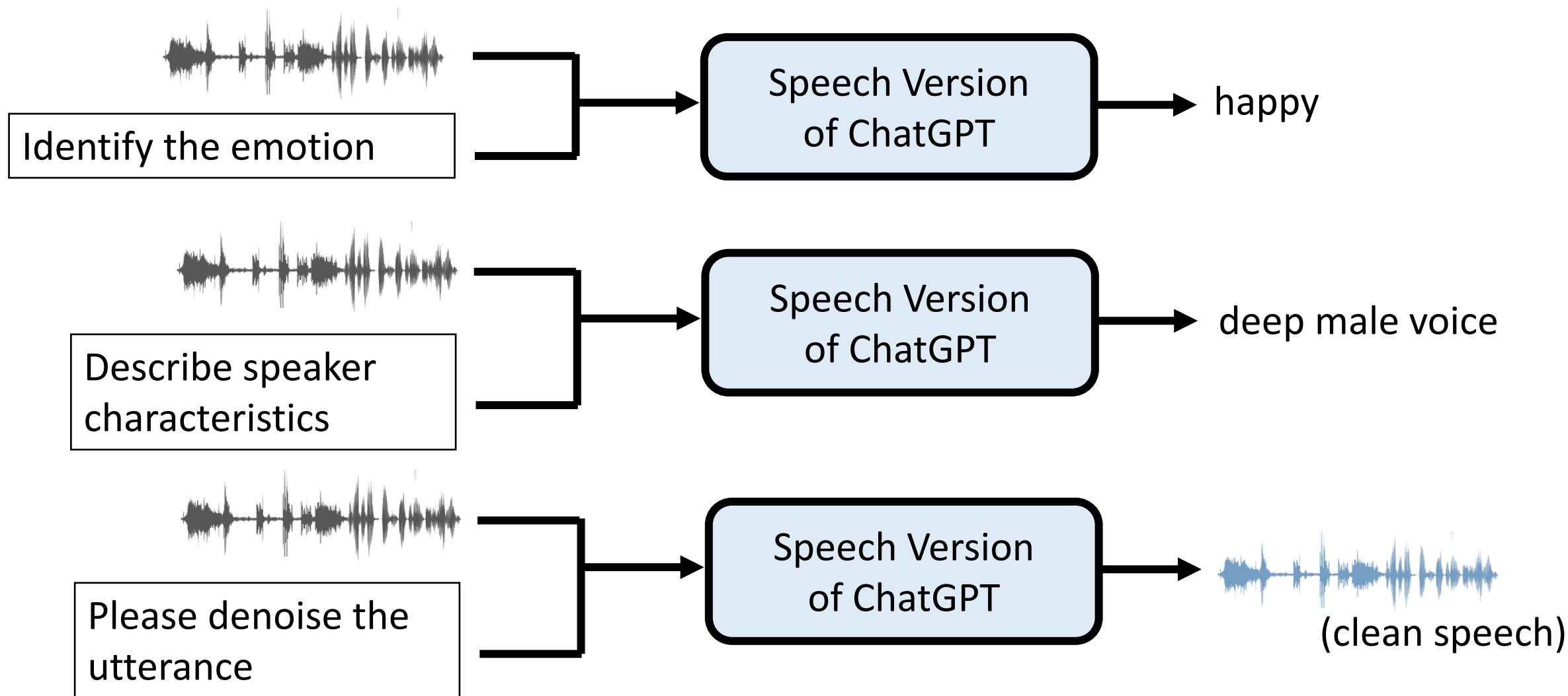
Motivation



Motivation

Evaluating **universal speech models** using **instruction tuning** to perform multiple tasks in a **zero-shot** fashion.

..... Speech Version of ChatGPT



What is still missing?

Speech

SUPERB series: 17 tasks

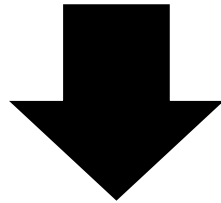
<https://arxiv.org/abs/2105.01051>

<https://arxiv.org/abs/2203.06849>

<https://arxiv.org/abs/2210.08634>

<https://arxiv.org/abs/2210.07185>

<https://arxiv.org/abs/2110.06280>



We need something bigger.

