# The Audio-Visual Conversational Graph: From an Egocentric-Exocentric Perspective

Wenqi Jia[1,2], Miao Liu[4], Hao Jiang,[2], Ishwarya Ananthabhotla[2], James M. Rehg[3], Vamsi Krishna Ithapu[2], Ruohan Gao[2]

[1] Georgia Institute of Technology, [2] Meta Reality Labs Research, [3] University of Illinois Urbana-Champaign, [4] GenAI, Meta

CVPR SEATTLE, WA JUNE 17-21, 2024

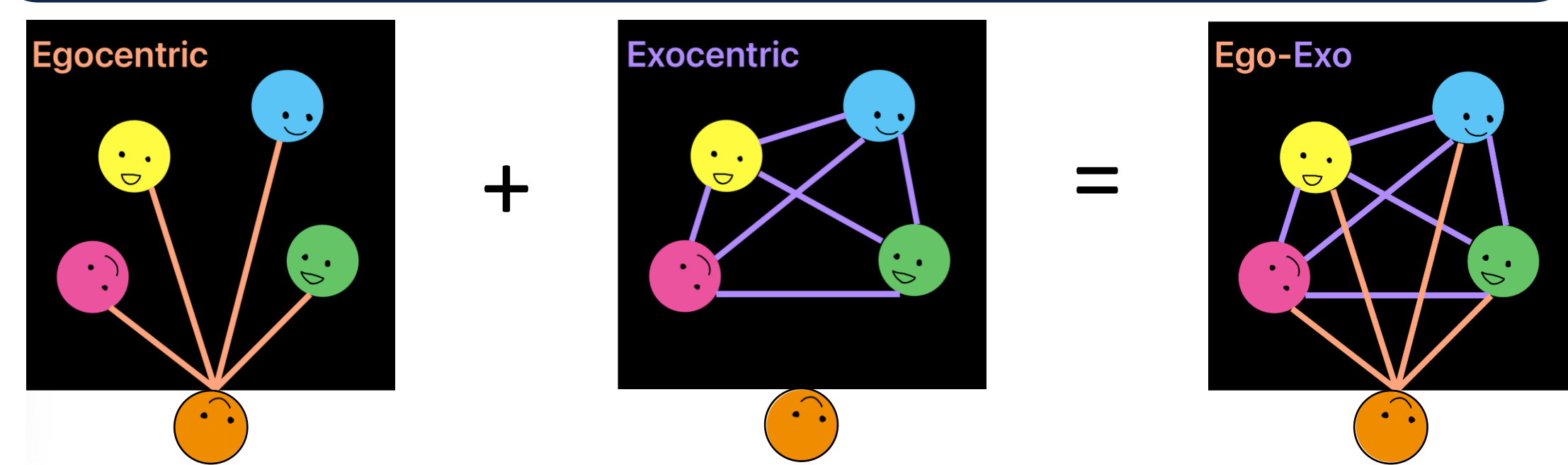*Find more on our project page!*

## Motivation

***Concurrent conversations are common in life***

➤ Could be noisy and ambiguous

➤ Capturing social states of participants helps decide which sound source to enhance for whom

➤ Facilitate effective and efficient communication

## Ego-Exo Conversational Graph



Egocentric Behavior | Exocentric Behavior

Camera Wearer as Observer (Ego)

*Humans can understand both **Egocentric** and **Exocentric** conversational behaviors*

