

HEADS-UP: Head-Mounted Egocentric Dataset for Trajectory Prediction in Blind Assistance Systems

Yasaman Haghighi¹, Celine Demonsant¹, Panagiotis Chalimourdas¹, Maryam Tavasoli Naeini¹, Jhon Kevin Munoz², Bladimir Bacca², Silvan Suter³, Matthieu Gani³ and Alexandre Alahi¹

Abstract—In this paper, we introduce HEADS-UP, the first egocentric dataset collected from head-mounted cameras, designed specifically for trajectory prediction in blind assistance systems. With the growing population of blind and visually impaired individuals, the need for intelligent assistive tools that provide real-time warnings about potential collisions with dynamic obstacles is becoming critical. These systems rely on algorithms capable of predicting the trajectories of moving objects, such as pedestrians, to issue timely hazard alerts. However, existing datasets fail to capture the necessary information from the perspective of a blind individual. To address this gap, HEADS-UP offers a novel dataset focused on trajectory prediction in this context. Leveraging this dataset, we propose a semi-local trajectory prediction approach to assess collision risks between blind individuals and pedestrians in dynamic environments. Unlike conventional methods that separately predict the trajectories of both the blind individual (ego agent) and pedestrians, our approach operates within a semi-local coordinate system—a rotated version of the camera’s coordinate system—facilitating the prediction process. We validate our method on the HEADS-UP dataset and implement the proposed solution in ROS, performing real-time tests on an NVIDIA Jetson GPU through a user study. Results from both dataset evaluations and live tests demonstrate the robustness and efficiency of our approach.

I. INTRODUCTION

Approximately 4% of the global population is affected by blindness or severe visual impairment, including 43.3 million who are blind and 295 million experiencing moderate to severe visual impairment [1]. Currently, the primary navigation aid for blind individuals is the white cane, which, while accessible and affordable, offers only limited assistance. It is primarily useful for detecting static obstacles in close proximity but provides no information about dynamic obstacles, such as pedestrians, that may pose a collision risk as shown in Figure 1. This limitation is especially problematic in complex, dynamic environments, where fast-moving obstacles can lead to dangerous situations. A proper solution is to use a combination of sensors, such as cameras for online sensing, and algorithms to accurately predict the future trajectory of pedestrians and inform the blind person about potential collisions.

While previous works have suggested mounting cameras on suitcases [2] or requiring users to continuously hold a



Fig. 1: With the increasing number of blind and visually impaired individuals, coupled with advancements in vision-based algorithms, there is a growing need for intelligent assistive tools that can inform the blind person in advance about potential collisions with dynamic obstacles, such as pedestrians. Unlike the traditional white cane, which is limited to detecting local collisions with static objects, these systems offer enhanced navigation safety by predicting and warning of dynamic threats.

smartphone [3], we argue that such designs are not user-friendly. Instead, we propose a head-mounted design that offers a more practical, hands-free solution for everyday navigation. Although this approach is more user-friendly, performing trajectory prediction from head-mounted images presents significant challenges due to the constant motion of the head. To effectively tackle this problem, an egocentric dataset is required. However, existing egocentric datasets are not optimized for the trajectory prediction task. They primarily focus on individual actions [4] or social interactions between people [5], [6], with no emphasis on potential collisions between the camera wearer and other individuals in the scene. Moreover, they lack the necessary trajectory labels for predicting movement and collision risks.

Addressing the challenge of predicting trajectories in dynamic environments requires specialized datasets that capture interactions between blind individuals and moving pedestrians. To fill this gap, we developed **HEADS-UP**, a custom dataset collected using a head-mounted stereo camera, specifically focused on scenarios where pedestrians could potentially collide with the camera wearer. **HEADS-UP** comprises more than 43,000 frames, including RGB, depth, point cloud data, IMU measurements, and trajectory labels of pedestrians. The dataset also contains approximately 1,000

¹ VITA laboratory at EPFL, Lausanne, Switzerland. Email: firstname.lastname@epfl.ch

² Universidad del Valle, Cali, Colombia.

³ EssentialTech Center at EPFL, Lausanne, Switzerland. Email: firstname.lastname@epfl.ch

Corresponding author: Y. Haghighi, yasaman.haghighi@epfl.ch

individual pedestrian tracks, providing a rich source of data for trajectory prediction and collision detection tasks.