with instruction 😊

NLP

General Language Understanding
Evaluation (GLUE): 9 tasks

<https://arxiv.org/abs/1804.07461>

Super GLUE: 8 tasks

<https://arxiv.org/abs/1905.00537>

FLAN: 62 tasks

<https://arxiv.org/abs/2109.01652>

CrossFit: 160 tasks

<https://arxiv.org/abs/2104.08835>

BIG-bench: 204 tasks

<https://arxiv.org/abs/2206.04615>

natural-instructions: 1616 tasks

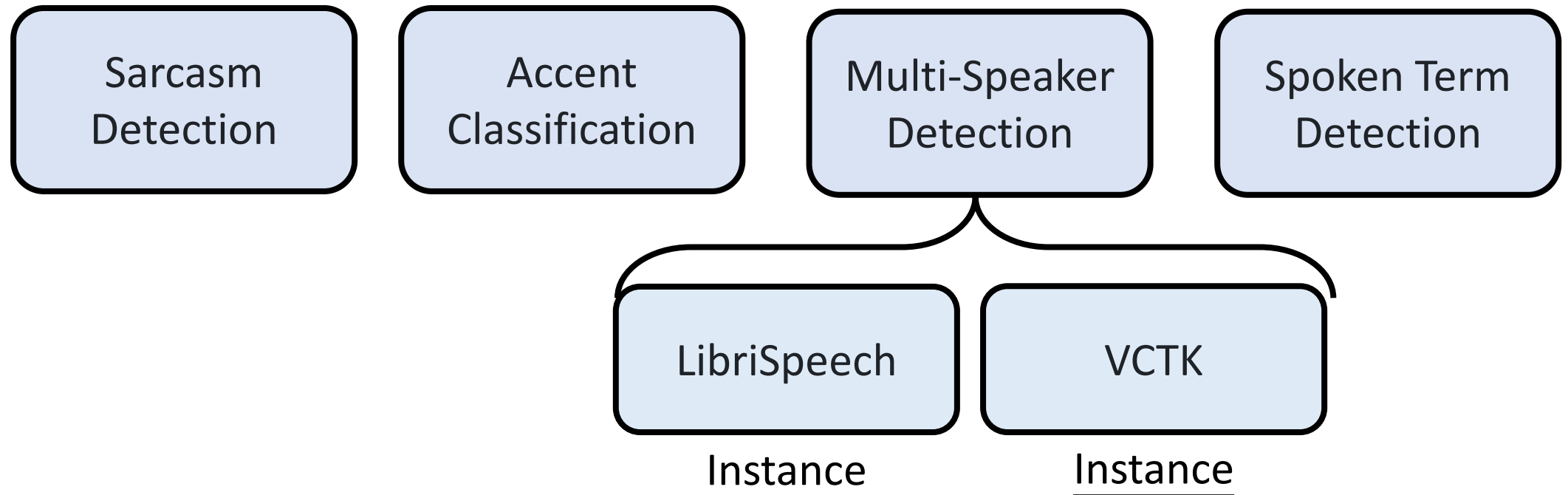
<https://arxiv.org/abs/2204.07705>



Setting

Terminology

Task: a specific type of processing or operation to be carried out



Instance refers to the specific combination of a task and a dataset

Examples in an Instance

Instance

Speaker
Count

LibriTTS

Text Instruction

Audio Input

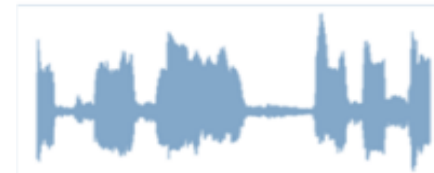
Output

Identify the total number of speakers in the audio. The answer could be one, two, three, four, or five.



one

Determine the number of speakers detected in the audio recording. The answer could be one, two, three, four, or five.



two

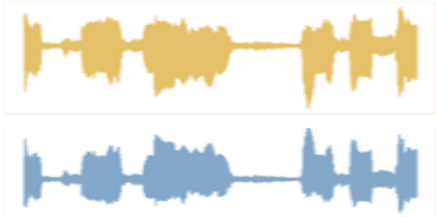
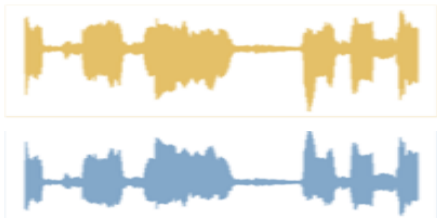
Count the distinct voices present in the audio recording. The answer could be one, two, three, four, or five.



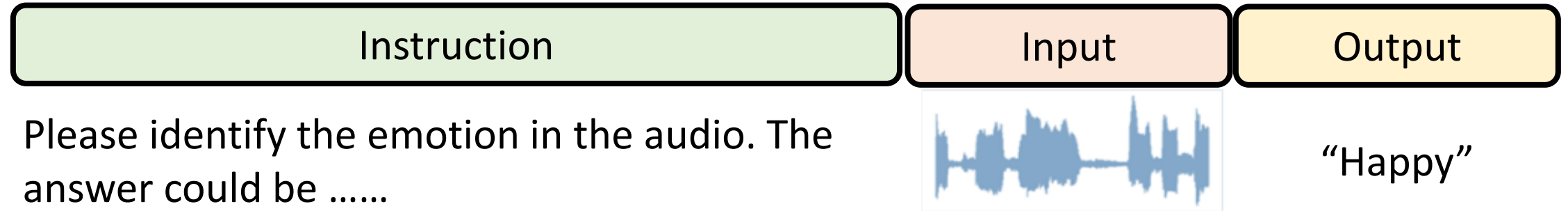
three

Different Example Formats

Speaker Verification

Text Instruction	Audio Input	Output
Listen carefully to both audio recordings and judge if they are spoken by the same individual. The answer could be yes or no.		Yes
Do the speech patterns in the two audio recordings come from the same speaker? The answer could be yes or no.		No