Egocentric + Exocentric = Ego-Exo

## Ego-Exo Directional Edge

➤ For each pair of nodes $(c, p_j)$ or $(p_i, p_j)$, we aim to determine:
  - If they are **Speaking To** ($S$) each other
  - If they are **Listening To** ($L$) each other
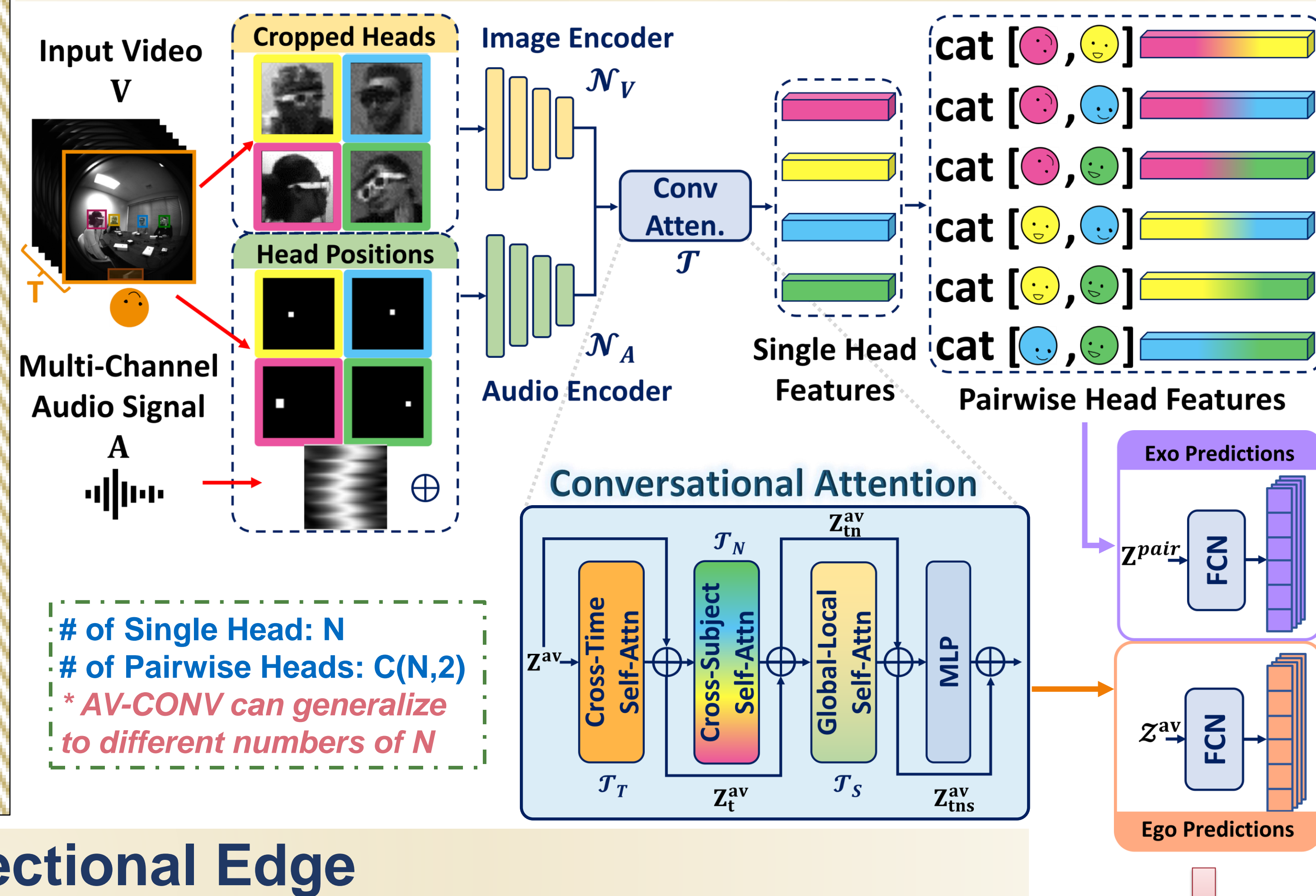➤ Results in four attributes:
  - For each **Egocentric** Edge:
    $$e^S_{c \to p_i} \quad e^S_{p_i \to c} \quad e^L_{c \to p_i} \quad e^L_{p_i \to c}$$
  - For each **Exocentric** Edge:
    $$e^S_{p_i \to p_j} \quad e^S_{p_j \to p_i} \quad e^L_{p_i \to p_j} \quad e^L_{p_j \to p_i}$$

$p_j$: Subject $j$
$p_i$: Subject $i$
$c$: Camera Wearer

$annot(c, p_i) \to$ **Ego Edge**
Is $c$ **Speaking** to $p_i$?
Is $c$ **Listening** to $p_i$?
Is $p_i$ **Speaking** to $c$?
Is $p_i$ **Listening** to $c$?

