

# ACUITY: CREATING REALISTIC DIGITAL TWINS THROUGH MULTI-RESOLUTION POINTCLOUD PROCESSING AND AUDIOVISUAL SENSOR FUSION

Jason Wu, Ziqi Wang, Ankur Sarker, Mani Srivastava (UCLA)



# 01

## INTRODUCTION

HOW TO REPRESENT HUMAN  
SUBJECTS AUDIOVISUALLY AT  
HIGH FIDELITY IN REAL TIME?



+

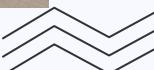
# VISUAL DOMAIN GOALS



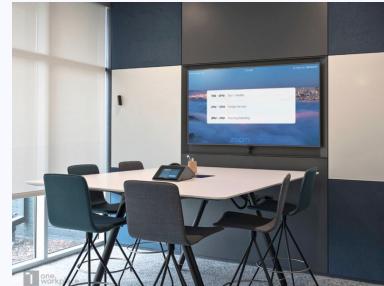
Point Clouds



Avatars



# VISUAL DOMAIN GOALS



Obtain High Resolution  
Point Clouds

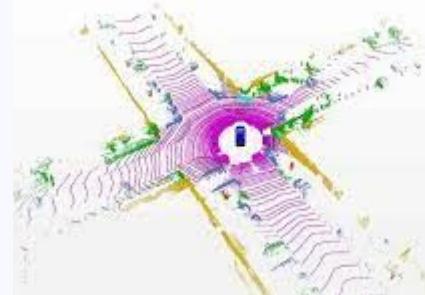
Isolate point clouds of  
human subjects

Place within  
virtual scene

# VISUAL: CURRENT CHALLENGES

+

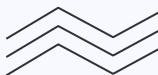
- Existing neural networks trained on sparse point clouds input from KITTI Dataset
- Suffer from runtime issues
- Generalize poorly to dense point cloud input



KITTI Point Cloud (Sparse)