Using the **HEADS-UP** dataset, we propose a semi-local trajectory prediction baseline for blind assistance. Traditional collision detection methods typically require predicting two separate trajectories: one for the blind individual (the ego agent) and one for surrounding pedestrians. Collisions are predicted by identifying intersections between these trajectories, which necessitates knowing the positions of both the ego agent and pedestrians in a unified global coordinate system. However, we propose a simplified approach: instead of performing two independent trajectory predictions, a single prediction in a semi-local coordinate system—a rotated version of the camera’s local coordinate system—can accurately assess collision risks. Our experiments demonstrate the effectiveness and practicality of this approach.

To further validate our method, we implemented the entire pipeline in ROS and tested it in real-time on an NVIDIA Jetson GPU, demonstrating the feasibility of our approach in dynamic, real-world settings. We aim to foster further research and development in this area by making both the **HEADS-UP** dataset and the ROS implementation publicly available to the research community. For additional details, please refer to the supplementary video.

II. RELATED WORK

A. Assistive tool for blinds

There are many blind individuals worldwide, highlighting the need for more advanced assistive tools beyond the traditional white cane. Recent advancements in computer vision algorithms suggest that using cameras to develop assistive tools has the potential to significantly improve navigation and safety for the visually impaired. Among the proposed solutions, [3] suggests using a combination of a smartphone and a white cane to build a 2D occupancy map of indoor environments. However, this approach focuses on static environments, and constantly holding a phone is not efficient for users.

[7] leverages recent advancements in vision-language models for video anomaly detection. While this approach can provide context to blind users, it lacks the accuracy required for detecting and predicting pedestrian trajectories, which is critical for providing real-time hazard warnings.

For dynamic environments, [2] suggests using two RGB-D cameras, a LiDAR sensor, and an IMU, all mounted on a suitcase. This setup detects pedestrians in a global coordinate system and predicts their future trajectories. However, the design is not user-friendly due to its bulky nature.

In contrast, we propose using a head-mounted stereo camera, which is more comfortable to wear. Additionally, we introduce a semi-local coordinate system, allowing for a single trajectory prediction. This method is more efficient than predicting two global trajectories, significantly improving the overall efficiency of the pipeline.

B. Related datasets

Among datasets captured using head-mounted cameras, Ego4D [4] is one of the largest, containing over 3,000 hours of video, primarily captured with GoPro cameras. This dataset focuses mainly on daily activities and is suitable for action recognition. Mo2Cap2 [8] is another egocentric dataset captured with a fisheye camera, but it focuses on the wearer of the camera rather than other people in the environment. You2Me [9] captures interactions between two people, while EgoBody [5] provides egocentric captures of two or more people interacting in indoor scenes. EgoHuman [6] offers synchronized egocentric views from four people interacting in various scenarios, including activities like playing volleyball, in both indoor and outdoor environments.

For our scenario, we require a dataset not only captured from an egocentric perspective, to resemble realistic head movements, but also where the camera wearer walks through an environment with the potential for collisions with pedestrians. None of the aforementioned datasets include such cases, nor do they provide pedestrian trajectory labels.

Among the datasets commonly used for human trajectory prediction, JRDB [10] is captured using a social mobile manipulator, and JTA [11] is synthetic. While both provide trajectory labels, neither is suitable for our scenario because they are not captured from head-mounted cameras and, thus, do not account for the natural head movements that are crucial in our context.

C. Pedestrian trajectory prediction

Kalman filter [12] is a basic approach for pedestrian motion estimation tasks. The Kalman filter is a recursive algorithm that estimates the current state of a pedestrian (e.g., position and velocity) based on noisy measurements and predicts the future state of the pedestrian’s trajectory using a motion model, then updating this prediction with incoming observations to minimize uncertainty. While effective for linear and smooth motions, the Kalman filter may struggle with non-linear or highly dynamic pedestrian behaviors, requiring more complex models for improved accuracy.

To improve the accuracy of trajectory predictions, early models focused on the attractive and repulsive forces between pedestrians [13] or used Bayesian inference to model human-environment interactions [14]. Later, data-driven methods [15], [2], [16], [17], [18], [19] were introduced to capture interactions between pedestrians, as their future trajectories are often influenced by their surroundings and other individuals. These approaches include modeling observed neighbor interactions through hidden states or directional grids.

Among various architectures like RNNs [15], GANs [20], and diffusion models [21], state-of-the-art methods have increasingly shifted toward using transformers [22]. In addition to leveraging pedestrians’ 2D positional data, these models often incorporate visual cues, such as detected bounding boxes, to enhance prediction accuracy.