$annot(p_i, p_j) \to$ **Exo Edge**
Is $p_i$ **Speaking** to $p_j$?
Is $p_i$ **Listening** to $p_j$?
Is $p_j$ **Speaking** to $p_i$?
Is $p_j$ **Listening** to $p_i$?

## Ego-Exocentric Conversational Graph Prediction

the *first* to explore Exocentric conversational interactions from Egocentric videos

✓ Jointly modeling talking *and* listening behaviors

✓ Jointly modeling Egocentric *and* Exocentric behaviors as graph

## Method



Input Video V
Cropped Heads
Head Positions
Image Encoder $\mathcal{N}_V$
Multi-Channel Audio Signal A
Audio Encoder $\mathcal{N}_A$
Conv Atten. $\mathcal{T}$
Single Head Features
Pairwise Head Features
cat [⬤, ⬤]

\# of Single Head: N
\# of Pairwise Heads: C(N,2)
* AV-CONV can generalize to different numbers of N

### Conversational Attention

$Z^{av}$ — Cross-Time Self-Attn $\mathcal{T}_T$ — $Z^{av}_t$ — Cross-Subject Self-Attn $\mathcal{T}_N$ — $Z^{av}_{tn}$ — Global-Local Self-Attn $\mathcal{T}_S$ — $Z^{av}_{tns}$ — MLP

Exo Predictions $Z^{pair}$ FCN
Ego Predictions $Z^{av}$ FCN

Output Conversational Graph
Speaking to → Listening to
Exo Edges
Ego Edges
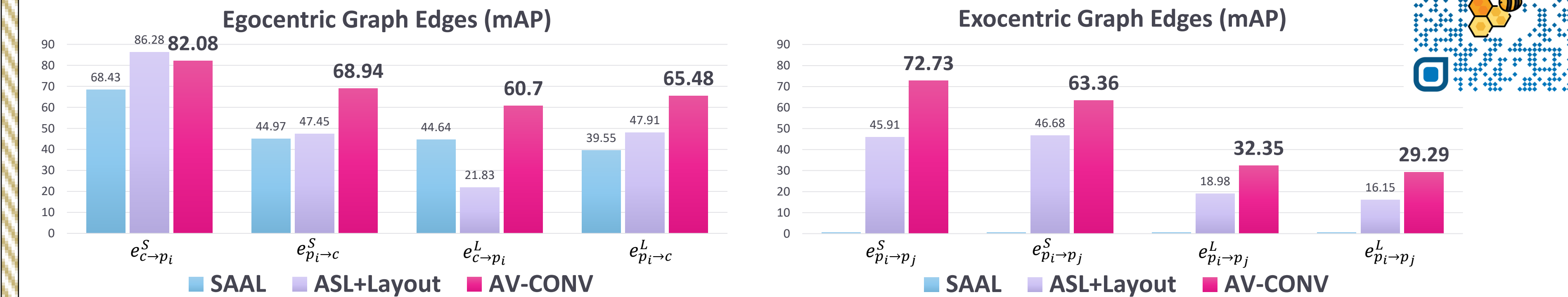
## Experiments and Results

➤ **Dataset:** Egocentric Concurrent Conversations Dataset (15,682/6,329 Train/test)

➤ **Baselines:** 1. SAAL (Ryan 2023)