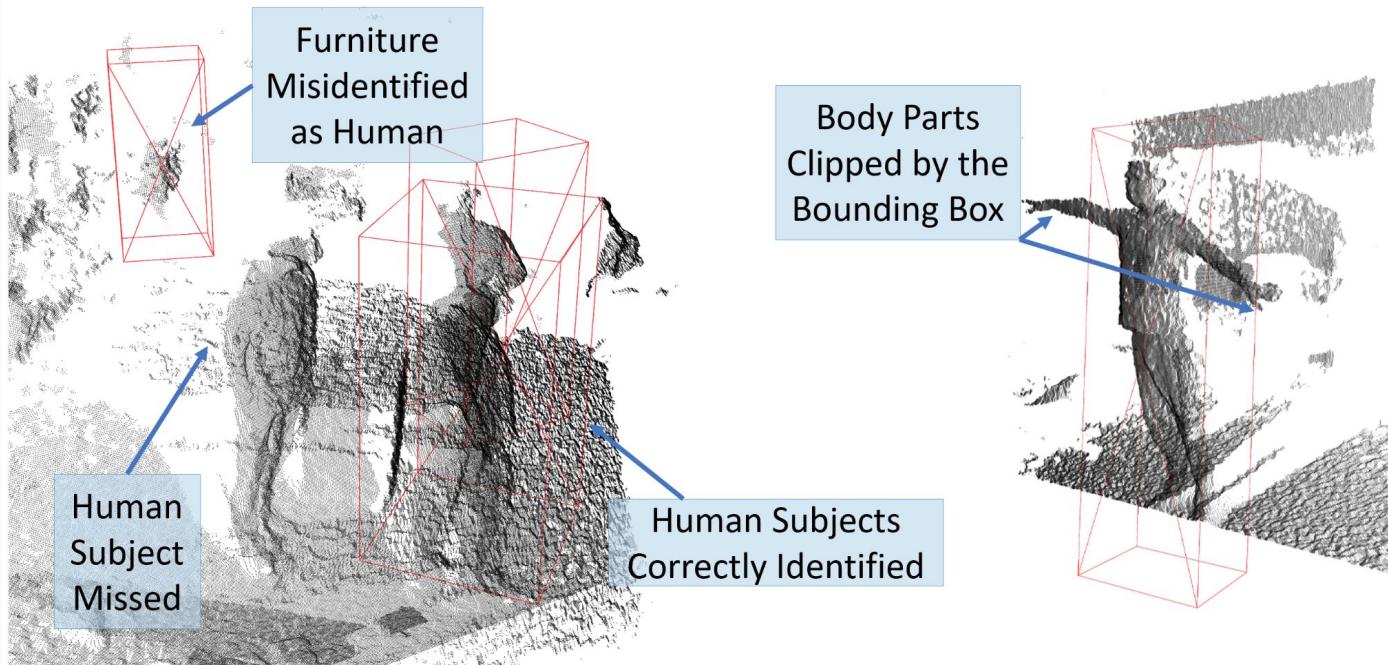


Realsense Point Cloud (Dense)

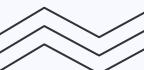


# VISUAL: CURRENT CHALLENGES

+



Example: Failure Cases of NN-based Human Subject Bounding Box Detection



---



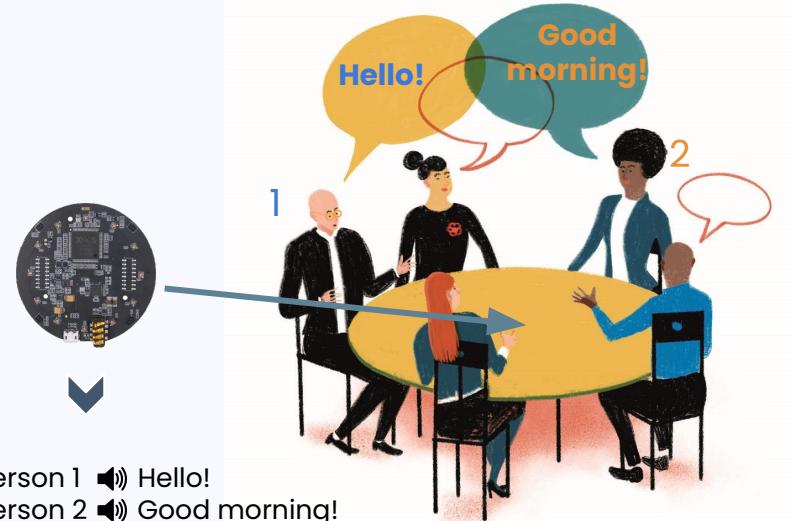
# CHALLENGE

**How do we isolated human  
subjects from dense point  
clouds in real time?**



# AUDIO DOMAIN GOALS

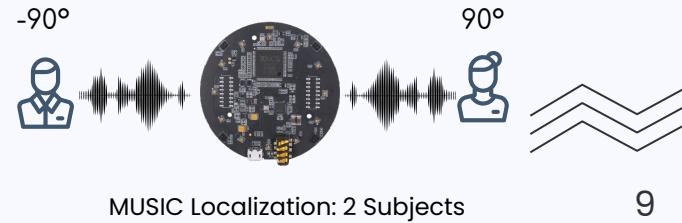
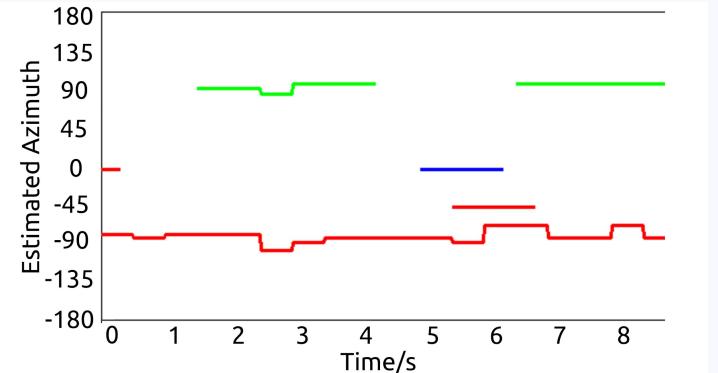
- Audio quality degrades with high background noise or background speech
- **Sound Source Separation:** Leverage beamforming with microphone array to isolate sound from particular direction
- Requires the angle of arrival (AoA) to be known



# AUDIO: CURRENT CHALLENGES

+

- AoA Estimation Methods: SRP-PHAT and MUSIC
- Learning based methods: We require real-time, mobile subjects, and variable number of subjects