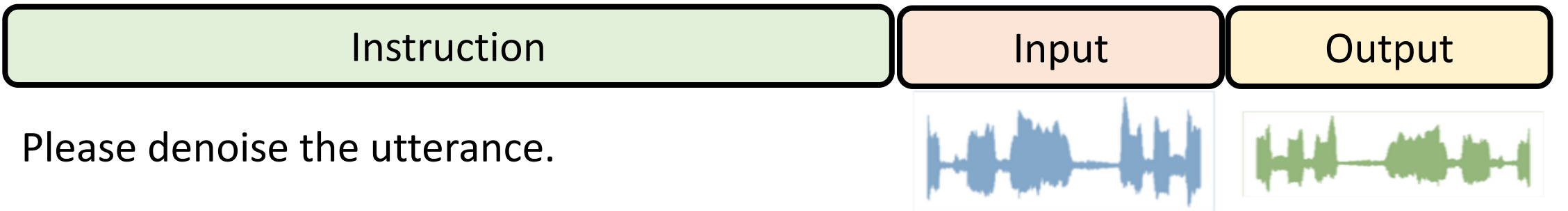
Emotion Recognition



ASR



Speech Enhancement

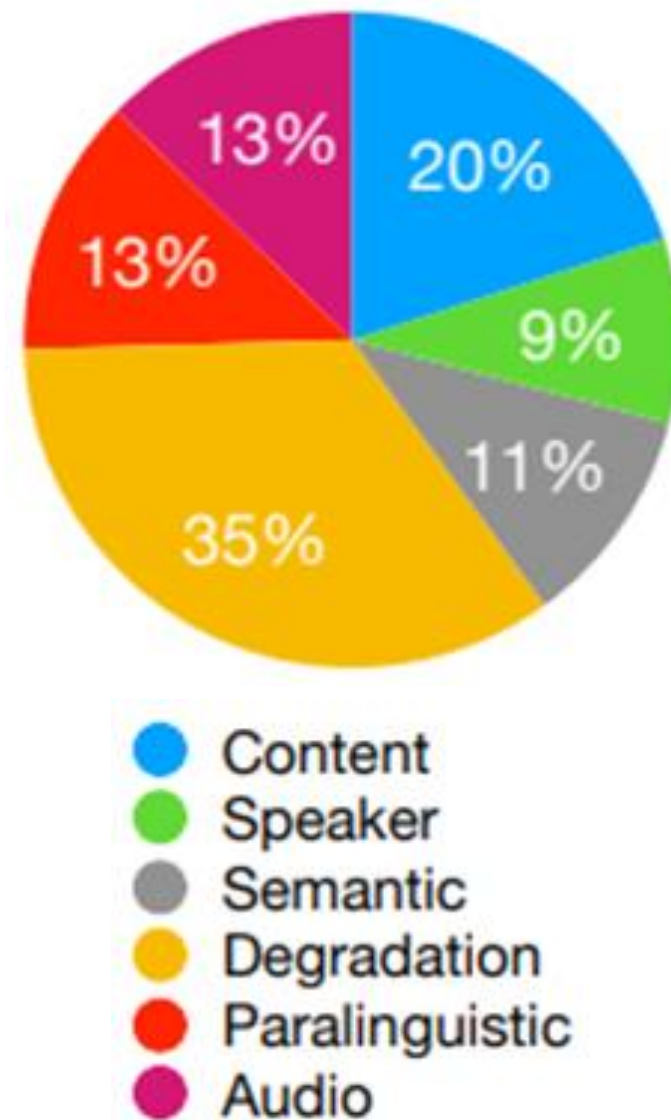


Current Status

- 55 instances
- Covering 6 dimensions
 - Content: speech command recognition
 - Speaker: speaker verification
 - Semantics: sarcasm detection
 - Degradation: noise SNR prediction
 - Paralinguistic: emotion recognition
 - Audio: environmental sound classification

55 is not a big number ...

They are all classification tasks ...



Let's work together!