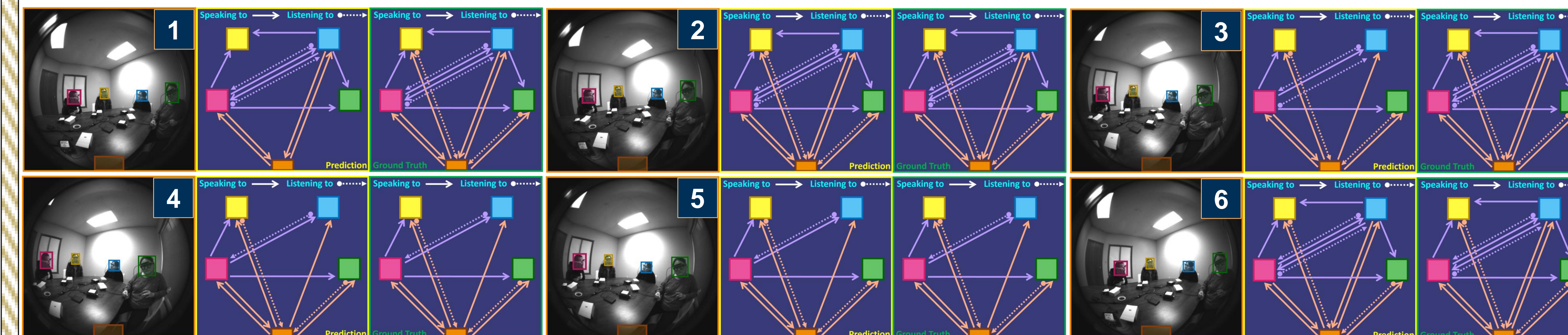2. Active Speaker Localization (Jiang 2022) + 3D person layout estimation



Egocentric Graph Edges (mAP) | Exocentric Graph Edges (mAP)

**Our AV-CONV consistently outperforms both baselines across all subtasks**

| | Egocentric Graph | | | | | Exocentric Graph | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | $e^S_{c \to p_i}$ | $e^S_{p_i \to c}$ | $e^L_{c \to p_i}$ | $e^L_{p_i \to c}$ | Ego Avg | $e^S_{p_i \to p_j}$ | $e^S_{p_j \to p_i}$ | $e^L_{p_i \to p_j}$ | $e^L_{p_j \to p_i}$ | Exo Avg |
| HEAD ONLY | 51.20 | 51.65 | 37.19 | 29.38 | 42.36 | 54.52 | 48.12 | 16.48 | 17.33 | 34.11 |
| AUDIO ONLY | 84.32 | 53.43 | 22.94 | 24.26 | 46.24 | 51.63 | 43.89 | 14.17 | 15.58 | 31.32 |
| MASK ONLY | 54.55 | 52.18 | 39.27 | 33.54 | 44.89 | 55.00 | 47.29 | 14.93 | 16.09 | 33.33 |
| HEAD+MASK | 47.84 | 50.28 | 35.80 | 22.38 | 39.08 | 52.85 | 45.90 | 14.83 | 15.89 | 32.37 |
| AUDIO+MASK | 45.83 | 47.40 | 22.83 | 21.31 | 34.34 | 50.40 | 43.86 | 14.76 | 15.95 | 31.24 |
| AV-CONV | 82.08 | **68.94** | **60.70** | **65.48** | **69.30** | **72.73** | **63.36** | **32.35** | **29.29** | **49.43** |

subject-specified → HEAD ONLY
global → AUDIO ONLY
subject-specified → MASK ONLY
global-local, subject-specified → AV-CONV

The combination of <u>subject-specified visual cues</u>, <u>global audio information</u>, and <u>spatial context through the positional mask</u> is essential for accurate prediction

## Visualization

➤ **Conversational Dynamics:** 6 frames with a temporal stride of 15, ~3 seconds



➤ **Future Work:** 1) extend our framework to other social behavior; 2) study more complex social relationships such as conversation groups' mobility

*This work was primarily done during an internship at Meta Reality Labs.*