

PEOPLE AWARE MOBILE ROBOT NAVIGATION

A Thesis
Presented to
The Academic Faculty

by

Akansel Cosgun

In Partial Fulfillment
of the Requirements for the Degree
Doctor of Philosophy in the
College of Computing

Georgia Institute of Technology
November 2015

PEOPLE AWARE MOBILE ROBOT NAVIGATION

Approved by:

Professor Henrik Iskov Christensen,
Advisor
College of Computing
Georgia Institute of Technology

Professor Andrea Thomaz
College of Computing
Georgia Institute of Technology

Professor Irfan Essa
College of Computing
Georgia Institute of Technology

Professor Ayanna Howard
School of Electrical and Computer
Engineering
Georgia Institute of Technology

Dr Emrah Akin Sisbot
Toyota InfoTechnology Center

Date Approved: X November 2015

PREFACE

Preface

ACKNOWLEDGEMENTS

Acknowledgements

TABLE OF CONTENTS

PREFACE	iii
ACKNOWLEDGEMENTS	iv
LIST OF TABLES	ix
LIST OF FIGURES	x
SUMMARY	xv
I INTRODUCTION	1
1.1 Robots in Human Environments	1
1.2 Robots that Navigate in Human Environments	1
1.3 The Science of Personal Spaces	2
1.4 Semantic Maps	2
1.5 Tour Scenario	3
1.6 Challenges	3
1.6.1 Robust People Tracking	4
1.6.2 Social Navigation	4
1.6.3 Interactive Map Annotation	5
1.6.4 Situational Awareness for Navigation	5
1.7 Scope and Context	6
1.8 Thesis Overview	7
1.8.1 Thesis Statement	7
1.8.2 Contributions	7
1.8.3 Document Outline	8
II MAP ANNOTATION	9
2.1 Related Work	10
2.2 Semantic Maps	10
2.2.1 Navigation Waypoints	11
2.2.2 Planar Landmarks	11

2.2.3	Objects	13
2.3	User Interface	13
2.4	Pointing Gestures for Human-Robot Interaction	15
2.4.1	Related Works	16
2.4.2	Pointing Gesture Recognition	17
2.4.3	Representing Pointing Directions	18
2.4.4	Determining Intended Target	19
2.4.5	Data Collection	21
2.4.6	Ground Truth Target Positions	21
2.4.7	Skeleton Data Capture	22
2.4.8	Pointing Gesture Annotation	24
2.4.9	Error Analysis	24
2.4.10	Evaluation	28
2.4.11	Object Separation Test	28
2.4.12	Results and Discussion	30
2.4.13	Conclusions	31
III	NAVIGATION AMONG PEOPLE	32
3.1	State-of-the-Art Approach in Autonomous Navigation	33
3.2	Related Work	35
3.2.1	Socially Acceptable Path Planning	35
3.2.2	Learning Navigation from Human Behavior	36
3.2.3	Human Cooperation in Robot Navigation	36
3.3	Goal Points for Navigation	37
3.3.1	Labeled Waypoints:	37
3.3.2	Labeled Planar Landmarks:	37
3.3.3	Labeled Objects:	39
3.4	People Aware Navigation	39
3.4.1	Global Planner	41

3.4.2	Local Planner	45
3.4.3	Results	47
3.5	Speed Limits for Safe Navigation	50
3.5.1	Results	54
IV	MULTIMODAL PERSON DETECTION AND TRACKING	57
4.1	Related Work	58
4.2	Person Detection	60
4.2.1	Leg Detection	60
4.2.2	Torso Detection	65
4.3	Person State Estimation	69
4.4	Face Recognition	73
V	PERSON FOLLOWING	75
5.1	Related Work	76
5.2	Basic Person Following	77
5.3	Situation Aware Person Following	80
5.3.1	Door Passing	80
5.3.2	User Activity Awareness	80
5.3.3	Corners	80
5.4	Application To Telepresence Robots	80
5.4.1	Robot Platform	81
5.4.2	User Study	83
5.4.3	Design Implications	86
5.4.4	Discussion	89
VI	PERSON GUIDANCE	90
6.1	Related Work	90
6.2	Guide Robot	91
6.2.1	Results	92
6.3	Application To Blind Users	94

6.3.1	Tactile Belt	94
6.3.2	Planning the Path of the User	95
6.3.3	Velocity to Vibration Mapping	96
6.3.4	Demonstration	96
VII	CONCLUSION	98
APPENDIX A	— QR CODE BASED LOCATION INITIALIZA-	
	TION	99
APPENDIX B	— ASSISTED REMOTE CONTROL	100
APPENDIX C	— TELEPRESENCE STUDY SURVEY QUESTIONS	101
APPENDIX D	— VIBRATION PATTERN ANALYSIS FOR HAP-	
	TIC BELTS	102
REFERENCES		108

LIST OF TABLES

1	μ and σ angular errors in degrees for each of Targets 1-4 (Figure 5), for each pointing method. The aggregate μ and σ is also shown.	25
2	μ and σ of angular error in degrees for Targets 5-7 (Figure 5), for each pointing method. The aggregate μ and σ is also shown.	25
3	Table shows average and standard deviations of geometric leg features calculated in our dataset.	63
4	Table shows average and standard deviations of geometric features for a human torso in laser scans.	67
5	Average orientation error of the torso detector with respect to distance from sensor and body pose in a study with 23 people	69
6	Survey results of the user study for person following for telepresence robots. Table displays survey question average and standard deviations for the two conditions: Autonomous Person Following and Manual Person Following.	86
7	Average recognition error and reaction times of directional patterns .	105
8	Average recognition accuracy and reaction times of rotational patterns	106

LIST OF FIGURES

1	Robot’s view of how users can add planar landmarks to the map	12
2	UI for labeling a navigation waypoint	14
3	(Left) Our approach allows a robot to detect when there is ambiguity on the pointing gesture targets. (Right) The point cloud view from robot’s perspective is shown. Both objects are identified as potential intended targets, therefore the robot decides that there is ambiguity.	15
4	Vertical (ψ) and horizontal (θ) angles in spherical coordinates are illustrated. A potential intended target is shown as a star. The z-axis of the hand coordinate frame is defined by either the Elbow-Hand (this example) or Head-Hand ray.	19
5	Our study involved 6 users that pointed to 7 targets while being recorded using 30 frames per target.	22
6	Data capturing pipeline for error analysis.	22
7	Euclidean distance error in cartesian coordinates for each method and target. The gesture ray intersection points in centimeters with the target plane, are shown here for each target (T1-T7) as the columns. Each subject’s points are shown in separate colors. There are 30 points from each subject, corresponding to the 30 frames recorded for the pointing gesture at each target. Axes are shown in centimeters. The circle drawn in the center of each plot has the same diameter (13 cm) as the physical target objects used.	23
8	Box plots of the errors in spherical coordinates θ and ψ for each pointing method.	26
9	Resulting Mahalanabis distances of pointing targets from the Object Separation Test is shown for a) Elbow-Hand and b) Head-Hand pointing methods. Intended object are shown in green and other object is in red. Solid lines show distances after correction is applied. Less Mahalanobis distance for intended object is better for reducing ambiguity.	28
10	Example scenarios from the object separation test is shown. Our experiments covered separations between 2cm (left images) and 30cm (right images). The object is comfortably distinguished for the 30cm case, whereas the intended target is ambiguous when the targets are 2cm apart. Second row shows the point cloud from Kinect’s view. Green lines show the Elbow-Hand and Head-Hand directions whereas green circles show the objects that are within the threshold $D_{mah} < 2$	29
11	Caption	34

12	Top down point cloud view of a room. A planar landmark with label <i>Table</i> has previously been annotated by a user. The convex hull for the planar landmark is shown in red lines. When asked to navigate to <i>Table</i> , the robot calculates a goal pose, which is shown as the yellow point.	38
13	Top down point cloud view of a hallway. The user has previously annotated two planar landmarks with the same label, <i>Hallway</i> . When asked to navigate to <i>Hallway</i> , the robot calculates a goal pose, which is shown as the yellow point.	39
14	Standard path planners fail to produce a solution to the 'room problem'. Our people-aware planner anticipates that the human can give way to the robot if it approaches towards its goal.	40
15	Disturbance costs in different human-human configurations and distances. A path that crosses the dashed lines incurs the disturbance cost calculated on the right side.	43
16	a) Social forces acting on the robot, including $F_{goal}, F_{social}, F_{obs}$, are shown at the first iteration of the dynamic planner. Note that $F_{group} = 0$ as the robot does not belong to a group. The group force (not shown) exists, however, for the humans as they are in the same group region. b) Social forces with respect to the distance towards the corresponding entity.	45
17	"Room Problem". The robot is outside a room and the goal is inside the room. Traditional planners can not solve the problem because two people are blocking the doorway. Our planner generates a tentative path, with the initial global plan shown in green and the dynamic refinements are shown in orange.	47
18	Paths differ drastically with the poses and grouping of humans. a) The robot takes shortest route, traveling in the vicinity of a group of two and another individual. b) third individual joins the group. Robot takes a longer path that doesn't have humans on path. c) fourth person changes his position, leading the robot to take the longest route. . .	48

- 19 The Hallway scenario. 2 runs are shown in first and second rows. The static plan (green line) and dynamic plan refinement (pink line) are shown. First run: a) Navigation starts. The dynamic planner anticipates that people will give way to the robot when it starts to move towards them. b) Humans notice the robot, and give way by increasing the separation between them. c) The robot continues towards its goal and humans regroup. Second run: d) both the static and dynamic plan involves going in between humans again e) human on the right gets closer to the other person. Since a human made significant movement, dynamic planner re-plans. Plan no longer involves going in between. f) static planner periodic re-plan triggers, leading to robot to stick to the wall to the right.

51

20 The Kitchen scenario. In the first run, there are two people blocking the path to the left and one person at the narrow corridor. a) robot decides to take the shorter route, because it would disturb one person instead of two. There is not enough space to pass, and dynamic planner assumes the person would get out of the bottleneck to give way. b) human behaves as robot anticipated and gets out of the narrow passage. robot slows down because it enters the human region. c) person gets back to his original position, robot reaches the goal. In the second run: d) there are two people at the narrow corridor and one person on the left. The robot decides to take the longer route and pass the third person from left. The safety cost from the two others would be too high if the robot took the direct route. e) the person steps back as he recognizes the robot. since the person has moved, the dynamic planner re-plans and decides to pass from right. f) after the robot passes the person, it proceeds to its goal.

52

21 Designed speed limits for. The robot has to be relatively slow in red zones, can have moderate speed in yellow zones and is allowed to move relatively faster in green zones.