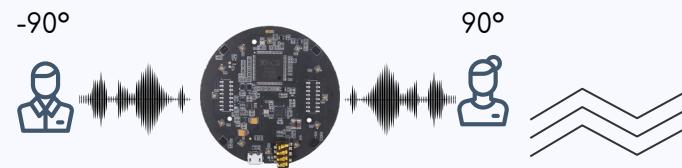
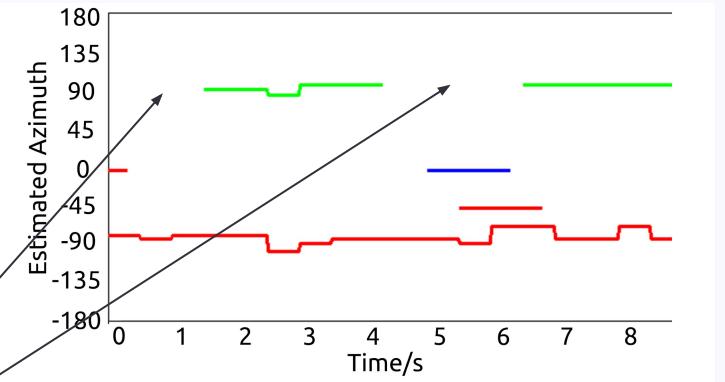


# AUDIO: CURRENT CHALLENGES

+

- AoA Estimation Methods: SRP-PHAT and MUSIC
- Learning based methods: We require real-time, mobile subjects, and variable number of subjects

Missing Subject

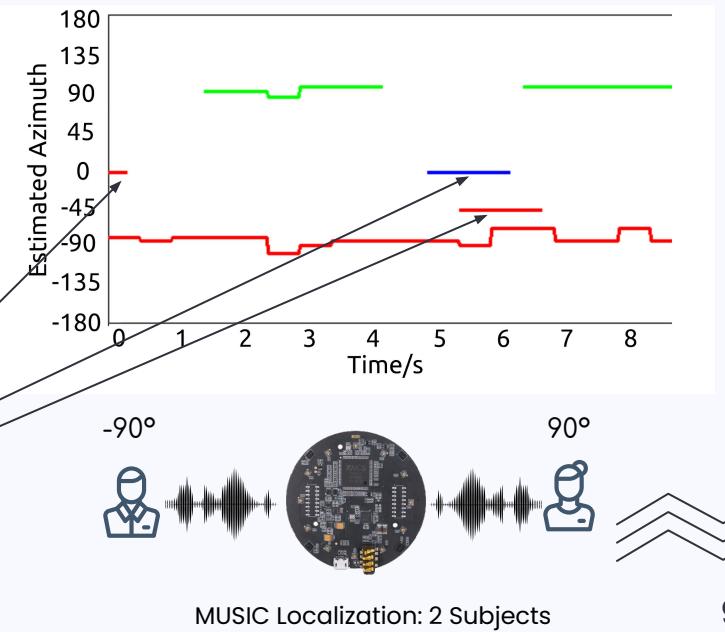


# AUDIO: CURRENT CHALLENGES

+

- AoA Estimation Methods: SRP-PHAT and MUSIC
- Learning based methods: We require real-time, mobile subjects, and variable number of subjects

“Ghost” Subjects



---



# CHALLENGE

**How do we obtain the AoA in real time?**



# ACUITY



## DOUBLE BACKGROUND SUBTRACTION PIPELINE

- Leverage **background subtraction** at two resolutions to efficiently isolate point clouds
- Runs at 30 fps with resolution of 640x480 and three subjects in the scene



## MULTIMODAL FUSION WITH VISUAL LOCALIZATION

- Leverages localization information from the visual pipeline to obtain AoA (**audiovisual fusion**)
- Uses the centroid of the isolated point cloud
- Latency of 30 ms (< 45 ms)