444 authors across 132 institutions

BEYOND THE IMITATION GAME: QUANTIFYING AND EXTRAPOLATING THE CAPABILITIES OF LANGUAGE MODELS

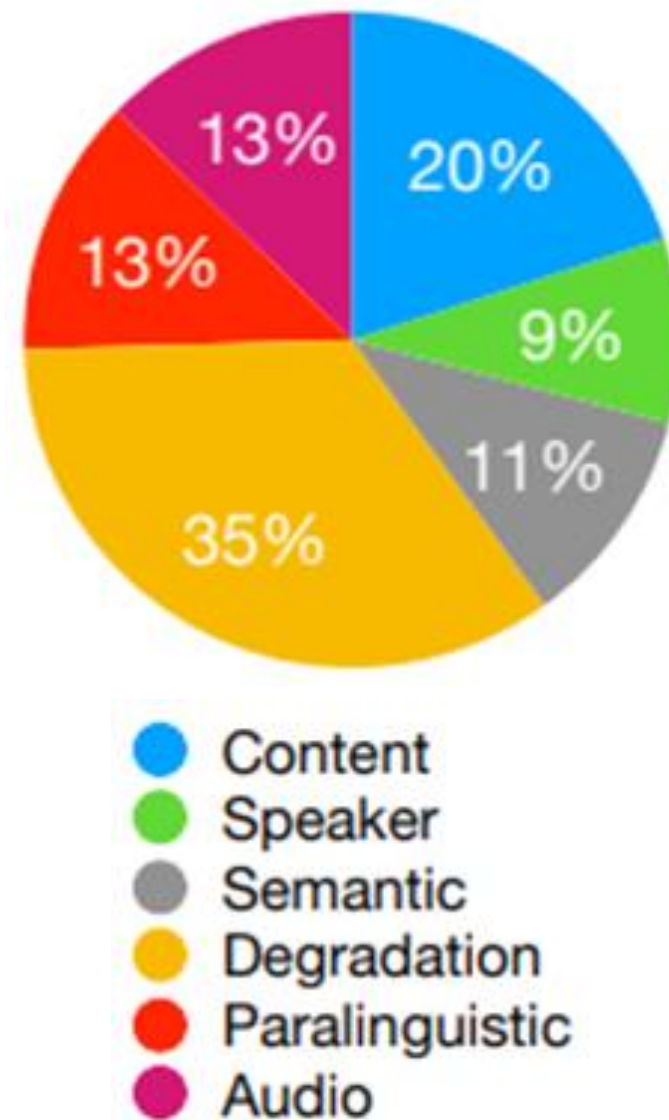
Aashu Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Asad Miri Sheikh, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santos, Aditya Gupta, Adithi Gargya-Alonzo, Agnieszka Kladka, Aitor Lantowicz, Akshar Agarwal, Alonza Power, Alon Ray, Alon Winstach, Alexander W. Kozurek, Ali Salehi, Ali Tazari, Alicia Kiang, Alicia Parrish, Allen Mo, Aman Hussain, Amanda Aschell, Amanda Dsouza, Ambrose Stone, Amotz Rubinfeld, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Smalchitsky, Andrew Dai, Andrew Lai, Andrew Lampinen, Andy Tsao, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gotardt, Annalisa Nonelli, Anna Vukobrat, Arash Ghahramanlou, Arfa Tabassam, Arif Mirza, Aron Kirshenbaum, Asher Mollendorn, Ashish Sabharwal, Austin Herrick, Avia Eilat, Aykut Erdem, Ayta Karakay, B. Ryan Roberts, Bao Sheng Lee, Barret Zoph, Bartholomaj Bojanowski, Benjamin Ouyang, Behnam Heydari, Behnam Neyshabur, Benjamin Peters, Benno Stein, Berk Elmekci, Bill Yuchen Lin, Blake Howald, Cameron Diaz, Cameron Dour, Catherine Stinson, Cedrick Argente, César Ferri Ramalho, Chandan Singh, Charles Rathkopf, Chelsea Wang, Chitra Baral, Chiya Wu, Chir Calhoun-Burch, Chris Walter, Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Sirio, Colin Raffel, Courtney Ashcraft, Cristina Garbacia, Damien Silvio, Dan Gamarnik, Dan Hendrycks, Dan Kilian, Dan Roth, Daniel Freeman, Daniel Khazanchi, Daniel Levy, Daniel Mesquita Gonzalez, Danielle Perczyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Der Gilboa, David Dohan, David Duikand, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleykin, Denis Yuret, Derek Chen, Derek Tan, Desiree Hapkins, Dignata Miara, Dilyar Baran, Dimitri Costin-Mollo, Diyi Yang, Dong-Ho Lee, Ekaterina Shumova, Elia Dognas Gabalt, Elad Segal, Elanor Hagerman, Elizabeth Barnes, Elizabeth Denny, Elie Fuxicik, Emanuele Rodola, Emma Lau, Eric Cha, Eric Tang, Erik Elden, Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Fersk, Ethan Kim, Ethan Engala Manyasi, Evgenij Zhelezniokov, Fanyue Xia, Fanrong Shao, Fernando Martinez-Plumed, Francesca Happel, Francois Chollet, Frieda Ring, Gaurav Mishra, Gema Duato Winans, Gerard de Melo, Germina Krawczyk, Gianbattista Parronandoli, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovich-Lopez, Gregor Betz, Guy Gur-Ari, Hana Galjanovic, Hannah Kim, Hannah Rucklin, Harmanesh Hajipour, Harsh Mehta, Hayden Bogar, Henry Shevlin, Herich Schirwa, Hooman Yekara, Hongming Zhang, Hugh Mo Wang, Ian Ng, Isaac Noble, Jaap Banaet, Jack Geisinger, Jackson Kerrison, Jacob Hilton, Jacobus Lee, Jaime Fernandez Piaz, James B. Simon, James Koppel, James Zheng, James Zou, Jan Kozak, Jan Thompson, Janek Kasper, Janina Radom, Jascha Sobel-Dickman, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, Jelle Boscher, Jennifer Marsh, Jeremy Kim, Jesse Tsal, Jesse Engel, Jengjia Ahbi, Kachang Xu, Jianing Song, Jiliang Tang, Joan Waweru, John Barden, John Miller, John U. Balis, Jonathan Herant, Jori Froberg, Jos Rozen, Jose Hernandez-Castro, Joseph Boudeman, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chan, Kamal Kancher, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Kaya Markov, Kenneth D. Doherty, Kevin