53

22	Comparison study of using a maximum top speed versus using location-dependent speed limits. Robot is given a fixed goal location. Right around the corner, there is a bystander human, who is not visible to the robot until the robot makes the turn. Points annotate robot position measured at fixed time intervals. a) Speed map of a corridor intersection at the second floor of College of Computing at Georgia Tech. b) Robot's top speed is fixed at $1.0m/s$. Note that the distance between robot positions are mostly constant. The robot gets very close to the bystander because it is moving relatively fast when it turned the corner. c) The robot is allowed to move with $1.5m/s$ in green, $0.5m/s$ in yellow and $0.15m/s$ in red zones. Colors of the sampled points on the path show the associated speed zone. It can be seen by looking at robot's positions that this approach was more gracious turning the corner and respecting human's personal space.	56
23	Circularity criterion in a perfect circle is: $ P_0P_n d_{mid} = 0.5$	62
24	Circularity criterion in a this laser segment is: $ P_0P_{10} /d_{mid}$	62
25	Inscribed angles of an arc are shown in the figure. Inscribed Angle Variance (IAV) is calculated by taking the average of all inscribed angles on a laser segment.	62
26	Two person detections are seen in this figure. Our leg segment association algorithm propagates pixels vertically from candidate leg segments and connects leg pairs.	64
27	Flow chart for determining if two leg segment candidates belong to a person.	65
28	Our torso detector fits and ellipse to the human torso and estimate its position and orientation.	66
29	Torso detection rate vs weighed Mahalanobis Distance Threshold in our dataset	68
30	Experimental setup for the evaluation study of the Torso Detector.	68
31	Example results of our person recognition method is shown in the image. We use <i>Eigenfaces</i> face recognition method and optionally shirt color recognition.	73
32	Overhead view of relevant ranges for person following. Robot is represented as the triangle in the middle.	78
33	An illustration of how the goal position is calculated when the user is in the social space $[1.2m - 3.5m]$	79
34	The telepresence robot platform we used for our experiments.	82

35	User Interface of the robot for the remote user.	83
36	Speed profile of a person guiding robot as a function of the distance to the user.	92
37	Comparison of robot and human speeds with respect to time. a) Standard ROS Navigation b) Our approach: accelerations are less steeper than a), which employs the dynamic speed adjustment for guidance.	93
38	The vibration pattern applied by the Tactile Belt to induce directional movements in the blind user. A motor is fired for a duration of <i>250ms</i> , inactivated for <i>250ms</i> and fired again for <i>250ms</i>	95
39	The vibration pattern applied by the Tactile Belt to induce rotational movements in the blind user. The consequent vibrations motors are fired consecutively, starting from left for CW and right for CCW rotation.	95
40	Autonomous guiding of a blindfolded person using the tactile belt. a) The guidance starts. The user is blindfolded and is standing at the left of the screen. The human detection system detects him and places an ellipse marker with an arrow depicting his orientation. The operator gives a goal point by clicking on the screen. The goal point is the right traffic cone, and given by the big arrow. b) The system autonomously generates a path for the user. As seen in the picture the path is collision free. At this stage the belt begin to vibrate towards the front of the user. c) An unexpected obstacle (another person) appears and stops in front of the user. The system detects the other person as an obstacle, and reevaluates the path. A new path going around the obstacle is immediately calculated and sent to the user by the belt. d) The user receives a rotation vibration modality, and begins to turn towards the new path. And follows this path from now on. e) The obstacle leaves. The path is then reevaluated and changed. The user receives forward directional belt signal, and advances towards the goal. f) The person reaches to the vicinity of the goal and stop signal is applied.	97
41	Evaluated vibration patterns. a) Directional b) Rotational	102
42	Confusion matrix of recognition accuracy of directional vibration patterns. Our results show that the recognition accuracy is highly dependent on the applied direction. The subjects recognised the front direction (Direction 1) with the highest accuracy (%97) whereas Direction 3 had the least accuracy (%55)	105
43	Results of post-study survey.	107

SUMMARY

Why should I provide a summary? Just read the thesis.

CHAPTER I

INTRODUCTION

1.1 *Robots in Human Environments*

The vision that was promised by Karel Čapek in his 1921 science fiction play R.U.R, in which the word “*robot*” was used the first time, featured intelligent autonomous systems that co-existed and worked for humans. The reality almost a century later today is that robots do exist, however almost all of them operate in factories, physically separated from human workers. With this separation, the safety of humans is ensured. However, it is now becoming clear that robots can provide much more value if they operate in human environments, assisting people in daily tasks. With the recent improvements in hardware and software systems and in robotics research, development of such robots are not science fiction anymore. This disruptive technology will make fascinating new applications possible, in contexts such as homes, hospitals, offices and factories, including but not limited to: delivery, elderly care, collaborating on an assembly line and household tasks such as cleaning. The feasibility and reliability of such applications will determine their business value, while potentially generating a new industry. Therefore, any development in this field is a step toward realization of *intelligent* and *social* robots.

1.2 *Robots that Navigate in Human Environments*

The most important difference of a robot from a computer is actuation: the ability to move around in the physical world. Perhaps the most fundamental task for a mobile robot is to be able to *navigate* from a location to another. Advances in mapping, localization and path planning research in the past decades made indoor navigation

a common capability for mobile robots.

Robots occupy the same physical space with us, therefore should be aware of the fact that their movements would be observed and reacted by humans. Investigating the context of robots in human environments gave rise to the sub-field of robotics: Human-Robot Interaction (HRI). This relatively new research field extends in many directions with a common goal: robots that will perceive its environment, reason and act in a safe way to facilitate people's lives. Thus, a robot that will serve to people should not only be a machine, but should respect social rules and protocols.

1.3 The Science of Personal Spaces

How should mobile agents adjust their spatial relationships with respect to humans? In fact, humans already have a framework for this problem. It is called the *proxemics*, coined by Hall [?], which studies of our use of space and the degree and effect of the spatial separation individuals naturally maintain. One of Hall's main findings is that humans adjust their distance in four distance levels: intimate, personal, social and public distances, depending on the intimacy level between individuals.

Personal spaces are influenced by other factors such as cultures, the range of their sensory systems, and postures of the two actors. Although it is not known whether it is best to model robot navigation behaviors after humans, using human-human studies to determine the spatial relationships to humans is a good start. This is especially true given that non-expert humans today had never interacted with a mobile robot and therefore likely expect the robot to behave within human social norms.

1.4 Semantic Maps

Robots keep a representation of its environment, and it is usually in the form of a metric *map*. It may not be the most intuitive way for human users to provide metric coordinates as goals points to robots. Robots that navigate in human environments will need to accept human-friendly navigation goals. A map for such robots should act

a *common ground* between the robot and its user, by enabling referencing to the same spatial elements. For example, the robot should be able to understand commands such as “go to the kitchen table” or “wait outside Joe’s room”. This is only possible with a richer map representation that includes *semantic* information. In particular, if the humans and robots need to communicate over navigational goals that involve objects, landmarks or rooms, a map should include information about these.

1.5 Tour Scenario

One method for a robot to collect semantic information in an environment is the *Tour Scenario*, or *Human Augmented Mapping (HAM)* [?], in which the task is to familiarize the robot to a new environment interactively. Here is how the Tour Scenario is defined: After the robot product is shipped and unpacked to a home or similar environment, a human takes the robot on a tour so the robot create a representation of the environment. In this interactive scenario, the robot is able to learn basic information about the surroundings, while following the user. During the tour, the user could point out and identify relevant places and objects that the robot will need to learn for future tasks. It is likely that the a commercial robot product that has the capability of the familiarizing task will include both human-interactive and automatic components (i.e. automatic object recognition). In this thesis we focus on the interactive components.

1.6 Challenges

There are four primary challenges in developing practical systems that allow a real robotic system to navigate using semantic information. From the robot’s perspective, the interesting research questions are:

- Where are the people?
- How should I move around humans?

- How can I acquire semantic information about this environment?
- How can I utilize the semantic information to improve navigation?

The answers to these questions are provided in the rest of this section.

1.6.1 Robust People Tracking

In order to be engaged in any kind of meaningful spatial interaction with people, the robot has to be able to track them. Robust tracking of people is necessary to ensure high task rates.

We use a person tracking system that fuses detections from three sources. These three detectors are: leg detector, torso detector and upper body detector. Usage of multiple modalities is aimed to increase the robustness of tracking and provide more coverage around the robot. The detections are consolidated into position estimations using a Kalman Filter.

1.6.2 Social Navigation

Robot path planning for navigation has been traditionally considered as a shortest-path problem. While this leads to correct and collision-free paths, the behavior of the robot may be perceived as unsafe or unnatural by human observers. When a robot is navigating in a human environment, the most important consideration is the safety of people. Moreover, even if a robot's motions are safe, they can still cause discomfort. For example, sudden appearance of a robot can cause fear in humans or cutting in between two people while they are in a conversation may be considered a rude behavior. Therefore, the robot motions should not only be effective but also be socially acceptable. Borrowing the idea from Sisbot [112], we embed personal spaces and social constraints as costs while planning a path. We further develop the planner by considering simulating people's reaction to robot's motion during planning. Humans routinely use anticipation during navigation. For example, in a

hallway encounter, one may lean towards the left of the corridor, expecting the other person to take right. We aim to give this anticipation ability to robots. Furthermore, we develop applications where the goal of the robot is to follow a person or a guide a person.

1.6.3 Interactive Map Annotation

The question of how semantic maps could be useful for various tasks is addressed in Section 1.5. Users may want to provide custom labels to landmarks, objects and waypoints instead of generic ones. By using unique labels, ambiguities can be resolved. For example, if an automatic object recognizer is utilized, the robot may be confused if it is given the command “bring my laptop” and there are multiple laptops in the environment.

Moreover, the ability to label objects and then refer to them relieves the robot the burden to have a strong recognizer. Without any object recognition capability, a user can point at and label an object and later refer to it while tasking. For example, an object that was previously labeled as “Joe’s mug” can be retrieved with a command such as “Fetch Joe’s mug from the kitchen table”.

A solution to enable a common ground between the user and the robot is via the Tour Scenario, as described in Section 1.5. For this scenario, the robot follows the user throughout the tour.

1.6.4 Situational Awareness for Navigation

Metric maps provide a single type of information: whether a cell is occupied with an obstacle or not. While this kind of a simplistic map representation may be sufficient for some tasks, if the robot needs to communicate with human users about its goals, it should have a richer map representation of the environment. The robot can leverage the semantic information embedded to the map and increase its usefulness for navigation tasks.

Situational awareness of the robot can be achieved in different ways. Person following is a widely studies behavior, however in the literature it is not usually considered *why* the robot is following the person and how the task context can be used to increase effectiveness. For example, during the Tour Scenario, the robot may employ a better following strategy if it knows the user will be pointing at landmarks and places. Furthermore, consider handling of the doors while following. If the robot does not know that it is standing on a doorway, it can block the path for others. If the door is added to the semantic map, the robot can have the opportunity to handle the door passing graciously. Moreover, the robot can utilize the map to adjust its speed. For example, if the robot is in a region where the chance to encounter a person is high, slowing down may be a sensible action to take.

1.7 Scope and Context

In this section, we provide an overview of the assumptions of our research. At a high level, autonomous robots perceive its environment, reason the situation, act in the world. In the HRI context, the research in perception is focused on improving the robot’s ability to “see” its environment, identify objects and humans. The reasoning part is focused on deciding on “what to do” and “how to do” in a given situation. Acting part is focused on the execution and control of motions from the reasoning part as well as the design of mechanisms. This thesis focuses on the reasoning part while also contributing to the perception part.

We assume that we have wheeled mobile robot platform. Even though some of our algorithms are derived for non-holonomic robots (differential drive), it is a technical limitation than a conceptual one, as they can be applied to holonomic robots with minor changes. We further assume that the robot is self-contained: equipped with on-board sensors, is capable of creating a map of the environment with an existing Simultaneous Localization and Mapping (SLAM) package, and then able to localize

itself in this created map as it moves in the environment.

Furthermore, we focus on scenarios where people in the environment are standing or walking, because our person tracking system is optimized for such cases.

The interaction between the robot and user is achieved via tablet or smartphone apps. Similar functionality could be achieved via a speech dialog system, however we chose touch-based interfaces due to their ubiquity and high reliability.

1.8 Thesis Overview

1.8.1 Thesis Statement

Non-expert users can label objects, landmarks and locations while giving a tour to a mobile robot, where the labeled entities enable situational awareness and serve as end goals, which robots can navigate to, by taking into account the personal spaces of nearby humans and social factors.

1.8.2 Contributions

- Development of a multimodal person detection and tracking system.
- A method for a human to take a mobile robot on a tour and interactively add labeled objects, landmarks and locations to the robot's map using natural deictic gestures and later use the labeled entities as end goals for navigation.
- Development of a navigation planner that takes into account social factors, and reactions of people to robot's motions.
- Development of a person guidance behavior and further investigation on the usage for blind users.
- Development of a person following behavior and further investigation on the usage for telepresence robots.

- Demonstration of situational awareness for person following behavior, targeted at the tour scenario.

1.8.3 Document Outline

CHAPTER II

MAP ANNOTATION

Robots that operate in human environments will need to be able to accept commands from human users. As mentioned in Chapter 1, the robot keeps a metric map in its own coordinate frame. Requiring users to understand the robot's representation of the world and provide metric commands may not be an interaction with the robot. For example, it is not easy to understand look at the map of the robot and provide an explicit goal with coordinates (1.2, 4.5, 0.0). The world representation of robots should enable communication between users and robots. Maps should also include spatial information to support the tasks. If the task is more complex than obstacle avoidance, the map should include more than just obstacle occupancy information. For example, if the tasks involve landmarks, objects, or people, a robot's map should be able to represent these information. For these reasons, adding semantic information to robot maps is essential.

If the map representation includes more than explicit coordinates, such as landmarks and objects, there is a need to establish a common ground between the robot and the users to the all can refer to the same entity. There are several methods to support the annotation. For example, while the robot is building its representation of the environment, it can recognize objects or landmarks such as doors, tables, rooms and automatically add these features to its map. However, such a system may result in ambiguities. For example, recognition errors may result in undesirable outcomes. Moreover, automatic recognition does not allow unique labels, such as "*John's Room*". Attaching custom labels to landmarks not only serves a common ground between the robot and users, but also frees the robot from having to be equipped with object and

landmark recognition capability.