---

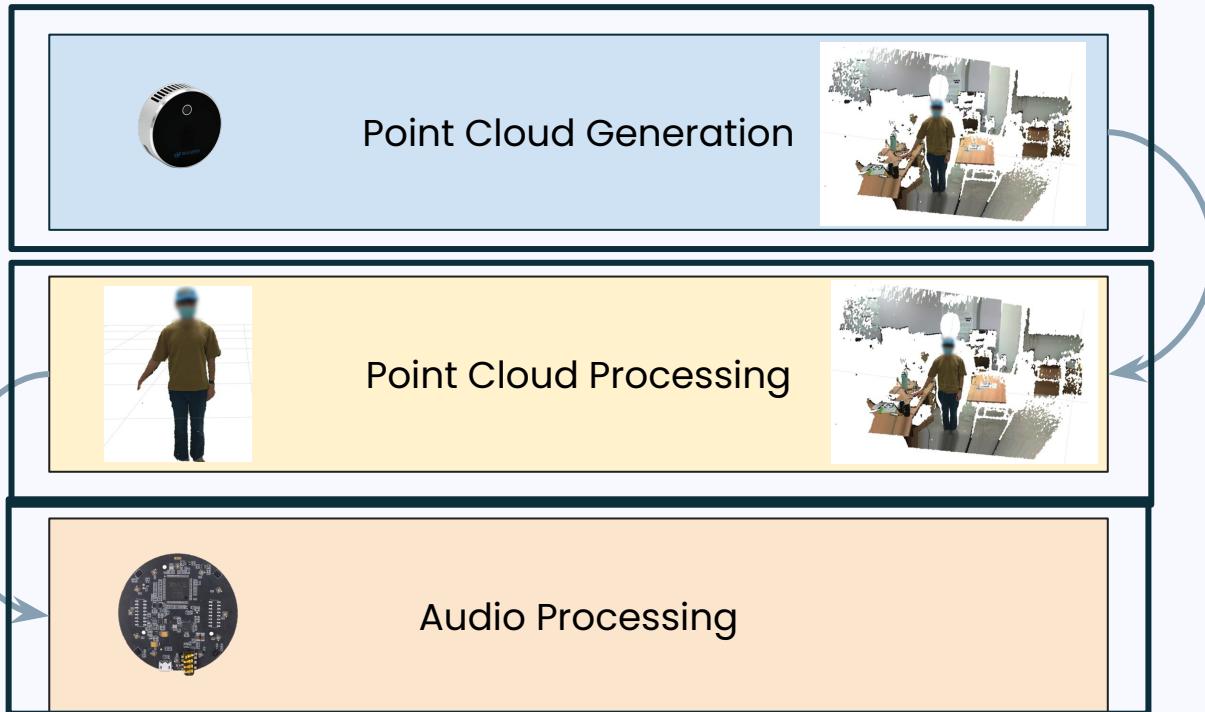
02

# SYSTEM DESIGN AND IMPLEMENTATION



+

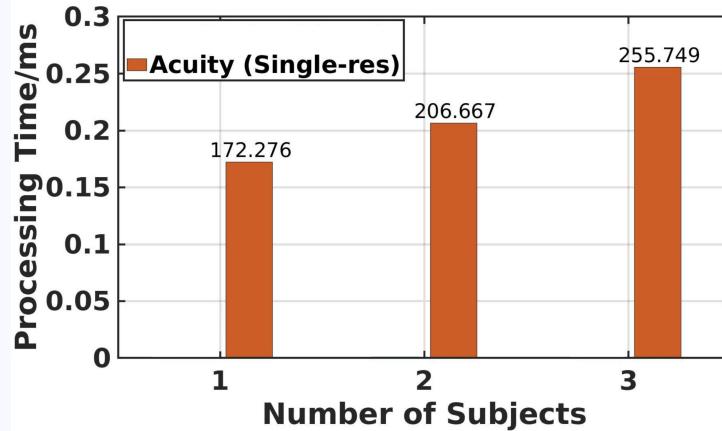
# ACUITY PIPELINE



# POINT CLOUD PROCESSING

+

- **Background subtraction:** Compare each new frame with reference frame without subjects. Removes all the voxels (points) that are identical
- Apply clustering to remove noise