Gimpel, Kevin Omadidi, Kyle Mathewson, Krizna Chaitanya, Ksenia Shkara, Ksenia Shkifdar, Kyle McDonnell, Kyle Richardson, Laris Reynolds, Leo Gao, Li Zhang, Lian Dognas, Lianhui Qie, Lidia Contreras-Ochando, Louis-Philippe Morneau, Luca Michalek, Lucas Lau, Lucy Noble, Ludvig Schenck, Ludwig He, Luis Oliveros Cubela, Luke Metz, Lutz Kervin Spel, Maarten Bosman, Maarten Sap, Maarten ter Hoeve, Mahesh Paragol, Manal Fawzi, Mantas Marozika, Marco Baran, Marco Marrelli, Marco Moro, Maria Jose Ramirez Quijano, Marie Tilkiche, Mario Givanelli, Martha Lewis, Martin Pothner, Matthew L. Lawrin, Matthias Hagen, Maylis Schubert, Medina Orduña Rautavaara, Melody Arnold, Melvin McElrath, Michael A. Yeo, Michael Cohen, Michael Gu, Michael Bontickly, Michael Starin, Michael Strube, Michal Smojkowski, Michele Bevilacqua, Michihito Yamanaka, Mihir Kale, Mike Cain, Minseu Kim, Minseu Seung, Mohit Bansal, Moira Aminzadeh, Mor Geva, Mostafid Ghafar, Mukund Varma T, Nanyang Peng, Nathan Chi, Nguyen Lau, Neta Gur-Ari Kralavetz, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nikita Nangia, Niklas Deckert, Niklas Muenighoff, Nirish Shitrik Kedar, Nivethita S. Iyer, Noah Constant, Noah Field, Nuan Wen, Oliver Zhang, Omar Aggar, Omar Elbaghdadi, Omar Levy, Oran Etan, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Wang, Paul Fu Liang, Paul Vicol, Piyush Alpoornalabadi, Polyna Liao, Potry Liang, Peter Chang, Peter Eckersley, Phu Hoa Huu Phung Huang, Piotr Mikulowski, Piyush Puri, Pooja Prasadipour, Priti Oli, Qianbo Mei, Qing Lyu, Qinfeng Chen, Robin Ranjale, Rachel Ema Radolph, Raefar Gabriel, Rafael Haberker, Ramiro Rizzo Delgado, Raphael Meltzer, Rhytham Gang, Richard Barnes, Rishi A. Santos, Riku Arakawa, Robbe Raymondker, Robert Frank, Rohan Sikand, Roman Novak, Roman Smolov, Roman Lofitus, Rosanna Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Tsohan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajjan Aam, Sam Dillavoss, Sam Shitler, Sam Wiseman, Samuel Gaurier, Samuel R. Bowman, Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rose, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hashemi, Sharon Zhou, Shaohui Srivastava, Sherry Shi, Shikhar Singh, Shima Asadi, Shizhang Shao, Shihui Pichler, Shihuan Tschaiwald, Shyam Upadhyay, Shyamkumar (Shamun) Debnath, Siamak Shakeri, Simon Thormeyer, Simona Melzi, Siva Reddy, Soha Pitschla Makini, Soa-Hyun Lee, Spencer Torres, Sriharsha Harvar, Stanislas Dehaene, Stefan Diehl, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Piantedosi, Stuart M. Shieber, Summer Mishra, Svetlana Kirichenko, Swapnil Mishra, Tai Linzon, Tai Schuster, Tao Li, Tao Yu, Tarik Ali, Tama Hashimoto, Te-Lin Wu, Telle Enderbich, Theodore Rutherford, Thomas Phan, Tianle Wang, Tobias Munkitell, Timo Schick, Timofei Kornev, Timothy Telleen-Lawton, Tim Tondary, Tobias Griesenber, Truong-Chan, Trishala Norraj, Tushar Khot, Tyler Shultz, Uri Shalit, Vidhan Mishra, Vera Demberg, Verónica Nyman, Vikas Rastaki, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Shrikumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tang, Xinran Zhao, Xinyi Wu, Xudong Shen, Yafolish Yaghoobzadeh, Yair Laksyetz, Yanguo Song, Yassman Bahet, Yefim Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hua, Yufang Hou, Yuntao Bai, Zachary Seld, Zhequn Zhao, Zijian Wang, Zijie J. Wang, Ziri Wang, Ziyi Wu.