We use a map representation that includes 2D waypoints, 3D polygons that represents the boundary of surfaces and objects along with a corresponding label that is provided by the user. We give more detail about waypoints in Section 2.2.1, planar surfaces in Section 2.2.2 and objects in Section 2.2.3.

To label landmarks and objects in the semantic map, we use an interactive system that uses a combination of a graphical user interface (GUI) and pointing gestures. Using the GUI, users can command the robot to follow them, add waypoints along the robot’s path and add landmarks or objects to the map by pointing at them and entering a label. Our GUI is shown in Section 2.3 and our work in pointing gestures is given in Section 2.4.

2.1 Related Work

One of the key concepts in semantic mapping is establishing *common ground* [19]. For humans and robots to refer to the same structures and objects, the references must be grounded. Many spatial tasks may require various terms to be grounded in the map.

TODO

2.2 Semantic Maps

Semantic mapping aims to create maps that include various types of semantic information to allow the robot to have a richer representation of the environment. In this section, we briefly present the semantic information we use to enrich the occupancy grid maps: waypoints, planar landmarks and objects. Although we describe three types of storing semantic information, one can come up with richer features and landmarks.

2.2.1 Navigation Waypoints

We define navigation waypoint as a navigable point represented in the map frame. We will refer navigation waypoint simply as a waypoint for the remainder of this text. After the robot is initialized, the user can save the current pose of the robot along with a label. This is done using the on-board Graphical User Interface (GUI) of the robot. This method is the most straightforward way to save a waypoint and use it later for a task. When the robot is asked to navigate to a labeled waypoint, the goal of the robot is simply the saved pose.

2.2.2 Planar Landmarks

Even though the waypoint method is easy to use and practical, it has shortcomings. It fails to capture the shape and extent of the structure designated by the label, which might be important for some tasks, such as object search. Moreover, point based references may be ambiguous to represent a region or volume in a map, such as a hallway or a room. For example, if the user wants to save hallway as a landmark, it is useful to have a representation of the location and the extent of the hallway. Another example is that the user can label a coffee table, which is a movable object. Instead of saving a fixed coordinate location for the landmark, robot can save the planar surface and can potentially detect that it moved. This gives the robot robustness in its operations and a method to potentially use the planar surfaces as features in localization [122]. For this reason, using semantic information is preferred whenever possible.

In this type of semantic information, we keep track of a set of observed planar surfaces as part of the map representation. Planes represented in the hessian normal form. Since the observed planes do not have infinite extent, we bound the observed area with a polygon - in this case, the convex hull of all the observed points on the plane. The convex hull is accompanied with a user-provided label.

We use a front-facing RGB-D camera to acquire suitable data. Planes are extracted from the point cloud by an iterative RANdom SAmple Consensus (RANSAC) method, which allows us to find all planes meeting constraints for size and number of inliers. A clustering step is performed on extracted planes to separate multiple coplanar surfaces, such as two tables with the same height, but at different locations. We make use of the Point Cloud Library (PCL) [100] for much of our point cloud processing.

In order to label a planar landmark, the user goes through the sequence:

1. User activates person following, goes nearby the planar landmark of interest (table)
2. Robot enters the label using UI, activates pointing gesture recognition
3. User performs a pointing gesture towards the landmarks of interest

The act of labeling a table from the view of the RGB-D camera planted on the robot is shown in Figure 1. After the labeling process, the planar landmark is added to the map representation.

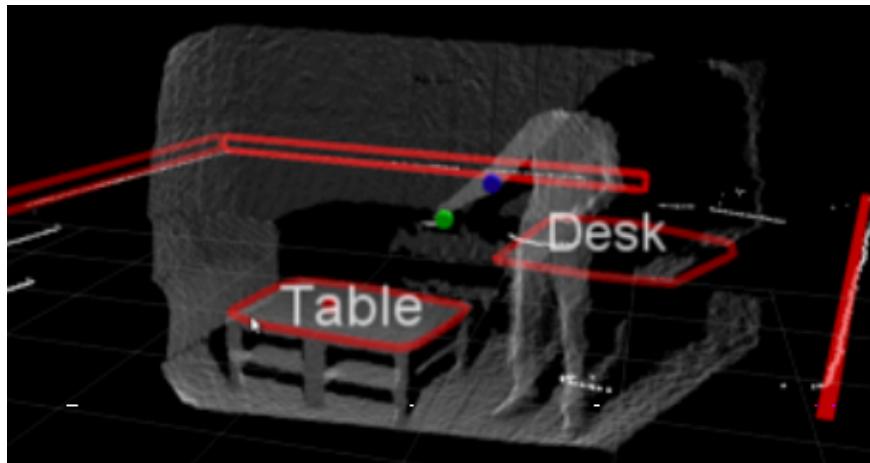


Figure 1: Robot's view of how users can add planar landmarks to the map

2.2.3 Objects

Our approach allows to interactively build object models for selected objects, and then annotate these models with a label. This enables the user to easily command the robot to build a small object database of the specific objects it needs to interact with. Note that the alternative, automatic recognition of objects, require robust detection and a capability to reliably recognize a large number of objects.

The interactive system works as described in the previous section. Apart from labeling planar landmarks, however, we need to do a little bit more work. Once a user has given a referenced point at which to detect objects, the robot moves its pan-tilt unit to aim the sensors at this location. Object segmentation is performed using point cloud data from an RGB-D sensor. First, the large planar surface corresponding to the table is detected. This is removed from the point cloud, and point clusters above this are detected. The cluster with a centroid nearest to the reference point is selected as the object to be modeled. The clusters points are projected into our camera image, and are used to generate a region of interest. SURF features are detected for the region of interest, and are stored as an object model along with the provided label.

2.3 User Interface

We developed an on-board tablet interface to enable interaction between the user and the robot. Similar functionality could be achieved with a speech recognition system, however we think a touch interface is a more robust type of communication. The interface has been implemented on a Nexus 7 Tablet running on Android Operating System. The tablet is equipped with WiFi and talks with the robot using a TCP-based server-client communication model. The client must have the the IP address of the server machine to establish the connection. The messages are sent back-and forth as via XML files. The XML files are serialized and de-serialized using *limXML++*

for *C++* and *XMLPullParser* for Android OS. With the tablet interface, a user can do the following:

- Add/Delete Waypoint
- Label a Landmark or Object
- Start Person Following
- Start Person Guidance to a Labeled Landmark/Object/Waypoint
- Stop Robot
- Modify IP Address and Port of the TCP connection

A screenshot of our tablet user interface is shown in Figure 2.

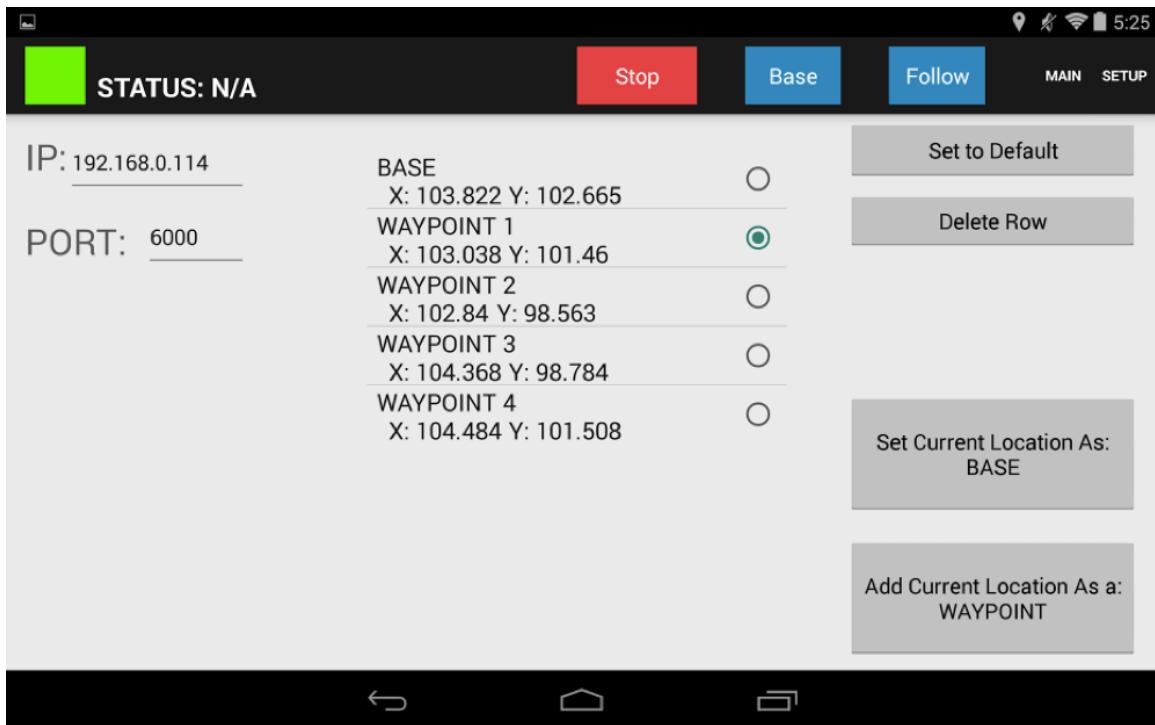


Figure 2: UI for labeling a navigation waypoint

2.4 Pointing Gestures for Human-Robot Interaction

A natural way to refer to objects and landmarks to label them is to use deictic gestures. In this work, we analyze the performance of our deictic gesture recognition, and present an uncertainty model that enables us to reason about ambiguity in pointing gestures, when objects are too close to one another to determine the referent. We model the uncertainty of pointing gestures in a spherical coordinate system, use this model to determine the correct pointing target, and detect when there is ambiguity. Two pointing methods are evaluated using two skeleton tracking algorithms: elbow-hand and head-hand rays, using both OpenNI NITE and Microsoft Kinect SDK [106]. A data collection with 6 users and 7 pointing targets was performed, and the data was used to model users pointing behavior. The resulting model was evaluated for its ability to distinguish between potential pointing targets, and to detect when the target is ambiguous. An example scenario is shown in Figure 3.

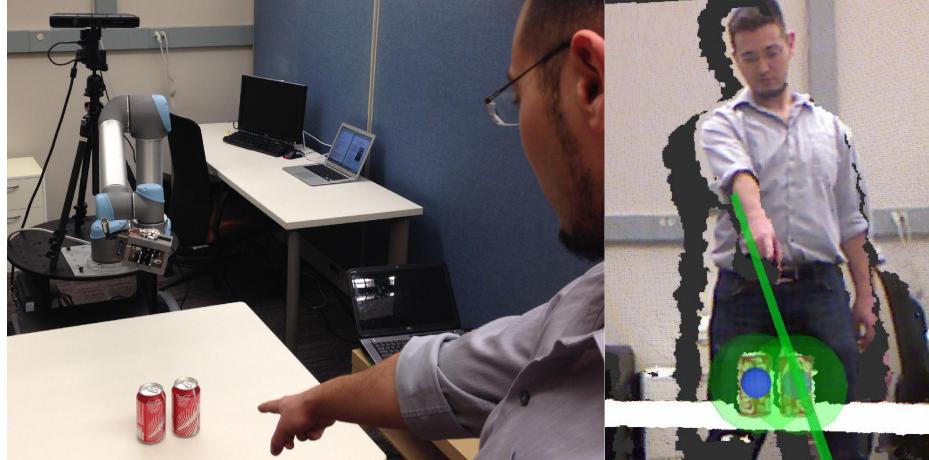


Figure 3: (Left) Our approach allows a robot to detect when there is ambiguity on the pointing gesture targets. (Right) The point cloud view from robot’s perspective is shown. Both objects are identified as potential intended targets, therefore the robot decides that there is ambiguity.

2.4.1 Related Works

Pointing gestures are widely used in Human-Robot Interaction applications. Examples include interpretation of spoken requests [135], pick and place tasks [11], joint attention [24], referring to places [37] and objects [104], instructing [71] and providing navigation goals [96] to a robot.