# CHALLENGE

Background Subtraction  
and Clustering do not run  
in real time



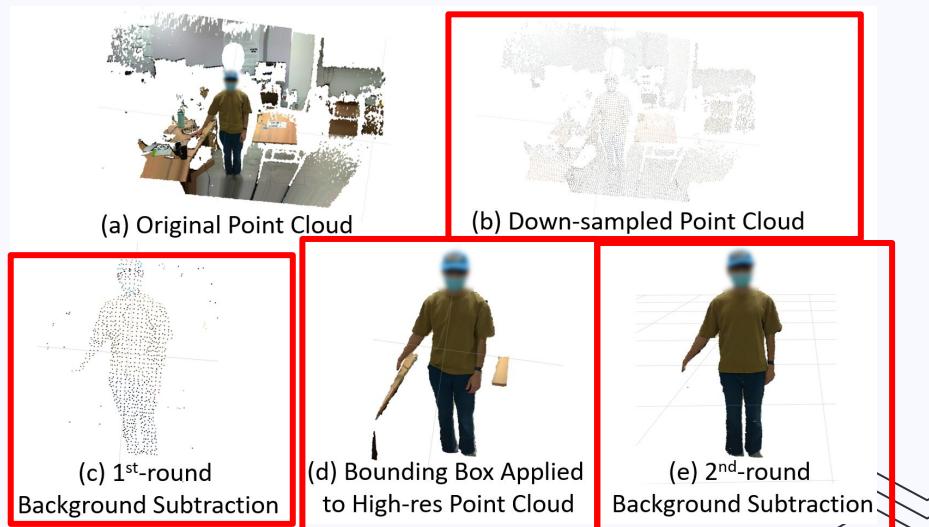
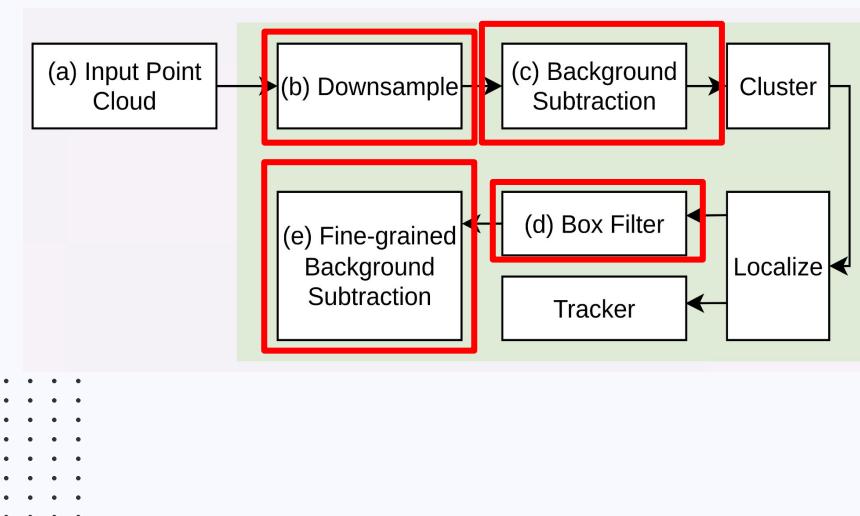
# SOLUTION

Utilize **multi-resolution**  
processing!



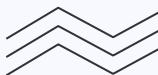
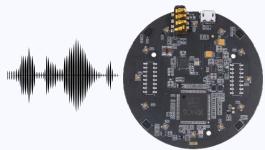
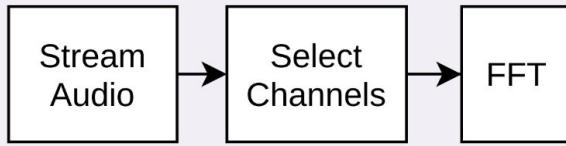
# MULTI RESOLUTION PROCESSING

+



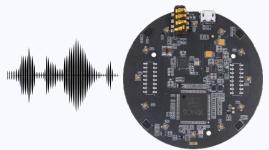
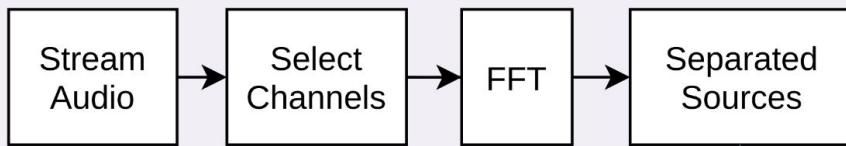
# AUDIO PROCESSING