Current Status

- 55 tasks created from 33 datasets
- Covering 6 dimensions
 - Content: speech command recognition
 - Speaker: speaker verification
 - Semantics: sarcasm detection
 - Degradation: noise SNR prediction
 - Paralinguistic: emotion recognition
 - Audio: environmental sound classification

Everyone can add new tasks! ➡ Dynamic

(We can write a big paper together like BIG-bench in the future.)

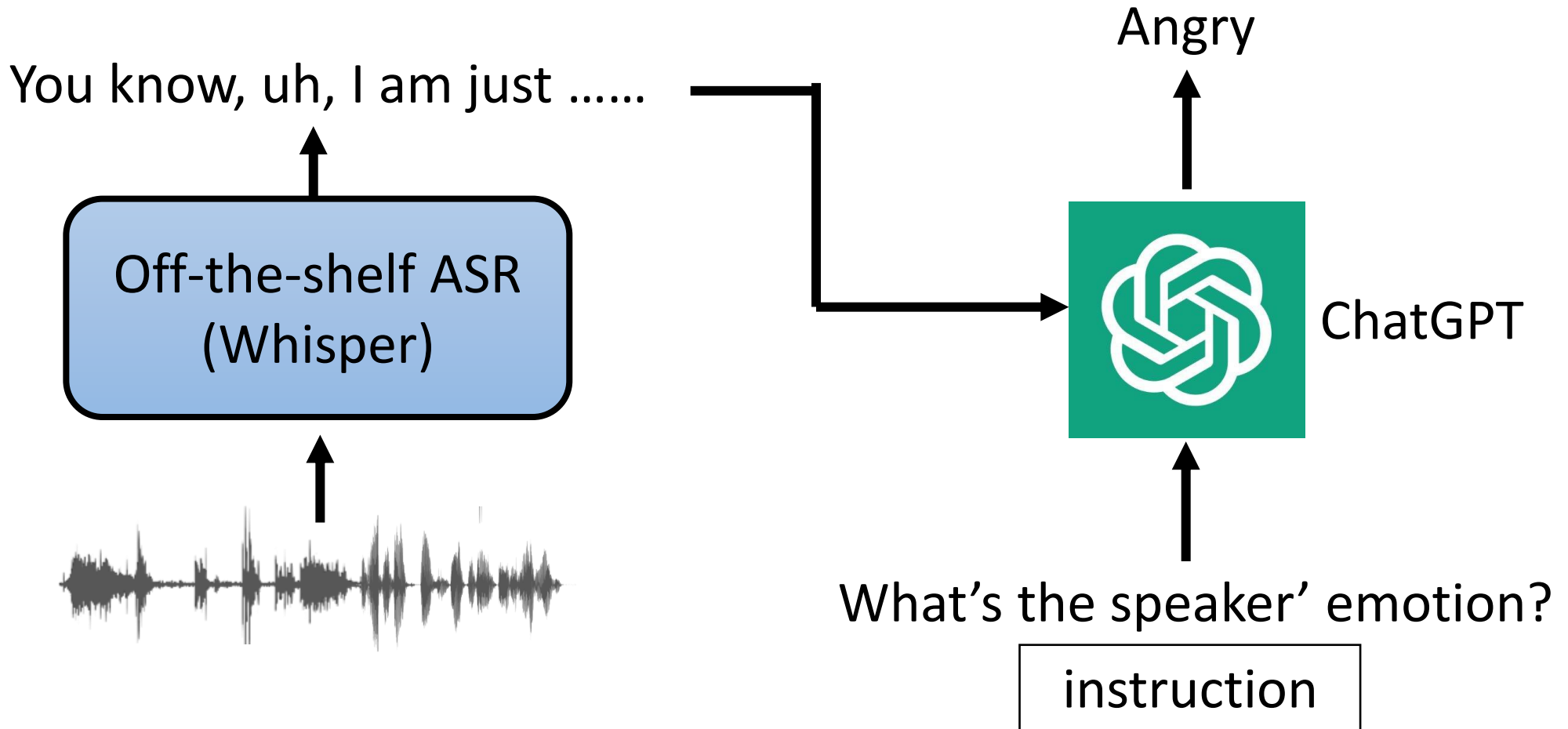


Let's work together!