Early works on recognizing pointing gestures in stereo vision utilized background subtraction [18, 44, 45]. Other popular methods include body silhouettes [47], hand poses [42], motion analysis [73] and Hidden Markov Models [128, 8, 65, 83, 23, 2]. Matuszek *et. al* presented a method for detecting deictic gestures given a set of detected tabletop objects, by first segmenting the users hands and computing the most distal point form the user, then applying a Hierarchical Matching Pursuit (HMP) approach on these points over time [74].

After deciding if a pointing gesture occurred or not, an algorithm must estimate the direction of pointing. This is typically done by extending a ray from one body part to another. Several body part pairs are used in the literature, such as eye-fingertip [47] and shoulder-hand [41]; with the two of most commonly used methods being elbow-hand [96, 12, 11] and head-hand [8, 104] rays. Some studies found that head-hand approach is a more accurate way for estimating pointing direction than elbow-hand [95, 24]. Recent works made use of skeleton tracking data from in depth images [95, 11, 96]. Other approaches, such as measuring head orientation with a magnetic sensor [83] and pointing with a laser pointer [16, 48] is reported to have a better estimation accuracy than the body parts method, but require additional hardware. We prefer not to use additional devices in order to the interaction as natural as possible.

Given a pointing direction, several methods have been proposed to determine which target or object is referred by the gesture, including euclidean distance on a planar surface [16], ray-to-cluster intersection in point clouds [11, 95] and searching

a region in interest around the intersection point [104]. Droeuschel [23] trains a function using head-hand, elbow-hand and shoulder-hand features with Gaussian Process Regression and reports a significant improvement on pointing accuracy. Some efforts fuse speech with pointing gestures for multi-modal Human-Robot Interaction [2, 58]. Aly [2] focuses on relation between non-verbal arm gestures and para-verbal communication based on a HMM approach.

To our knowledge, only Zukerman [136] and Kowadlo [58] considered a probabilistic model for determining the referred object for a pointing gesture. The probability that the user intended an object is calculated using the 2D distance to the object and the occlusion factor. Objects that reside in a Gaussian cone emanating from the user’s hand are considered as candidates in the model. The approach is implemented in [136], where it is reported that due to high variance of the gesture recognition system, the Gaussian cone typically encompassed about five objects in cluttered settings. Our work addresses the confusion in such settings. In contrast to their work, we use a 3D elliptical Gaussian cone where its shape is extracted using a prior error analysis, and use point cloud data to account for the size of the objects.

2.4.2 Pointing Gesture Recognition

Our approach to pointing gesture recognition is to use a skeleton tracking algorithm as implemented by OpenNI NITE 1.5 (OpenNI) or Microsoft Kinect SDK 1.5 (MS-SDK). Skeleton tracking software produces 3D positions for several important points on the user’s body, including hands, elbows, shoulders and head. Our points of interests are user’s hands, elbows, and head for the recognition of pointing gestures.

We are primarily interested in deictic gestures generated by pointing with one’s arm. We consider two rays for determining the pointing direction: elbow-hand and head-hand. Both of these methods were evaluated with the two skeleton tracking implementations. For each depth frame, this yields two rays for each of the OpenNI

and MS-SDK trackers:

- $\vec{v}_{eh} := \vec{p}_{elbow}\vec{p}_{hand}$

- $\vec{v}_{hh} := \vec{p}_{head}\vec{p}_{hand}$

When a pointing gesture recognition request is received from a higher level process, the gesture is searched in a time window of T seconds. Two conditions must be met to trigger a pointing gesture:

- Forearm makes an angle more than ϕ_g with the vertical axis
- Elbow and hand points stays near-constant for duration Δt_g

The first condition requires the arm of the person to be extended, while the second ensures that the gesture is consistent for some time period. The parameters are empirically determined as: $T = 30s$, $\phi_g = 45^\circ$ and $\Delta t_g = 0.5s$.

2.4.3 Representing Pointing Directions

We represent a pointing ray in two angles: a “horizontal” / “azimuth” sense we denote as θ and a “vertical” / “altitude” sense we denote as ψ . We first attach a coordinate frame to the hand point, with its z-axis oriented in either Elbow-hand \vec{v}_{eh} or Head-Hand \vec{v}_{hh} directions. The hand was chosen as the origin for this coordinate system because both of head-hand and elbow-hand pointing methods include the user’s hand. The transformation between the sensor frame and the hand frame ${}^{sensor}T_{hand}$ is calculated by using an angle-axis rotation. An illustration of the hand coordinate frame for Elbow-Hand method and corresponding angles are shown graphically in Figure 4.

Given this coordinate frame and a potential target point P, we first transform it to the hand frame by:

$${}^{hand}p_{target} = T_{hand} * p_{target}$$

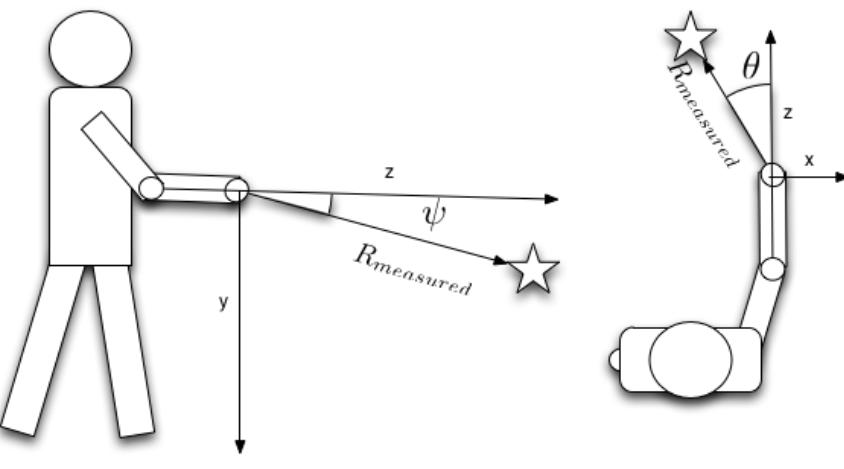


Figure 4: Vertical (ψ) and horizontal (θ) angles in spherical coordinates are illustrated. A potential intended target is shown as a star. The z-axis of the hand coordinate frame is defined by either the Elbow-Hand (this example) or Head-Hand ray.

We calculate the horizontal and vertical angles for a target point as ${}^{hand}p_{target} = (x_{targ}, y_{targ}, z_{targ})$ follows:

$$[\theta_{target} \psi_{target}] = [\text{atan2}(x_{targ}, z_{targ}) \text{ atan2}(y_{targ}, z_{targ})]$$

Where $\text{atan2}(y, x)$ is a function returns the value of the angle $\arctan(\frac{y}{x})$ with the correct sign. This representation allows us to calculate the angular errors in our error analysis experiments in Section 2.4.9. The angles for each object is then used to find the intended target, as explained in the following section.

2.4.4 Determining Intended Target

We propose a probabilistic approach to determine the referred target by using statistical data from previous pointing gesture observations. We observed that head-hand and elbow-hand methods, implemented using two skeleton trackers, returned different angle errors in spherical coordinates. Our approach relies on learning statistics of each of these approaches, and compensating the error when the target object is searched for. First, given a set of prior pointing observations, we calculate the mean and variance of the vertical and horizontal angle errors for each pointing method.

This analysis will be presented in Section 2.4.9. Given an input gesture, we apply correction to the pointing direction and find the Mahalanobis distance to each object in the scene.

When a pointing gesture is recognized, and the angle pair

$$[\theta_{target} \psi_{target}]$$

is found as described in the previous section, we first apply a correction by subtracting the mean terms from measured angles:

$$[\theta_{cor} \psi_{cor}] = [\theta_{target} - \mu_\theta \psi_{target} - \mu_\psi]$$

We also compute a covariance matrix for angle errors in this spherical coordinate system:

$$S_{type} = \begin{bmatrix} \sigma_\theta & 0 \\ 0 & \sigma_\psi \end{bmatrix}$$

We get the values for $\mu_\theta, \mu_\psi, \sigma_\theta, \sigma_\psi$ from Tables 1 and 2 for the corresponding gesture type and target. We then compute the mahalanobis distance to the target by:

$$D_{mah}(target, method) = \sqrt{[\theta_{cor} \psi_{cor}]^T S_{method}^{-1} [\theta_{cor} \psi_{cor}]}$$

We use D_{mah} to estimate which target or object is intended. We consider two use cases: the objects are represented as a 3D point or a point cloud. For point targets, we first filter out targets that have a Mahalanobis distance larger than a threshold $D_{mah} > D_{thresh}$. If none of the targets has a D_{mah} lower than the threshold, then we say the user did not point to any targets. If there are multiple targets that has $D_{mah} <= D_{thresh}$, then we determine ambiguity by employing a ratio test. The ratio of the least D_{mah} and the second-least D_{mah} among all targets is compared with a threshold to determine if there is ambiguity. If the ratio is higher than a threshold, then the robot can ask the person to solve the ambiguity.

If the target objects are represented as point clouds, we then compute the horizontal and vertical angles for every point in the point cloud and find the minimum mahalanobis distance among all. The distance to an object is then represented by this minimum value. Usage of the point cloud instead of the centroid for determining the intended object has several advantages. First, it yields better estimations due to the coverage of the full point cloud. Second, it takes into account the size of the object. For example, if a person is pointing to a chair or door, it is very unlikely that he/she will target the centroid of it. If the point cloud is used, then we can easily tell that the object is targeted.

2.4.5 Data Collection

To evaluate the accuracy of pointing gestures, we created a test environment with 7 targets placed on planar surfaces in view of a Kinect sensor (Figure 5). Depth data was collected from six users, who pointed at each of the seven targets with their right arm while standing at 2 meters away from the sensor. Targets 1 through 4 were on tables positioned around the user, while targets 5 through 7 were located on a wall to the user’s right. Our use case is on a mobile robot platform capable of positioning itself relative to the user. For this reason, we can assume that the user is always centered in the image, as the robot can easily rotate to face the user and can position itself at a desired distance from the user.

2.4.6 Ground Truth Target Positions

We make use of plane extraction technique in a point cloud to have an accurate ground truth measurement. First, the two table and wall planes are extracted from the point cloud data using the planar segmentation technique described in [121]. We then find the pixel coordinates of corners on targets in RGB images, using Harris corner detection [36], which produces calculated corners in image coordinates with sub-pixel accuracy. The pixel coordinate corresponding to each target defines a ray in

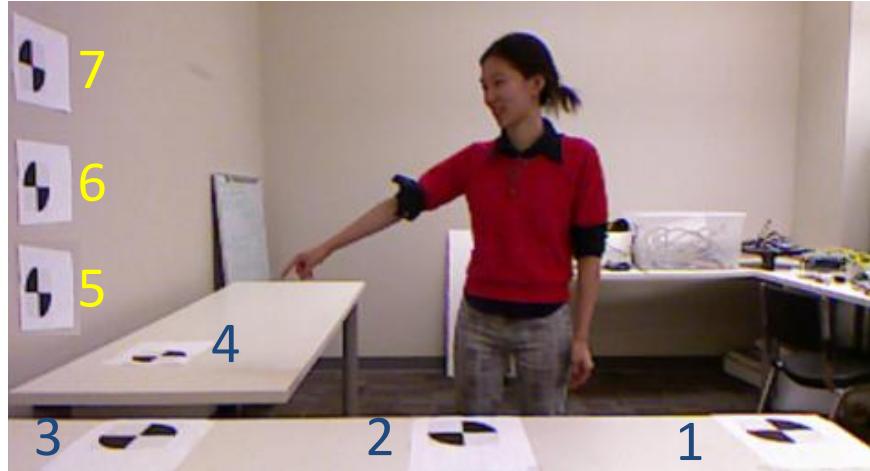


Figure 5: Our study involved 6 users that pointed to 7 targets while being recorded using 30 frames per target.

3D space relative to the Kinect's RGB camera. These rays are then intersected with the planes detected from the depth data, yielding the 3D coordinates of the targets.