## Audio Processing

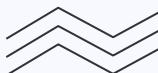


# AUDIO PROCESSING

## Audio Processing

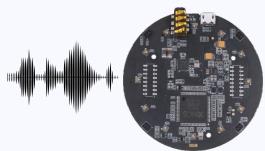
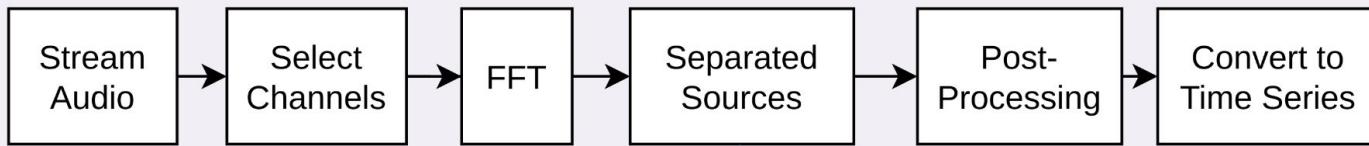


Human Subject  
Centroids Information  
(From Vision Pipeline)

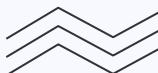


# AUDIO PROCESSING

## Audio Processing



Human Subject  
Centroids Information  
(From Vision Pipeline)



---

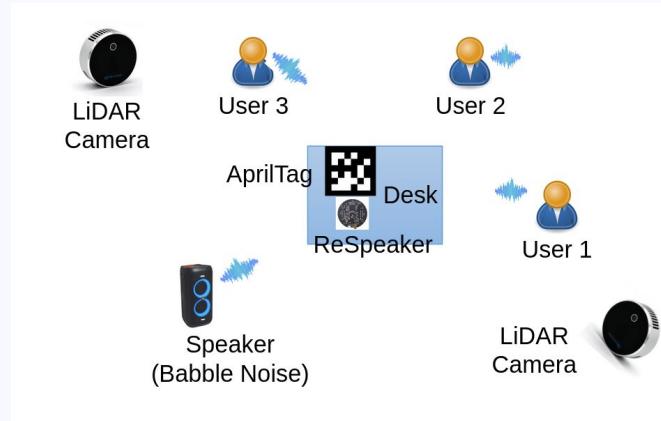
03

# EXPERIMENTS AND RESULTS

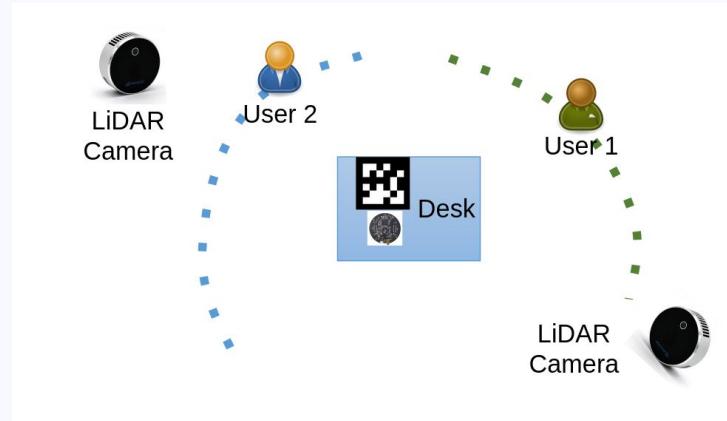


+

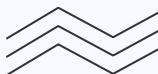
# EXPERIMENTAL SETUP



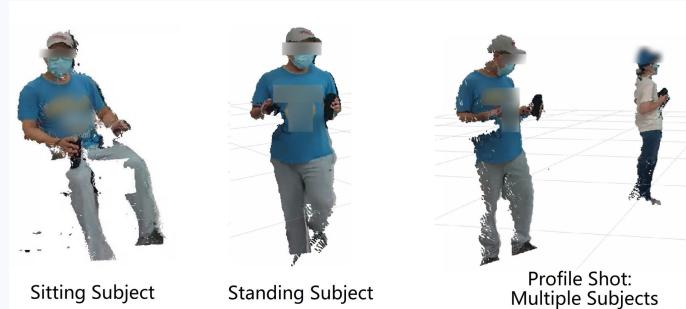
Experiment 1: Three Stationary Subjects



Experiment 2: Two Mobile Subjects



# VISUAL DOMAIN RESULTS



Visual Point Cloud Output

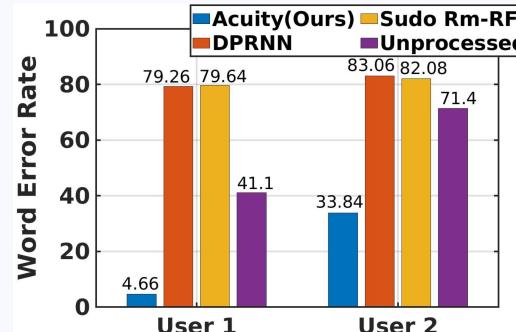
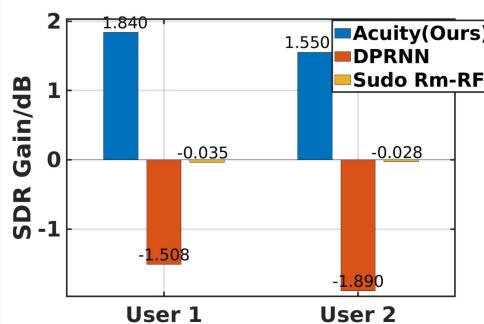
| Number of Subjects | Acuity (Ours) |          |                  | PV-RCNN++    |          |                  | SECOND       |          |                  |
|--------------------|---------------|----------|------------------|--------------|----------|------------------|--------------|----------|------------------|
|                    | Latency (ms)  | Accuracy | Average F1 Score | Latency (ms) | Accuracy | Average F1 Score | Latency (ms) | Accuracy | Average F1 Score |
| 1                  | 34            | 100%     | 1.00             | 1584         | 53.30%   | 0.33             | 33           | 88.50%   | 0.6              |
| 2                  | 54            | 100%     | 1.00             | 1635         | 60%      | 0.26             | 32           | 90.80%   | 0.78             |
| 3                  | 30            | 100%     | 1.00             | 690          | 76%      | 0.38             | 31           | 94.40%   | 0.83             |
| 4                  | 44            | 97.40%   | 0.99             | 691          | 68.80%   | 0.33             | 31           | 71.60%   | 0.63             |
| 5                  | 61            | 93.30%   | 0.96             | 699          | 75%      | 0.37             | 31           | 92.80%   | 0.77             |