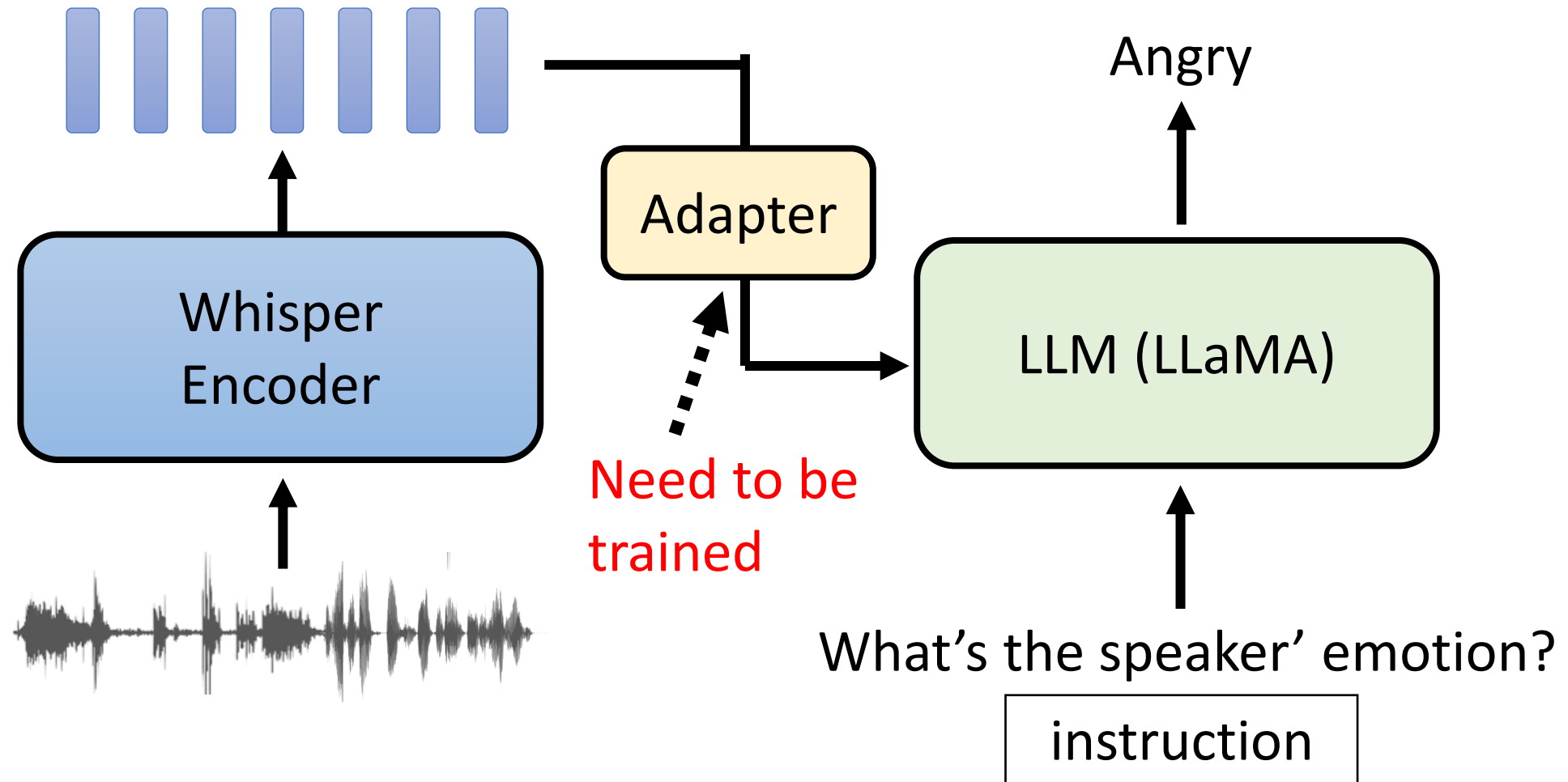
Baselines



Example Baseline: ASR + ChatGPT

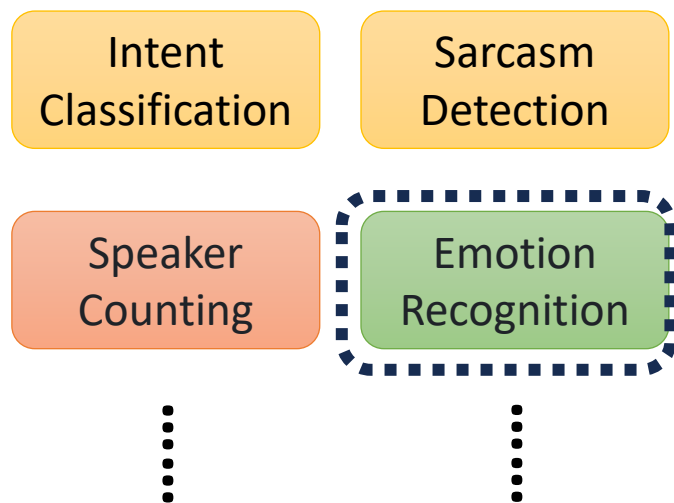


Example Baseline: Whisper Encoder + LLM



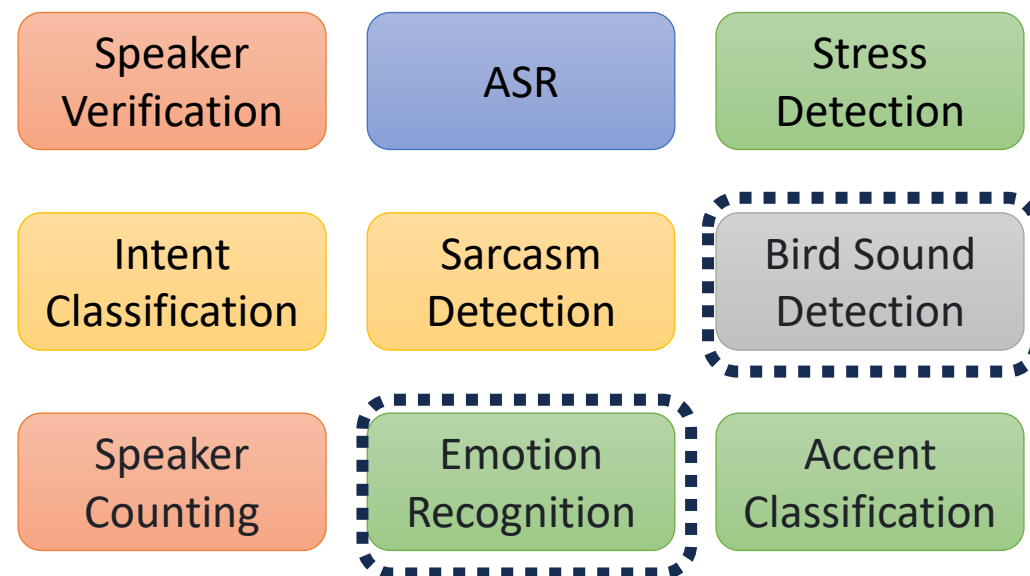
Training Data vs. Testing Data

Training Data (23 instances)



**Training model
parameters**

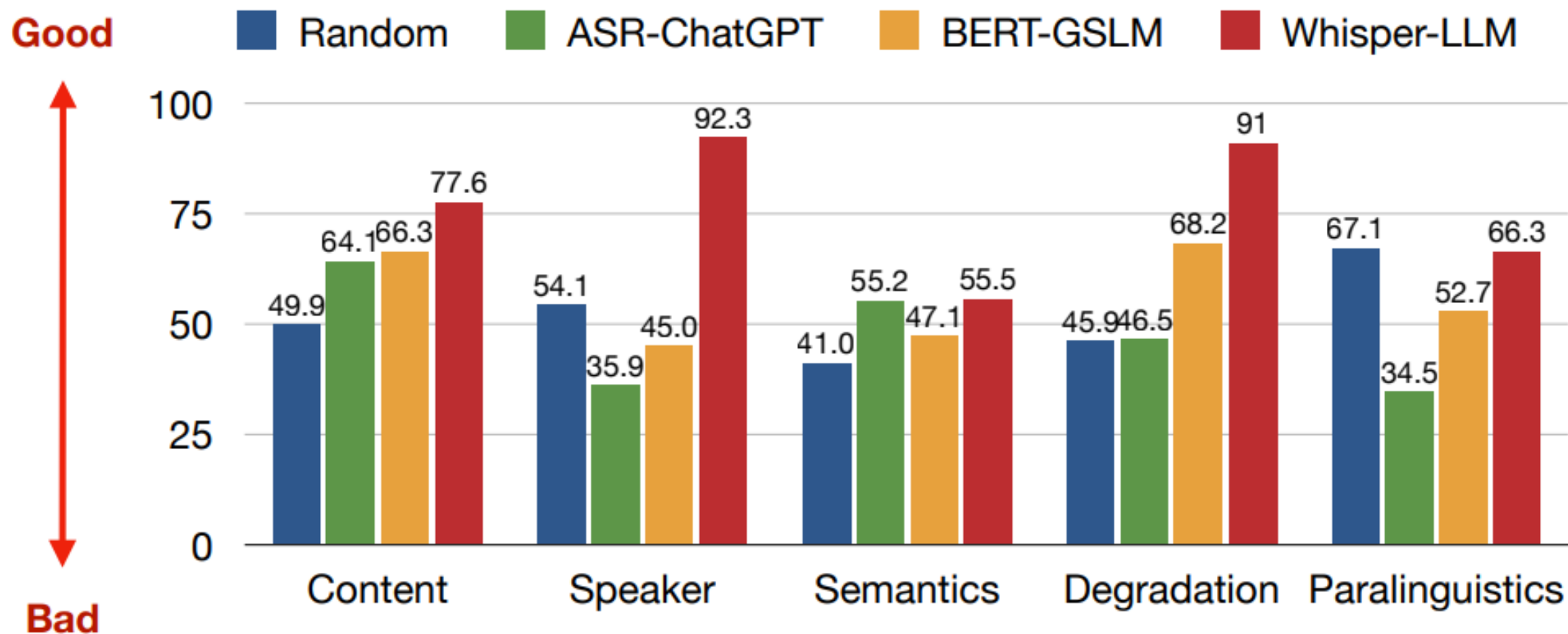
Dynamic-SUPERB (55 instances)