2.4.7 Skeleton Data Capture

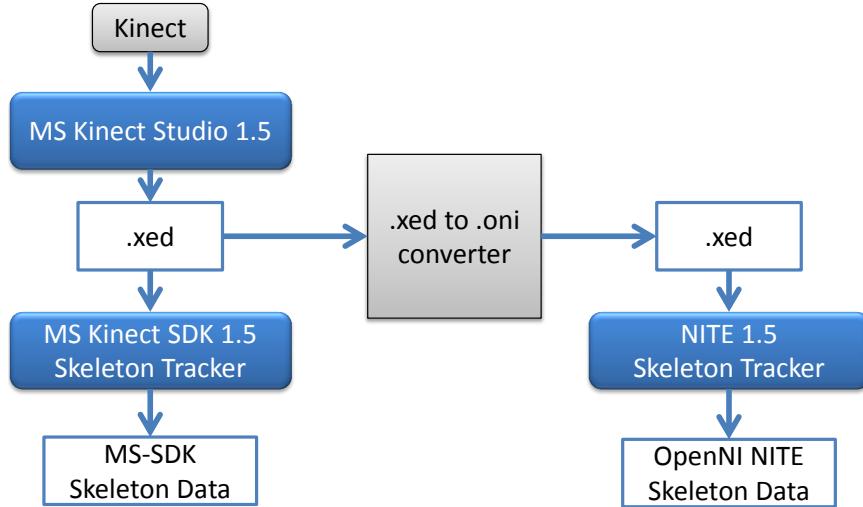


Figure 6: Data capturing pipeline for error analysis.

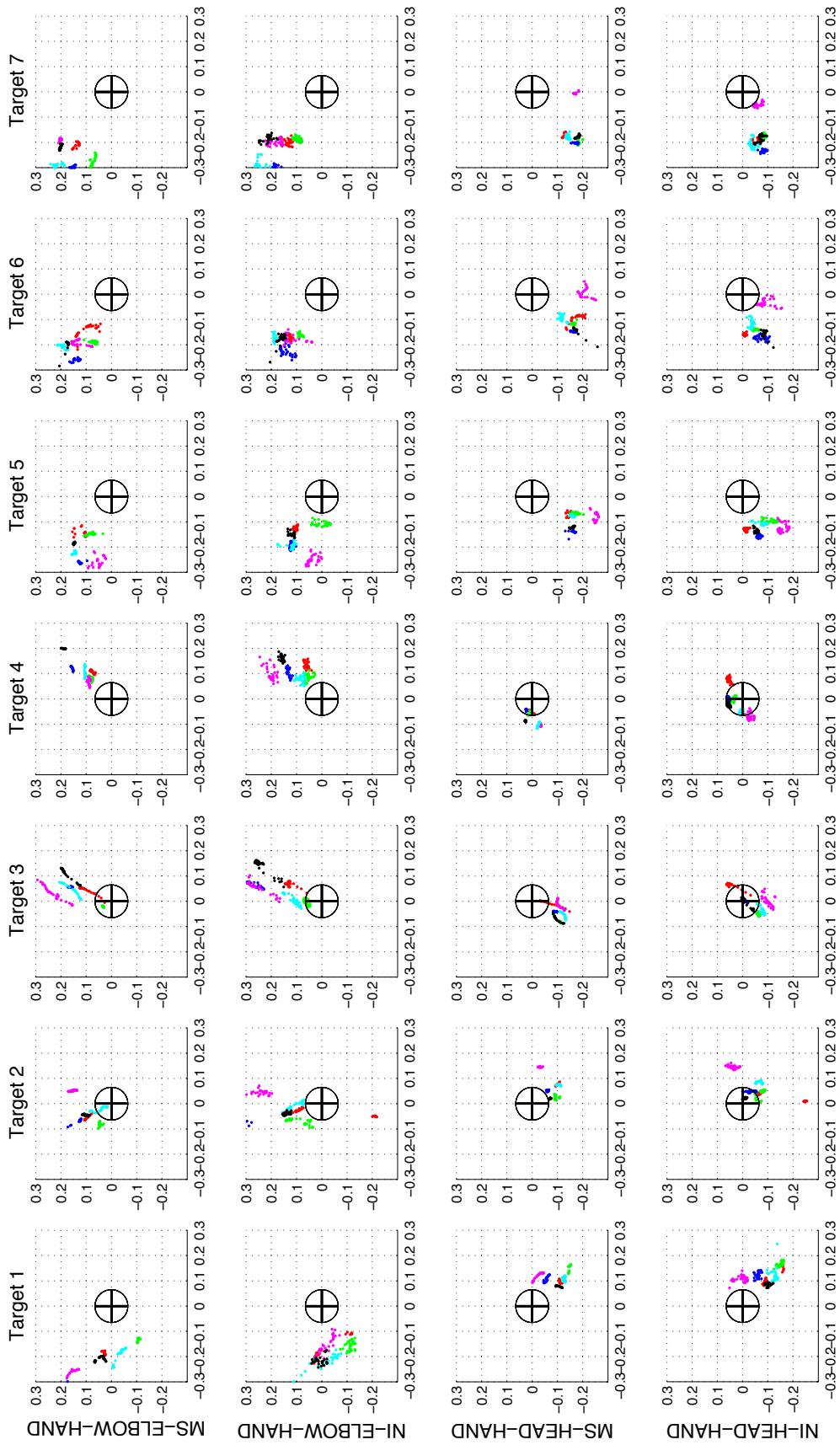


Figure 7: Euclidean distance error in cartesian coordinates for each method and target. The gesture ray intersection points in centimeters with the target plane, are shown here for each target (T1-T7) as the columns. Each subject's points are shown in separate colors. There are 30 points from each subject, corresponding to the 30 frames recorded for the pointing gesture at each target. Axes are shown in centimeters. The circle drawn in the center of each plot has the same diameter (13 cm) as the physical target objects used.

In order to be able to do a fair comparison between MS-SDK and OpenNI skeleton trackers, we used the same dataset for both. MS-SDK and OpenNI use different device drivers, therefore can not be directly used on the live depth stream at the same time. Because of this, we extract the skeleton data offline in multiple stages. The pipeline for the data capture procedure is illustrated in Figure 6. We first save the depth streams as .xed files using Microsoft Kinect Studio program. The acquired .xed file is converted to .oni in a OpenNI recorder application by streaming the depth stream to through Kinect Studio. The .oni is then played back in skeleton tracking application in OpenNI, which writes the OpenNI skeleton data to a .txt file. MS-SDK skeleton data is obtained by playing back the original .xed in the skeleton tracking application.

2.4.8 Pointing Gesture Annotation

To factor out the effectiveness of our pointing gesture recognition method described in Section 2.4.2, we manually labeled when each pointing gesture began for data collection. Starting from the onset of the pointing gesture as annotated by the experimenter, the following 30 sensor frames were used as the pointing gesture. This corresponds to a recording of 1 second in the Kinect sensor stream. For each frame, we extracted skeleton data using both the MS-SDK and the OpenNI.

2.4.9 Error Analysis

The four rays corresponding to the four different pointing approaches described in Section 2.4.2 were used for our error analysis. As described in Section 2.4.6, the ground truth target positions are available. We computed two types of errors for each gesture and target:

- Euclidean error of ray intersections with target planes (Figure 7)
- Angular error in spherical coordinates (Tables 1,2 and Figure 8)

We elaborate our error analysis subsequent sections:

	Target 1			Target 2			Target 3			Target 4		
	θ	ψ										
	μ	σ										
MS-Elbow-Hand	-15.7	2.9	5.5	6.1	-4.3	6.6	7.8	3.1	-3.7	2.4	7.4	3.1
NI-Elbow-Hand	-16.4	2.9	4.3	7.0	-3.8	6.6	11.3	10.9	-4.8	2.6	9.7	3.4
MS-Head-Hand	7.7	2.6	-12.0	5.3	10.8	6.4	-9.1	3.2	2.2	2.0	-8.3	3.2
NI-Head-Hand	8.5	2.6	-11.7	6.1	10.2	6.7	-5.7	8.0	2.0	2.3	-2.9	4.7

Table 1: μ and σ angular errors in degrees for each of Targets 1-4 (Figure 5), for each pointing method. The aggregate μ and σ is also shown.

	Target 5			Target 6			Target 7			ALL TARGETS		
	θ	ψ	θ	ψ								
	μ	σ	μ	σ								
MS-Elbow-Hand	-14.5	4.1	11.6	3.7	-16.6	3.3	9.9	3.7	-20.8	3.7	5.7	3.4
NI-Elbow-Hand	-12.9	4.2	9.7	5.1	-16.2	2.6	11.7	3.7	-20.2	4.5	8.1	-3.0
MS-Head-Hand	-9.4	1.9	-11.6	3.0	-7.1	4.7	-13.8	1.8	-8.0	5.5	-15.6	-2.9
NI-Head-Hand	-11.5	1.5	-4.7	4.8	-11.2	4.8	-4.9	2.9	-12.3	5.2	-8.9	-2.5

Table 2: μ and σ of angular error in degrees for Targets 5-7 (Figure 5), for each pointing method. The aggregate μ and σ is also shown.

2.4.9.1 Euclidean error

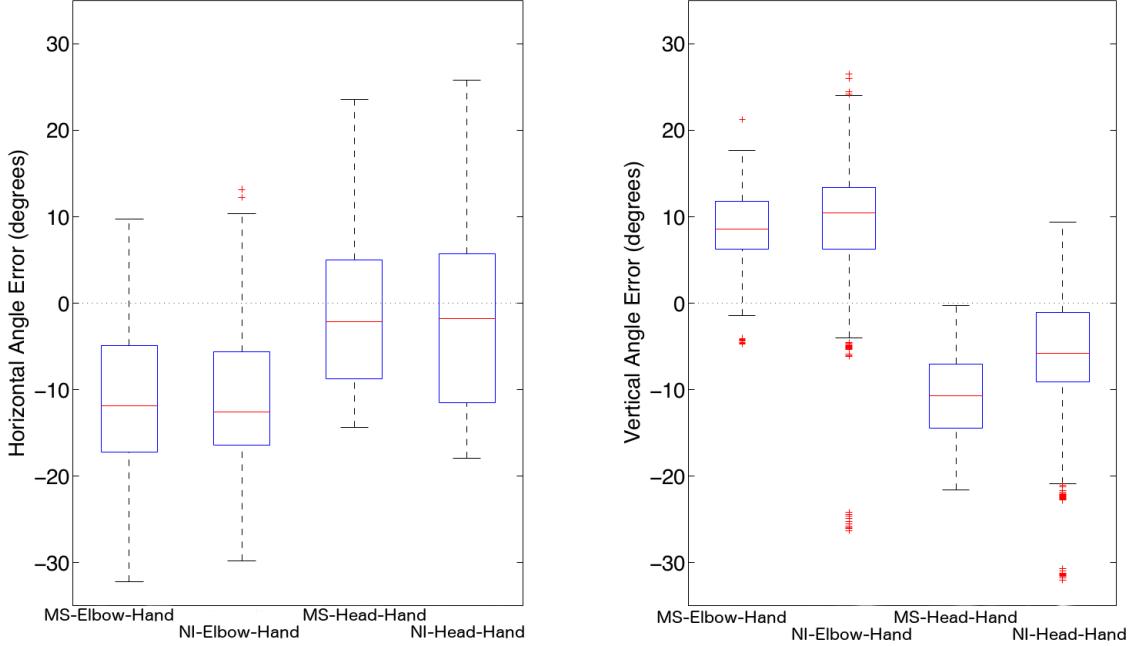


Figure 8: Box plots of the errors in spherical coordinates θ and ψ for each pointing method.

Given a ray \vec{v} in the sensors frame from one of the pointing gesture approaches, and a ground truth target point p_{target} lying on a target planar surface, the ray-plane intersection between \vec{v} and plane was computed for each ray, resulting in a 3D point lying on the plane. Figure 7 shows the 2D projections for all 30 measurements from each subject (shown in different colors) and each target. For ease of display, the 3D intersection coordinates with the target planes are displayed in a 2D coordinate system attached to the target plane, with the ground truth target point as the origin.

As can be seen in Figure 7, the pointing gesture intersection points were fairly consistent across all users, but varied per target location. The elbow-hand method produced similar results for MS-SDK and OpenNI. The same is true for the head-hand method. It is also noteworthy that the data for each target tends to be quite clustered for all methods, and typically not centered on the target location.

2.4.9.2 Angular Error

We computed the mean and standard deviations of the angular errors in the spherical coordinate system for each pointing gesture method and target. Section 2.4.3 describes in detail how the angles (θ, ψ) are found. The mean and standard deviation values are given in Tables 1 and 2. The aggregate errors are also displayed as a box plot in Figure 8.