While these methods are highly accurate, applying them to collision detection between pedestrians and blind individuals



Fig. 2: We use a ZED Mini stereo camera [23] to capture the dataset, securely mounted on a cap using a custom 3D-printed attachment to ensure stable positioning during data collection.

poses challenges. Specifically, it requires representing the positions of both the pedestrian and the blind person in a unified global coordinate system. Moreover, it necessitates two separate predictions, one for the blind person and one for the pedestrian. Afterward, the system must determine whether the predicted trajectories will collide and issue a hazard warning to the blind person.

In contrast, we propose using a semi-local coordinate system to jointly predict the future relative trajectories of both agents. We demonstrate that existing approaches can perform effectively in this semi-local coordinate system, achieving results comparable to global prediction methods.

III. DATASET

In this section, we describe our data capture setup and collection methodology, detailing the process of annotating the data, as well as providing key dataset statistics.

A. Data collection setup

We captured the dataset using a ZED mini camera [23] attached to a cap worn by the blind individual. An illustration of our capture setup is shown in Figure 2. The ZED mini camera provided synchronized RGB images, depth maps, and IMU measurements, all recorded at a frame rate of 30 FPS. The camera sensor specifications are summarized in Table I.

B. Data collection protocol

To tackle the task of navigation for blind individuals using a head-mounted camera, we focused on outdoor scenarios with varying levels of collision risk. The dataset captures realistic pedestrian interactions, particularly in situations where collisions with the blind individual were likely.

We considered three primary experimental setups:

- **Easy setup:** The blind person walked slowly with minimal head movement, creating a controlled environment with limited distractions or challenges for trajectory prediction.
- **Hard setup:** The blind individual made rapid and drastic head movements, reflecting more realistic real-world behavior. This setup introduces greater complexity, as it involves more erratic head motion and pose changes, which can impact pedestrian detection and tracking.

TABLE I: Specifications of the ZED Mini camera used for dataset acquisition.

Specification	Details
2x Camera	RGB, 30 Hz, 1920×1080
Depth Range	0.5 m to 25 m
Depth Accuracy	< 2% up to 3 m < 4% up to 15 m
IMU	ZED Mini built-in Gyroscope and Accelerometer, 800 Hz

TABLE II: Dataset modalities and specifications across different subsets: Easy, Hard, and Uncontrolled.

Subset	RGB	Depth	Point Cloud	IMU	Camera Pose	Total Frames	# Agent Tracks
Easy	✓	✓	✓	✓	✓	15,510	224
Hard	✓	✓	✓	✓	✓	14,038	169
Uncontrolled	✓	✓	✓	✓	✓	13,665	566

- **Uncontrolled setup:** In this scenario, the blind individual navigated in an environment with multiple pedestrians, increasing the likelihood of collisions. This scenario simulates a highly dynamic and uncontrolled real-world environment with high pedestrian density and movement.

Examples of RGB and depth images from each of the subsets are shown in Figure 3.

C. Data annotation

To obtain camera poses, we employed the Visual-Inertial Odometry (VIO) [24] module of the ZED SDK [25]. For trajectory labeling, pedestrians were first detected using the YOLOv8 [26] object detection model. ByteTrack [27] was then employed to track the pedestrians over time, representing each by the center of their detected bounding box.

Due to the head-mounted nature of the system, the resulting trajectories were affected by noise from camera motion, depth capture variability, and detection model inconsistencies. To mitigate the impact of noise, we first downsampled the detected pedestrian positions to 2.5 FPS by averaging positions over consecutive frames, reducing frame-to-frame jitter. Then, we applied smoothing Kalman filter [12] to further refine the pedestrian paths. Pedestrians who were tracked for fewer than 6 frames were excluded from the dataset to ensure reliable trajectory data. An example of the trajectories before and after smoothing is shown in Figure 4. Finally, for privacy reasons, we blurred the faces of all pedestrians in the sequences.

D. Dataset statistics

Our dataset comprises over 43,000 frames, including RGB images, depth data, IMU readings, point clouds, and camera poses, along with approximately 1,000 agent tracks. This provides a comprehensive and rich resource for addressing the task of navigation for blind individuals using head-mounted cameras. Further details are provided in Table II.