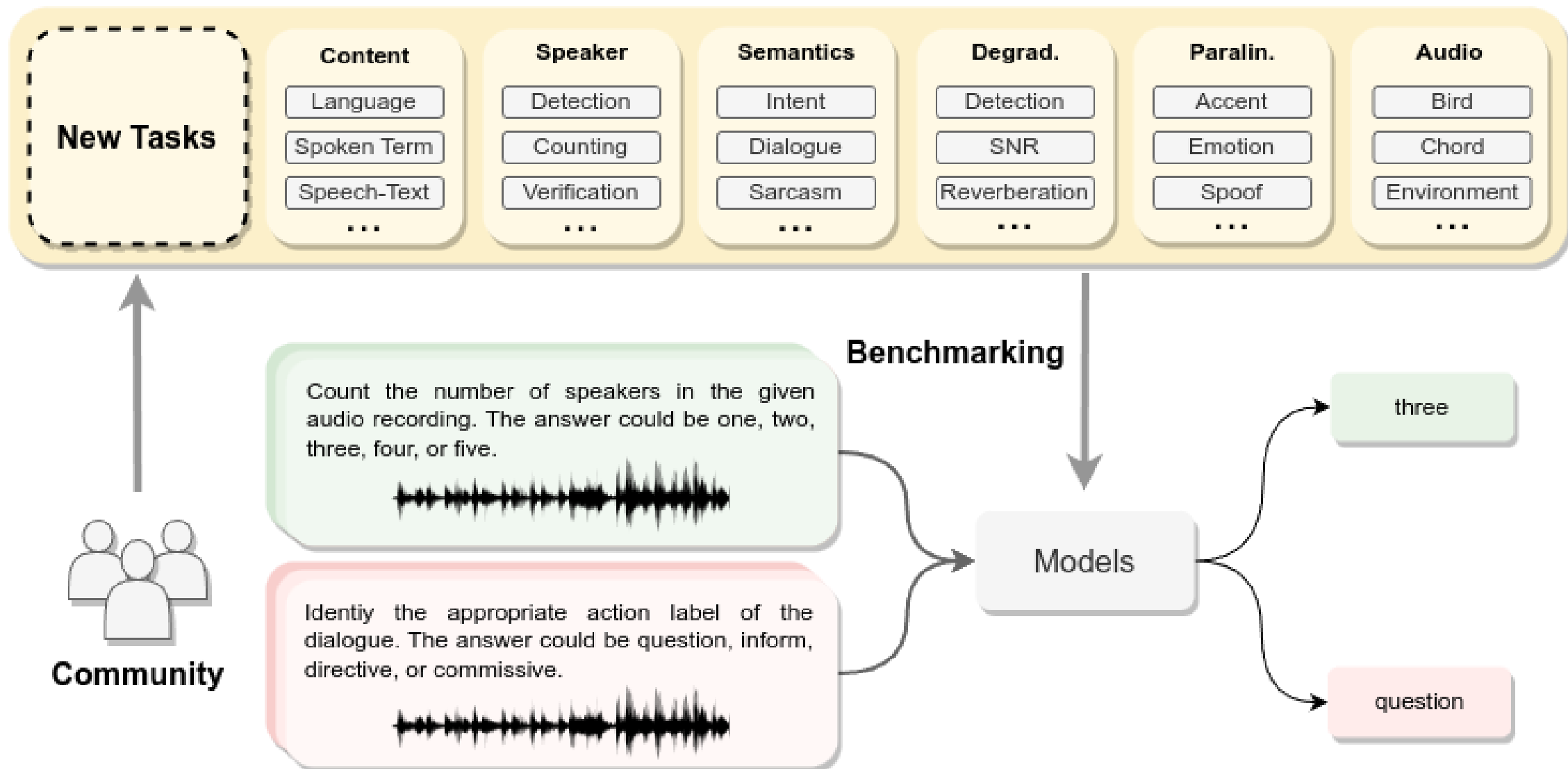
Seen
(not the
same data)

Unseen

Overall Results



Dynamic-SUPERB



Action Item

- ASRU 2023 satellite workshop - **Speech foundation models and their performance benchmarks** (SPARKS)
- Submit a paper to Dynamic-SUPERB session
 - Short paper: 2-3 pages, detailing your methodology, findings, in-depth analysis or new tasks with the Dynamic-SUPERB benchmark.
 - Deadline: **Nov. 24, 2023** (not the workshop deadline)
 - Present at the Dynamic-SUPERB session

[https://sites.google.com/
g.ntu.edu.tw/sparks](https://sites.google.com/g.ntu.edu.tw/sparks)