Several observations can be made from these results. The data from the elbow-hand pointing method reports that users typically point about 11 degrees to the left of the intended target direction, and about 9 above the target direction. Similarly, the data from the head-hand pointing method reports that users typically point about 2 degrees to the left of the intended pointing direction, but with a higher standard deviation than the elbow-hand method. On average, the vertical angle ψ was about 5 degrees below the intended direction for the OpenNI tracker and 10 degrees below for the MS-SDK tracker, with a higher standard deviation than the elbow-hand methods. The disparity between the two skeleton trackers for this pointing method is because they report different points for the head position, with the MS-SDK head position typically being reported higher than the OpenNI head position. The overall performance of the OpenNI and MS-SDK skeleton trackers, however, is fairly similar, with the MS-SDK having slightly less variation for our dataset.

As can be seen in the aggregate box plot in Figure 8, the horizontal angle θ has a higher variation than the vertical angle ψ . Examining the errors for individual target locations shows that this error changes significantly with the target location. As future work, it would be interesting to collect data for a higher density of target locations to attempt to parameterize any angular correction that might be applied.

2.4.10 Evaluation

Using the error analysis and pointing gesture model we presented in previous sections, we conducted an experiment to determine how our approach distinguished two potentially ambiguous pointing target objects. We use the results of the angular error analysis results and not the euclidean error analysis for the remainder of the paper because of our angular representation of pointing gestures.

2.4.11 Object Separation Test

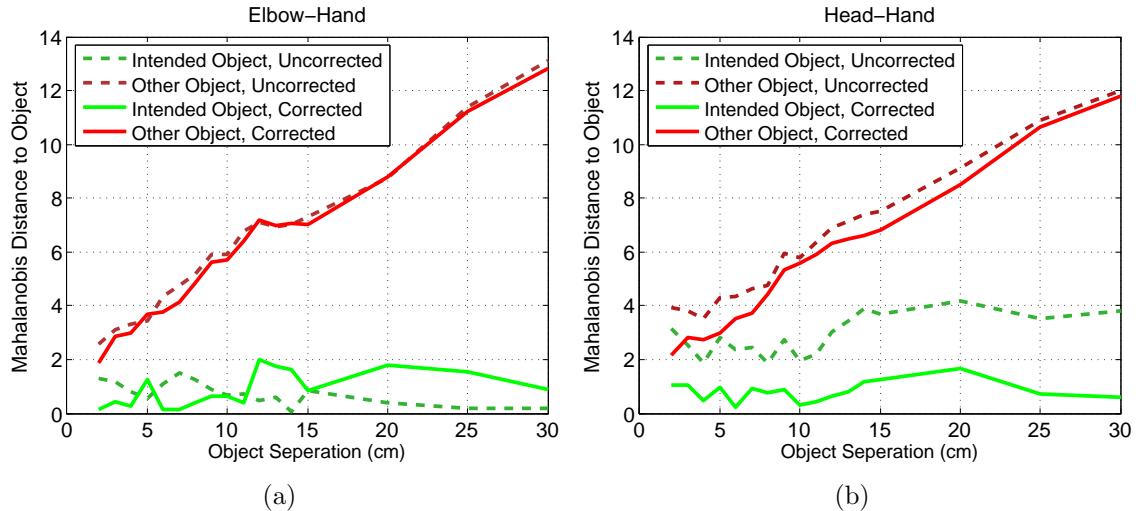


Figure 9: Resulting Mahalanobis distances of pointing targets from the Object Separation Test is shown for a) Elbow-Hand and b) Head-Hand pointing methods. Intended object are shown in green and other object is in red. Solid lines show distances after correction is applied. Less Mahalanobis distance for intended object is better for reducing ambiguity.

The setup consisted of a table between the robot and the person and two coke cans on the table (Figure 10) where the separation between objects was varied. To detect the table plane and segment out the objects on top of it, we used the segmentation approach presented in [121]. The objects centroid positions, along with their point clouds were calculated in real-time using our segmentation algorithm. The separation between objects were varied with 1 cm increments from 2 – 15cm and with 5cm increments between 15 – 30cm. We could not conduct the experiment below 2 cm

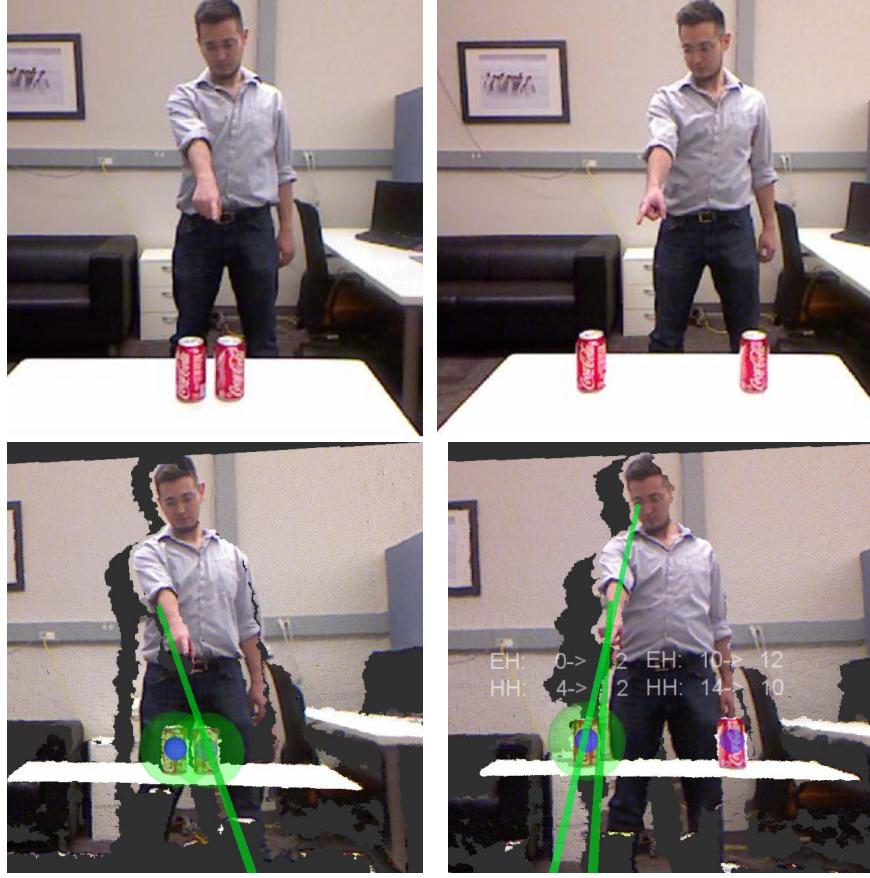


Figure 10: Example scenarios from the object separation test is shown. Our experiments covered separations between 2cm (left images) and 30cm (right images). The object is comfortably distinguished for the 30cm case, whereas the intended target is ambiguous when the targets are 2cm apart. Second row shows the point cloud from Kinect’s view. Green lines show the Elbow-Hand and Head-Hand directions whereas green circles show the objects that are within the threshold $D_{mah} < 2$.

separation because of the limitations of our perception system. We used the OpenNI skeleton tracker because rest of our system is based in Linux, and we already found that performance difference with MS-SDK for pointing angle errors is insignificant. The experiment was conducted with one user, who was not in the training dataset. For each separation, the user performed 5 pointing gestures to the object on the right and 5 to the object on the left. The person pointed to one of the objects and the Mahalanobis distance D_{mah} to the intended object and the other object is calculated using the approach in Section 2.4.4. We used the mean and standard deviation values

of Target 2 (Figure 5) for this experiment because the objects were located between the robot and the person.

2.4.12 Results and Discussion

The results of the object separation experiment is given for Elbow-Hand (Figure 9(a)) and Head-Hand (Figure 9(b)) methods. The graphs plot object separation versus the Mahalanobis distance for the pointed object and the other object for corrected and uncorrected pointing direction. There are several observations we make by looking at these results.

First, the Mahalanobis distance D_{mah} for the intended object was always lower than the other object. The corrected D_{mah} for both Elbow-Hand and Head-Hand methods for the intended object was always below 2, therefore selecting the threshold $D_{thres} = 2$ is a reasonable choice. We notice that some distances for the unintended object at $2cm$ separation is also below $D_{mah} < 2$. Therefore, when the objects are $2cm$ apart, then the pointing target becomes ambiguous for this setup. For separations of $3cm$ or more, D_{mah} of the unintended object is always over the threshold, therefore there is no ambiguity.

Second, correction significantly improved Head-Hand accuracy at all separations, slightly improved Elbow-Hand between $2 - 12cm$ but slightly worsened Elbow-Hand after $12cm$. We attribute this to the fact that the angles we receive is heavily user-dependent and can have a significant variance across methods as showed in Figure 8. Moreover, the user was not in the training set.

Third, the Mahalanobis distance stayed generally constant for the intended object, which was expected. It linearly increased with separation distance for the other object.

Fourth, patterns for both methods are fairly similar to each other, other than Head-Hand uncorrected distances being higher than Elbow-Hand.

2.4.13 Conclusions

Humans and robots engaged in joint actions that involve objects will need to be able to communicate about these objects. Deictic gestures can play a crucial part such communication, especially when establishing goals and negotiating responsibilities. A pointing target detection approach that returns a single hypothesis can lead to failures due to perception error. On the other hand, estimation of a pointing likelihood of nearby objects can give valuable insight to a robot. For example, a robot can ask the user to disambiguate the objects if it perceives more than one object that has a significant likelihood of being referred to.

In this work, we model the uncertainty of pointing gestures in a spherical coordinate system, use this model to determine the correct pointing target, and detect when there is ambiguity. Two pointing methods are evaluated using two skeleton tracking algorithms: Elbow-Hand and Head-Hand rays, using both OpenNI NITE and Microsoft Kinect SDK. A data collection with 6 users and 7 pointing targets was performed, and the data was used to model users pointing behavior. The resulting model was evaluated for its ability to distinguish between potential pointing targets, and to detect when the target is ambiguous. Our evaluation showed that in a scenario where the separation between two objects were varied, our system was able to identify that there is ambiguity for 2 cm separation and comfortably distinguished the intended object for 3 cm or more separation.

CHAPTER III

NAVIGATION AMONG PEOPLE

Autonomous navigation is one of the most fundamental tasks for a mobile robot. For a mobile robot with adequate actuation and sensing, collision-free navigation is considered a solved problem. There are many algorithms that achieve point-to-point autonomous navigation thanks to the advances in the motion planning community. Many of these algorithms are optimized to find the least-cost path, or the shortest path. However, when there are humans in the environment, such algorithms suddenly become inefficient or insufficient. For example, while it is acceptable for a robot to get inches close to a wall, doing so to a human is socially unacceptable and potentially dangerous. Similarly, sudden appearance of a robot can surprise or shock humans. There are many other social scenarios where the shortest path may not be optimal.

In addition to sub-optimality, these approaches may be incomplete in the sense that they can not find a solution even though there is a feasible one. This is because shortest-path navigation algorithms treat every object in the environment as an obstacle. This assumption does not hold when intelligent agents are present in the environment. Therefore navigation should differentiate humans and obstacles for more intelligent robot behavior.

Another aspect to spatial interaction between humans and robots is the dynamics of the robot motion. For example, people may feel uncomfortable and unsafe when they are in close proximity to high-speed agents or objects. Therefore, for a robot in a human environment, while it may be acceptable to speed up in dedicated regions, its speed should be limited in places where there is a significant possibility of encountering a human.

In this Chapter, we first provide a background and present the most common approach in contemporary autonomous navigation methods in Section 3.1. Second, we provide relevant works on navigation among people in Section 3.2. Third, in Section 3.3, we present how the goal points for navigation are determined. We then present our people-aware navigation method in Section 3.4. Lastly, we touch to the subject of introducing speed limits for all robots in a human environment in Section 3.5.

3.1 State-of-the-Art Approach in Autonomous Navigation

There are two prerequisites that enables autonomous navigation:

1. The map of the environment, usually in the form of a discrete grid, that represents static objects in the environment
2. A way to localize the robot in the map using sensory information as it moves in the environment

Robot navigation involves finding a collision-free path from a start pose (x_0, y_0, θ_0) to a goal pose (x_g, y_g, θ_g) . In real-time operation, (x_0, y_0, θ_0) is the robot's current pose as the robot tries to reach to the goal pose from where it currently is. θ_g is optional as the goal of the robot could be to reach the goal position regardless of its orientation. The goal position is provided from an external process, and we will touch upon how the goal positions are calculated in Section 3.3.