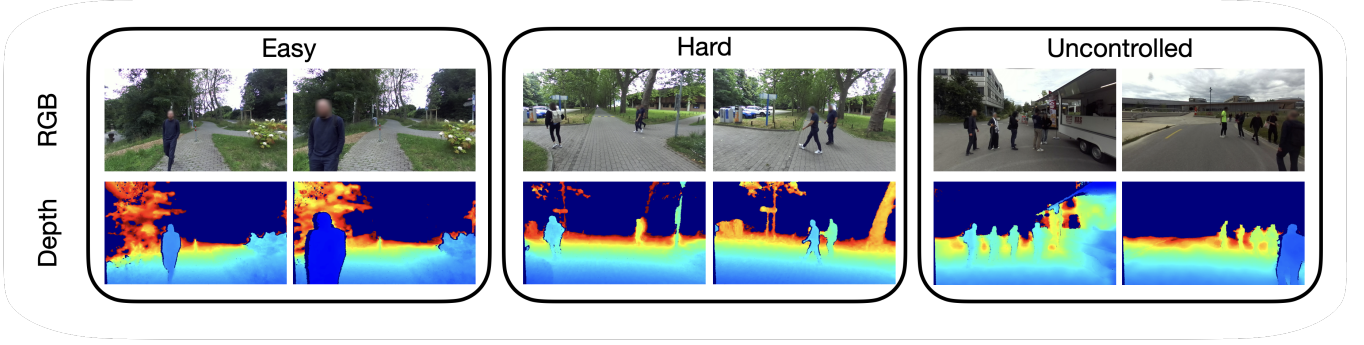


Fig. 3: Example of RGB and depth images from each subset of the dataset. In the Easy subset, head movement is limited, while in the Hard subset, more drastic head movements are present. Both subsets have a controlled setup with a limited number of pedestrians. In contrast, the Uncontrolled setup features multiple pedestrians, a higher possibility of collisions, and drastic head movements. These diverse subsets enable the development and evaluation of algorithms in a variety of scenarios, facilitating research in blind navigation systems.

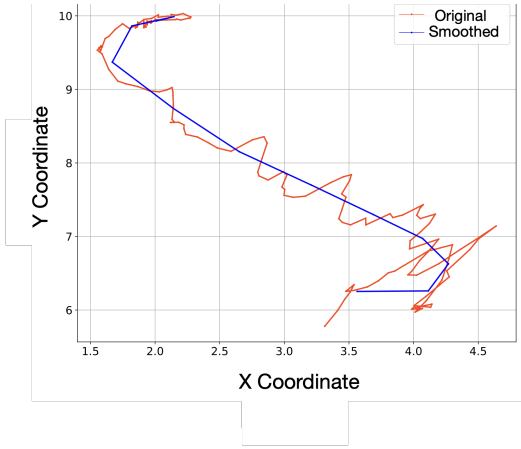


Fig. 4: An example of pedestrian trajectories before (red) and after (blue) applying smoothing techniques, demonstrating reduced noise and improved trajectory reliability. The X and Y coordinates are in meters.

IV. SEMI-LOCAL PREDICTION

In this section, we present our semi-local trajectory prediction baseline. To inform the blind person about potential collisions, we need to predict the future trajectories of both the pedestrians and the blind individual. If their predicted paths intersect, we can alert the blind person to the risk of collision. Traditionally, this requires knowing the positions of both the pedestrian and the blind person in a unified global coordinate system.

However, we propose a simpler approach: using a single trajectory prediction in a semi-local coordinate system. Unlike a purely local coordinate system, where everything is relative to the camera (or the blind individual), the semi-local system takes into account the camera’s rotation. This adjustment is crucial, as the blind person frequently rotates their head, and failing to consider this would result in inaccurate predictions. We run a live user study to demonstrate the effectiveness of this baseline. Below, we provide the

mathematical reasoning that supports the validity of this approach:

Assume each detected and tracked pedestrian in the RGB image is represented by the center of the bounding box, denoted as $\mathbf{p}_{uv} = [u \ v \ 1]^T$, where u and v are the pixel coordinates of the bounding box center. Using the corresponding depth d_{uv} and the perspective camera model, we backproject the 2D pixel coordinates into camera coordinates. The backprojection from the pixel coordinates in the image plane to the camera coordinate system can be expressed as:

$$\mathbf{p}_{\text{camera}} = d_{uv} \cdot \mathbf{K}^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}, \quad (1)$$

where $\mathbf{p}_{\text{camera}}$ is the 3D position of the pedestrian in the camera coordinate system, \mathbf{K} is the intrinsic camera matrix, and d_{uv} is the depth at pixel (u, v) .