| 6                  | 69            | 91.70%   | 0.95             | 693          | 70.80%   | 0.41             | 31           | 63.10%   | 0.61             |

Comparison to NN

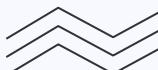
# AUDIO DOMAIN RESULTS

|               |          | SDR Gain (db) |        |          | Word Error Rate (%) |               |       |          |
|---------------|----------|---------------|--------|----------|---------------------|---------------|-------|----------|
|               |          | Acuity        | DPRNN  | SuDORMRF | Raw                 | Acuity        | DPRNN | SuDORMRF |
| Single Source | Source 1 | <b>1.92</b>   | -1.803 | -0.249   | 6.02                | <b>1.95</b>   | 60.43 | 47.075   |
| Two Sources   | Source 1 | <b>3.194</b>  | -1.094 | 0.296    | 79.02               | <b>5.08</b>   | 95.5  | 95.95    |
|               | Source 2 | <b>10.822</b> | 1.107  | 0.104    | 87.82               | <b>5.34</b>   | 78.13 | 75.4     |
| Three Sources | Source 1 | <b>5.309</b>  | 0.2899 | -        | 100                 | <b>15.49</b>  | 100   | -        |
|               | Source 2 | <b>7.025</b>  | -0.192 | -        | 100                 | <b>3.13</b>   | 92.1  | -        |
|               | Source 3 | <b>4.835</b>  | -0.145 | -        | 100                 | <b>28.425</b> | 95.5  | -        |

Experiment 1: Static Subject Audio

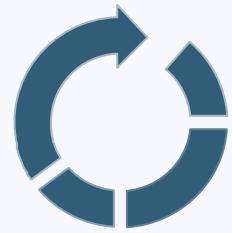


Experiment 2: Mobile Subject Audio



04

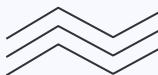
# CONCLUSION



# LIMITATIONS AND FUTURE WORK



- **Real-time Point Cloud Streaming and Rendering:** Acuity does not address issues of streaming point clouds to the end user or rendering point clouds for viewing
- **Scaling up Acuity:** Acuity currently utilizes a two camera + one microphone setup, and may benefit from the introduction of additional sensors
- **Environmental Conditions:** The LiDAR camera performs poorly in low light situations, and saturates in the presence of direct sunlight



# ACKNOWLEDGEMENTS

+