A common approach to path planning is to divide the path planning into two parts: *global* and *local*. Global planning aims to find a path from the start position to the goal position. The global path is a set of consecutive positions that connect the start to goal position. A global path is usually found with a search algorithm executed on a graph of points. The search heuristics is dependent on specific global planners, and in most cases collision-free shorter paths are favored. The local planner is responsible

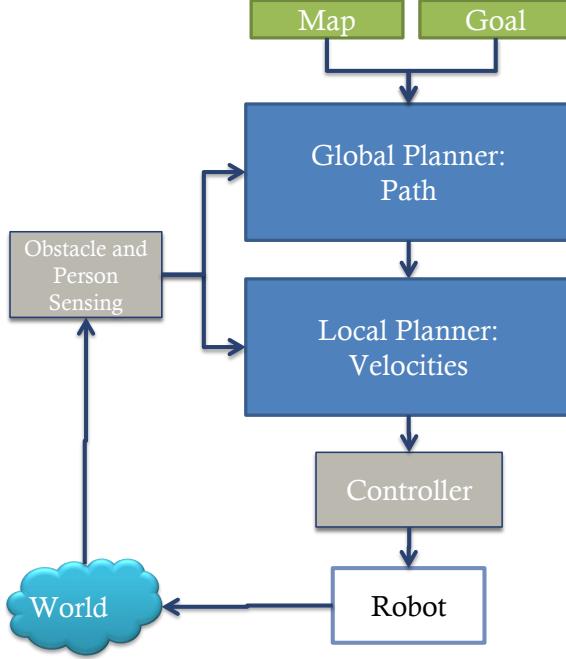


Figure 11: Caption

to execute the global path by calculating a trajectory and sending velocity commands to motor controllers. As the robot acts in the environment, its sensors sense the new state of the robot and people, and the new iteration begins. This cycle is shown in Figure 11.

A popular method to implement the global and local planners is by using a *costmap*. A *costmap* not only has the same representation as the map, however collision-free positions can have non-zero costs. A lower cost cell is more favored to be in to a higher cost cell. After the calculation of all cells, the least-cost path is found that connects the start position to the goal position.

Note that this approach assumes the robot is able to execute any path provided to it. Holonomic robots can move in any direction, however non-holonomic robots has limitations in their movements. For example, two wheel robots can not move sideways. Two common approaches to solve this problem are: to implement trajectory planners that can handle imperfect control or to embed the robot's dynamics into sampling for global and local planning.

3.2 Related Work

In this section, we review the literature on robot navigation in human environments including socially acceptable navigation, learning behaviors from humans and cooperative navigation.

3.2.1 Socially Acceptable Path Planning

Socially acceptable robot navigation is considered in different applications such as free navigation [112], approaching people [102] and evacuating buildings [85]. Some works used the personal space concept in cost-based general path planners [112, 54]. Sisbot [112] models the social spaces as a ellipse-shaped Gaussian, and takes into account the safety, preferences and vision fields of humans for a robot that navigates from a location to another. Kirby [54] presents a path planner that takes into account social conventions such as tending to one side of the hallways. A potential field based trajectory planner for dynamic human environments is presented by Svenstrup [116]. Rios-Martinez [97] presents a RRT-based planner that considers not just safety but also the disturbance of humans. In simulation, if interaction within a group of people is detected, the robot can either not disturb the interaction or join the group. This approach is implemented on a wheelchair robot [125]. Althaus [1] presents a robot that can join a group of people and adjust to the formation reactively. The scenario where a robot encounters a human in a hallway is studied by Pacchierotti [89]. Parameters such as the distance between the human and the robot when the robot begins to deviate from its path and lateral distance that robot should be placed when it is passing the human are found from experiments. Recent work by Lu [68] showed that using gaze cues and social navigation makes robot-human hallway passing more efficient.

3.2.2 Learning Navigation from Human Behavior

Behaving human-like in robot navigation is usually favored in the literature [101]. One way to simulate human navigation behavior is to use social cost maps that capture social conventions [103, 69]. Contrary to the imitation approach, [9] tries to avoid predicted paths, with the goal to minimize the risk of interference. Kuderer [61] presents a tele-operated robot that computes the policy of a desired interactive navigation by learning from observations of pedestrians. Pellegrini [92] trains a dynamic social behavior, that account for social interactions, using pedestrian data.

3.2.3 Human Cooperation in Robot Navigation

Robots can exploit human cooperation in certain scenarios. In populated environments, one way to move with the crowd is to follow individuals that move towards the robot's goal [114, 80].

Some of the recent works in the literature claim that the robot motions should be predictable so the human observers can judge the motive and future behavior of the robot. Observational study in [66] claims that three features can increase the predictability of robot navigation: straight lines, stereotypical motions and usage of additional gestures. In a user study conducted by Gockley [33], humans observers watched two ways of person following. People found direction-following more natural than exact path following. Kruse [59] observes that when paths of two humans are crossed at a right angle, they adapt their velocity rather than the path. This behavior is implemented on a robot, resulting in more predictable motions.

Trautman [120] introduces the 'freezing problem', where traditional path planners fail to produce a feasible solution in crowded human environments. Muller [80] briefly mentions a 'shooing away' behavior, where the robot accelerates towards a human, hoping that he/she will get out of the way. Kruse [60] introduces an optimistic planner, which assumes that people will cooperate with robot movements. Their

approach relies on assigning a non-infinite cost if a robot enters to a human’s personal space, however the plan fails if humans doesn’t move as expected because of the lack a local planner.

3.3 Goal Points for Navigation

As presented in Section 2, our interactive system allows a user to annotate landmarks. After completing the *HomeTour*, the robot can navigate to or towards the labeled entities. A user can enter a navigation destination to the robot in three distinct ways: via labeled waypoints, planar landmarks or objects.

3.3.1 Labeled Waypoints:

If a waypoint is labeled and saved, the robot attaches that label to the explicit coordinates, namely position and orientation. Therefore, if the robot is instructed to navigate to a labeled waypoint, then the goal is readily the pose of the waypoint.

3.3.2 Labeled Planar Landmarks:

If the label is attached to planar landmark, or a set of planar landmarks, we use the following methodology depending on the number of landmarks attached to the corresponding label:

3.3.2.1 Only a single plane has the corresponding label:

We assume that the robot should navigate to the closest edge of the plane, so we select the closest vertex on the landmark’s boundary to the robot’s current position. This point is projected down to the ground plane, as our robot navigates on the floor. We calculate a line between this point and the robot’s current pose, and navigate to a point on this line a meter away from the point, and facing this point. This results in the robot navigating to near the desired landmark, and facing it. This method is suitable for both horizontal planes such as tables, or vertical planes such as doors.

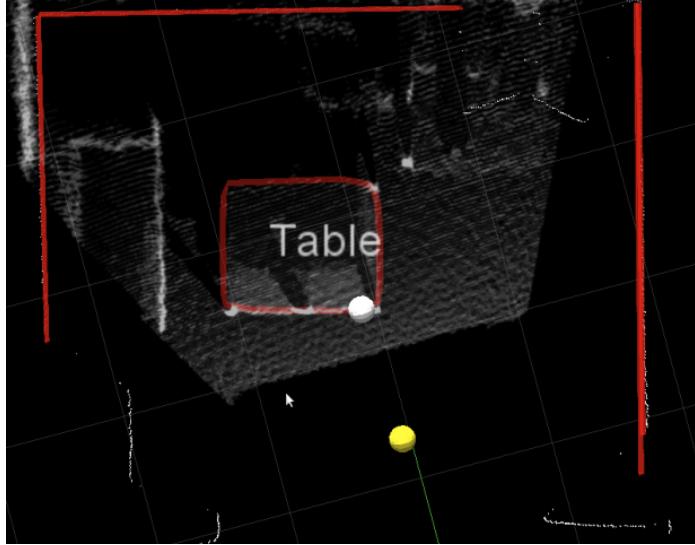


Figure 12: Top down point cloud view of a room. A planar landmark with label *Table* has previously been annotated by a user. The convex hull for the planar landmark is shown in red lines. When asked to navigate to *Table*, the robot calculates a goal pose, which is shown as the yellow point.

An example for calculating a goal for a single labeled planar landmark is shown in Figure 12.

where a the goal point corresponding to the singular label *Table*.

3.3.2.2 Multiple planes are attached to the same label:

We assume that the requested label corresponds to a region of space such as a room or corridor. In this case, we project the points of all planes with this label to the ground plane, and compute the convex hull. For the purposes of navigating to this label, we simply navigate to the centroid of the hull. While navigating to a labeled region is a simple task, this labeled region could also be helpful in the context of more complex tasks, such as specifying a finite region of the map for an object search task. An example for calculating a goal for a two labeled planar landmarks is shown in Figure 13.

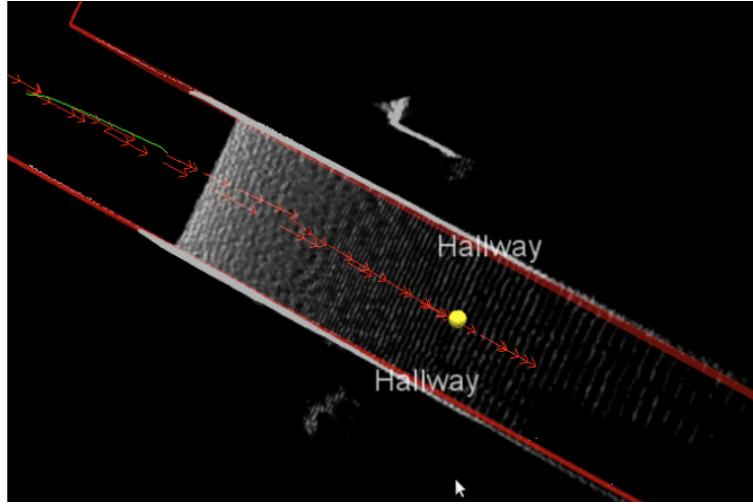


Figure 13: Top down point cloud view of a hallway. The user has previously annotated two planar landmarks with the same label, *Hallway*. When asked to navigate to *Hallway*, the robot calculates a goal pose, which is shown as the yellow point.

3.3.3 Labeled Objects:

As discussed in Section 2.2.3, we first perform planar surface detection before detecting tabletop objects. When the robot is asked to navigate to a labeled object, the planar surface that the object lies on is given as the goal landmark. The robot calculates the goal position as described in the single labeled landmark case in Section 3.3.2.1.

3.4 People Aware Navigation

As extensively reviewed in Section 3.2, people-aware navigation algorithms aim to generate human-friendly paths that consider the safety and comfort of people. A common assumption for point-to-point people aware navigation is that humans are independent agents and that robot's motions have no effect on people's motions. However, humans navigate by constantly anticipating other people's reactions. Similarly, mere presence of a robot in motion is likely to influence how nearby humans would move.

Robots can potentially use this implicit cooperation between moving embodied agents. For example, consider a robot that is outside a room and given a goal pose

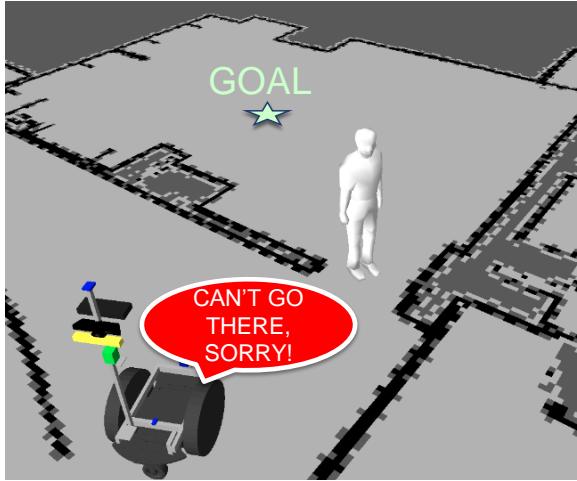


Figure 14: Standard path planners fail to produce a solution to the 'room problem'. Our people-aware planner anticipates that the human can give way to the robot if it approaches towards its goal.

in the room. There is a person standing at the door and blocking the path. Such an example is illustrated in Figure 14. Standard path planners, as well as planners that consider dynamic objects fail to produce a solution to this problem. The role of physical embodiment in human-robot interaction is significant [127], however it is commonly ignored in robot navigation. A people-aware planner should anticipate that the human may give way to the robot if it expresses its intent to go inside the room. Extending this idea, by using anticipation a robot can reduce its time of travel and behave more human-like in general cases.