We then transform the 3D point from the camera frame to the global frame by applying the following transformation:

$$\mathbf{p}_{\text{global}} = \mathbf{R} \cdot \mathbf{p}_{\text{camera}} + \mathbf{t}, \quad (2)$$

where \mathbf{R} is the camera rotation matrix and \mathbf{t} is the translation vector, calculated by algorithms such as [24]. Additionally, \mathbf{t} represents the position of the ego agent (blind individual) in the global coordinate system.

Assume trajectory prediction algorithm is denoted by the function $f(\cdot)$, to detect potential collisions, at each time step i we calculate:

$$f_i(\underbrace{\mathbf{R} \cdot \mathbf{p}_{\text{camera}}}_{\text{Pedestrian}} + \mathbf{t}) - f_i(\underbrace{\mathbf{t}}_{\text{Blind}}) = 0. \quad (3)$$

We argue that for our setup, this can be approximated as:

$$f_i(\underbrace{\mathbf{R} \cdot \mathbf{p}_{\text{camera}}}_{\text{Semi-Local}}) \approx 0. \quad (4)$$

This holds when the function f is linear, such as in the case of a Kalman filter [12]. In our experiments, we observe that this approximation extends to non-linear neural network based algorithms such as [22].

V. LIMITATIONS AND FUTURE WORK

Although we have proposed a baseline for collision detection in blind navigation and demonstrated its effectiveness, there are still challenges that need to be addressed. First, semi-local trajectory prediction requires accurate estimation of the camera's rotation. In our current ROS implementation, we rely on IMU rotation measurements, which are prone to error. An end-to-end pipeline that directly estimates the possibility of collision using only RGB and depth images would be more robust and reliable.

With the rich annotations and multimodal data in our dataset, future work could focus on training models that enhance both the efficiency and performance of collision detection for blind navigation. This would help overcome the limitations of relying on external sensors like the IMU and provide a more accurate solution. Ultimately, this could pave the way for more advanced assistive systems that offer real-time, reliable support for visually impaired individuals in navigating complex and dynamic environments.

VI. CONCLUSION

In this paper, we introduced **HEADS-UP**, the first ego-centric dataset captured using a stereo head-mounted camera, specifically designed for blind navigation. The dataset intentionally simulates scenarios with potential collisions between a blind individual and pedestrians in three Easy, Hard and Uncontrolled setups. It offers multiple modalities, including RGB, depth, IMU measurements, camera poses and trajectory labels, making it a valuable resource for the novel task of collision detection in blind navigation.

Additionally, we proposed a semi-local trajectory prediction baseline and demonstrated its effectiveness both on our dataset and in online user study evaluations. Our results show that the semi-local approach is a viable alternative to traditional global methods, providing accurate and efficient collision detection. This work paves the way for future advancements in assistive technologies for visually impaired individuals, with potential for further improvements using our dataset.

REFERENCES

- [1] V. L. E. G. of the Global Burden of Disease Study *et al.*, "Global estimates on the number of people blind or visually impaired by cataract: a meta-analysis from 2000 to 2020," *Eye*, vol. 38, no. 11, p. 2156, 2024.
- [2] S. Kayukawa, T. Ishihara, H. Takagi, S. Morishima, and C. Asakawa, "Guiding blind pedestrians in public spaces by understanding walking behavior of nearby pedestrians," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 4, no. 3, pp. 1–22, 2020.
- [3] M. Kuribayashi, S. Kayukawa, J. Vongkulbhisal, C. Asakawa, D. Sato, H. Takagi, and S. Morishima, "Corridor-walker: Mobile indoor walking assistance for blind people to avoid obstacles and recognize intersections," *Proceedings of the ACM on Human-Computer Interaction*, vol. 6, no. MCHI, pp. 1–22, 2022.
- [4] K. Grauman, A. Westbury, E. Byrne, Z. Chavis, A. Furnari, R. Girdhar, J. Hamburger, H. Jiang, M. Liu, X. Liu *et al.*, "Ego4d: Around the world in 3,000 hours of egocentric video," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 18995–19012.
- [5] S. Zhang, Q. Ma, Y. Zhang, Z. Qian, T. Kwon, M. Pollefeys, F. Bogo, and S. Tang, "Egobody: Human body shape and motion of interacting people from head-mounted devices," in *European conference on computer vision*. Springer, 2022, pp. 180–200.
- [6] R. Khirrodar, A. Bansal, L. Ma, R. Newcombe, M. Vo, and K. Kitani, "Ego-humans: An ego-centric 3d multi-human benchmark," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 19807–19819.
- [7] H. Wang, J. Qin, A. Bastola, X. Chen, J. Suchanek, Z. Gong, and A. Razi, "Visiongpt: Llm-assisted real-time anomaly detection for safe visual navigation," *arXiv preprint arXiv:2403.12415*, 2024.
- [8] W. Xu, A. Chatterjee, M. Zollhoefer, H. Rhodin, P. Fua, H.-P. Seidel, and C. Theobalt, "Mo 2 cap 2: Real-time mobile 3d motion capture with a cap-mounted fisheye camera," *IEEE transactions on visualization and computer graphics*, vol. 25, no. 5, pp. 2093–2101, 2019.
- [9] E. Ng, D. Xiang, H. Joo, and K. Grauman, "You2me: Inferring body pose in egocentric video via first and second person interactions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9890–9900.
- [10] R. Martin-Martin, M. Patel, H. Rezafofighi, A. Sheno, J. Gwak, E. Frankel, A. Sadeghian, and S. Savarese, "Jrdb: A dataset and benchmark of egocentric robot visual perception of humans in built environments," *IEEE transactions on pattern analysis and machine intelligence*, vol. 45, no. 6, pp. 6748–6765, 2021.
- [11] M. Fabbri, F. Lanzi, S. Calderara, A. Palazzi, R. Vezzani, and R. Cucchiara, "Learning to detect and track visible and occluded body joints in a virtual world," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 430–446.
- [12] R. E. Kalman, "A new approach to linear filtering and prediction problems," *Journal of Basic Engineering*, vol. 82, no. 1, pp. 35–45, 1960.
- [13] D. Helbing and P. Molnar, "Social force model for pedestrian dynamics," *Physical review E*, vol. 51, no. 5, p. 4282, 1995.
- [14] G. Best and R. Fitch, "Bayesian intention inference for trajectory prediction with an unknown goal destination," in *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2015, pp. 5817–5823.
- [15] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, "Social lstm: Human trajectory prediction in crowded spaces," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 961–971.
- [16] F. Giuliani, I. Hasan, M. Cristani, and F. Galasso, "Transformer networks for trajectory forecasting," in *2020 25th international conference on pattern recognition (ICPR)*. IEEE, 2021, pp. 10335–10342.
- [17] P. Kothari, S. Kreiss, and A. Alahi, "Human trajectory forecasting in crowds: A deep learning perspective," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 7, pp. 7386–7400, 2021.
- [18] G. Chen, Z. Chen, S. Fan, and K. Zhang, "Unsupervised sampling promoting for stochastic human trajectory prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 17874–17884.
- [19] J. Sun, Y. Li, L. Chai, H.-S. Fang, Y.-L. Li, and C. Lu, "Human trajectory prediction with momentary observation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 6467–6476.
- [20] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi, "Social gan: Socially acceptable trajectories with generative adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2255–2264.
- [21] T. Gu, G. Chen, J. Li, C. Lin, Y. Rao, J. Zhou, and J. Lu, "Stochastic trajectory prediction via motion indeterminacy diffusion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 17113–17122.
- [22] S. Saadatnejad, Y. Gao, K. Messaoud, and A. Alahi, "Social-transmotion: Promptable human trajectory prediction," *arXiv preprint arXiv:2312.16168*, 2023.
- [23] Stereolabs, "Zed mini stereo camera," n.d., accessed: 2024-09-15. [Online]. Available: <https://www.stereolabs.com/en-ch/store/products/zed-mini>

- [24] S. Leutenegger, S. Lynen, M. Bosse, R. Siegwart, and P. Furgale, "Keyframe-based visual-inertial odometry using nonlinear optimization," *The International Journal of Robotics Research*, vol. 34, no. 3, pp. 314–334, 2015.
- [25] Stereolabs, "Zed sdk," n.d., accessed: 2024-09-15. [Online]. Available: <https://www.stereolabs.com/en-ch/developers/release>
- [26] G. Jocher, A. Chaurasia, and J. Qiu, "Ultralytics YOLO," Jan 2023. [Online]. Available: <https://github.com/ultralytics/ultralytics>
- [27] Y. Zhang, P. Sun, Y. Jiang, D. Yu, F. Weng, Z. Yuan, P. Luo, W. Liu, and X. Wang, "Bytetrack: Multi-object tracking by associating every detection box," in *European conference on computer vision*. Springer, 2022, pp. 1–21.