In this section, we propose a people-aware navigation planner that considers reactions of humans to robot motion. Our planner first finds the least-cost map in the costmap that considers safety and disturbance of people. The costmap definition is discussed in Section 3.4.1. Then the path is refined by simulating people's reaction to robot's motion using Social Forces Model [38]. The path refinement will be discussed in Section 3.4.1.1. In dynamic simulation, robots and humans repulse each other, and additional forces helps to stay away from obstacles and conserve formation in groups. Paths are re-planned when the world state changes or humans does not move

as anticipated. In Section 3.4.2, we discuss our local planner. We then discuss the implementation of the system in Section 3.4.3, demonstrate two example scenarios in simulation in Section 3.4.3.1 and two on the real system in Section 3.4.3.2.

3.4.1 Global Planner

The global planner takes the start and goal positions and a 2D grid map as input and aims to find a set of waypoints that connects the start and goal cells. The output path has the minimum cost with regards to a cost function with 3 parameters: path length, safety and disturbance. We use A* search with Euclidean heuristics on a 8-connected grid map to find the minimum cost path. The configuration space obstacles are found by inflating the map obstacles for as much as the radius of the robot with the assumption that the robot is circular.

Path length cost: Each action a of the robot (moving to one of the 8 adjacent cells) has a non-negative action cost $Cost_a(x_i, y_i, a)$. If the destination cell is occupied by a configuration space obstacle, then the action cost is infinite. Otherwise, it is the distance in meters. The action cost is thus defined as:

$$Cost_a(x_i, y_i, a) = \begin{cases} u & \text{if } a = \text{N, E, S, W} \\ u\sqrt{2} & \text{if } a = \text{NW, NE, SW, SE} \\ \infty & \text{if } Cell(x_{i+1}, y_{i+1}) \text{ in obstacle} \end{cases} \quad (1)$$

where N,NW,.. are the grid cell expansion directions and u is the grid cell size. The resulting path length cost of a path P is then the sum of all action costs:

$$Cost_{path}(P) = \sum_{a \in P} Cost_a(x_i, y_i, a) \quad (2)$$

Safety cost: The notion of safety is the absolute need of any human-robot interaction scenario. This cost is a human centered 2D Gaussian form of cost distribution and aims to keep a distance between the robot and the humans in the environment. While some approaches used un-isotropic cost functions to account for human orientation, we use a isotropic Gaussian for its simplicity. Each cell coordinate around a

human contains a cost inversely proportional to the distance. Since the safety loses its importance when the robot is sufficiently far away from the human, safety cost becomes zero after a threshold distance. If there are multiple people in an environment, the safety cost of a cell takes its value from the closest human.

$$Cost_{safety}(x, y) = \begin{cases} u \max_{h \in H} (\mathcal{N}(\mu_h, \Sigma)) & \text{if } d < d_{max} \\ 0 & \text{if } d \geq d_{max} \end{cases} \quad (3)$$

where d is the distance to the closest human, H is all humans, $\mu_h = (|x - h.x|, |y - h.y|)$ and $\Sigma = 0.5I_2$ is a fixed covariance matrix. The multiplication by the grid cell size compensates for the grid map resolution. Otherwise, for example, if a very fine map was used, safety cost would dominate the path length and disturbance costs, which are independent of the map resolution.

Disturbance cost: This cost is aimed to represent the cases where the robot potentially disturbs the interaction of a group of humans. For example, if two people are facing each other and talking, then the robot should not cross between them. We model this with a disturbance cost that is introduced if a path crosses between two people. We do not detect if there actually is conversation between the people, but estimate the disturbance cost using body poses of agents. This cost increases if body orientations of two people are facing each other and is inversely proportional on the distance between the two humans.

For each step taken in the grid, we check if the line segment from the current position to the projected position intersects a line segment between all pairs of humans. To illustrate, let's assume the robot crosses the line between human A and human B in Figure 15(a).

The disturbance cost is calculated as:

$$\begin{aligned} Cost_{dist}(x, y, a) &= \max(0, f(d).(\vec{AA'}.\vec{AB} + \vec{BB'}.\vec{BA})) \\ f(d) &= \frac{1}{d} - \frac{1}{d_{max}} \end{aligned} \quad (4)$$

where all the vectors are normalized and d_{max} is the maximum distance between the humans that returns a disturbance cost. Figure 15(b) illustrates several examples of disturbance costs with $d_{max} = 3$ meters.

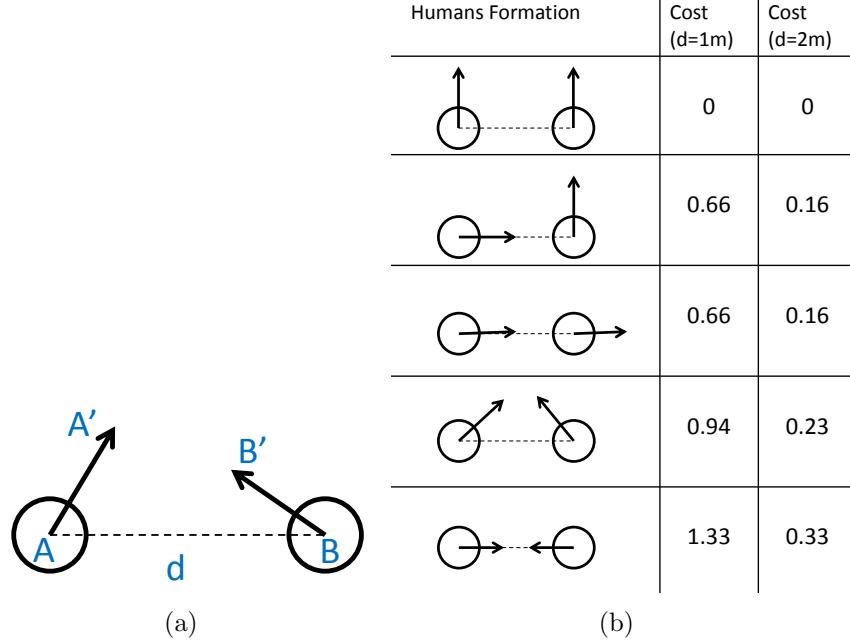


Figure 15: Disturbance costs in different human-human configurations and distances. A path that crosses the dashed lines incurs the disturbance cost calculated on the right side.

Total Cost: The total cost of a path P is computed with a weighted average of path length, safety and disturbance costs. We use A* search to find the least-cost path.

$$Cost_{Total}(P) = Cost_{path} + w_s.Cost_{safety} + w_d.Cost_{dist} \quad (5)$$

3.4.1.1 Path Refinement using Social Forces

In this section, we describe the path refinement process that is applied to the global path. The initial geometric path generated by the global planner is not smooth, therefore robot motion might not be easy to interpret for human observers. The path refinement processes the global plan and simulates the parts of the path where group of humans are closeby. We use Social Forces Model (SFM) [38] to simulate the

motions of both humans and the robot. Interaction between people are modeled as attractive and repulsive forces in SFM, similar to the Potential Field Method [51] for robot navigation. The forces are recomputed iteratively and the resulting simulated paths replaces the corresponding path sections in the global plan.

First, groups of people are found by clustering with respect to their positions. Simple euclidean distance thresholding is used for clustering. In our current implementation, a group region is defined as a rectangle, although other shapes are also possible. The path refinement process receives the global plan and finds out where it enters and exits each group region if it intersects the region. Goal of the dynamic iterative simulation is to find a sub-plan between those two points. Forces apply to all agents, including the robot and humans. We define 4 forces acting on the agents:

- F_{goal} : attraction towards a sub-goal
- F_{social} : repulsion from other agents
- F_{obs} : repulsion from nearest obstacle
- F_{group} : attraction or repulsion towards group members

The forces acting on the robot at the first iteration of forces simulation are illustrated on the robot in Figure 16(a). The force magnitudes with respect to distances between entities are plotted in Figure 16(b).

Starting from the first group region that intersects the static plan, the following procedure is applied within every group region: At every iteration, first the resultant force vector acting on the robot is found. Then the planner takes a step in the direction of the F vector for a fixed step size. Then each of the humans in the group takes a step towards the resultant force that is acting on them. The planner continues the iterations until a solution is found. If a solution is found, the calculated sub-plan replaces the static plan in this group region. Potential fields are known to stuck to

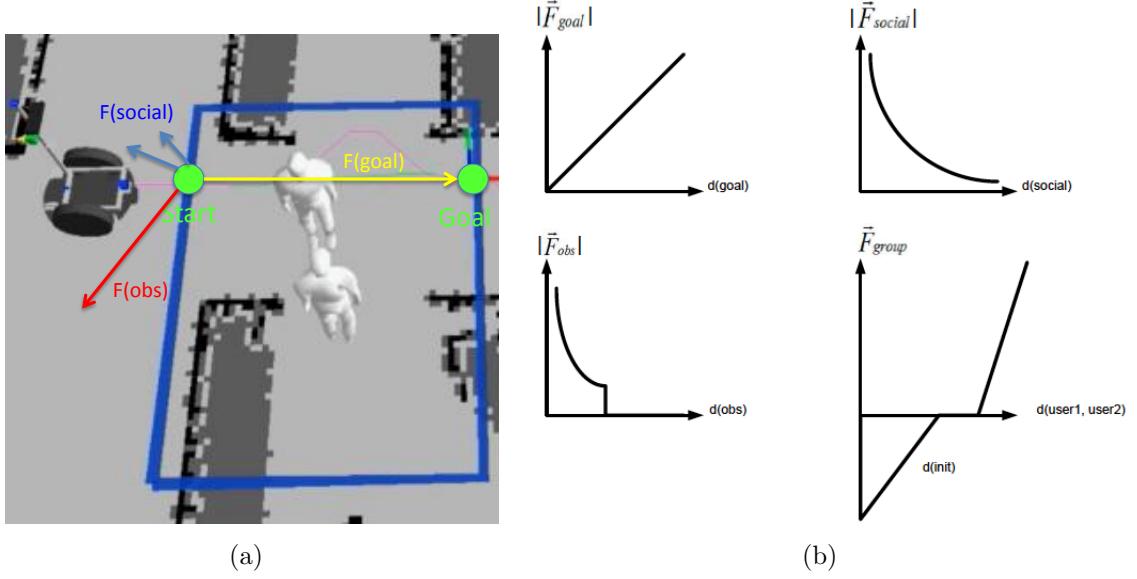


Figure 16: a) Social forces acting on the robot, including $F_{goal}, F_{social}, F_{obs}$, are shown at the first iteration of the dynamic planner. Note that $F_{group} = 0$ as the robot does not belong to a group. The group force (not shown) exists, however, for the humans as they are in the same group region. b) Social forces with respect to the distance towards the corresponding entity.

local minima [57], and the planner might go into infinite loop. We stop the planner after a number of iterations and accept the static plan in the corresponding group region if that happens.

We assume that humans have a cognitive model of the robot, by thinking that the robot has a limited Field of View (FOV). When the robot has gone past a human (out of the FOV), then we make the repulsion force $F_{social} = 0$. We think that humans behave that way: as someone walks past, there are no social constraints resulting from that individual any more.

3.4.2 Local Planner

The local planner is responsible for finding the trajectory that the robot is capable of executing. It accepts a geometric global path as input and computes the linear and angular velocity necessary to follow the dynamic path. We adopt a local planner inspired by Dynamic Window Approach (DWA) by Fox [26]. In the original DWA

approach, only circular trajectories are considered, defined by pairs (v, w) of linear and angular velocities. An objective function, consisting of target heading, clearance from obstacles and velocity of the robot is maximized by sampling admissible velocities.

Our approach also samples admissible velocities, but the optimization criteria we use consists only of the Euclidean distance to a sub-goal point chosen on the path that is ahead of the robot. The velocity pair that resulted in the closest proximity to the sub-goal is chosen and sent to robot controllers. At every control iteration, the sub-goal is chosen as the first point ahead of the robot that is further than a distance threshold. We found that a threshold of 0.25 meters was sufficient to choose the sub-goal. After the local planner calculates the output velocities, they are applied to the robot and the iterative process continues until the the robot reaches the goal. Since the goal is a singular point, it is impossible for the robot to be exactly at the goal. Therefore, a tolerance around the goal point, defined as a circle around the goal is defined.

Given the robot's current pose and an applied velocity, the DWA approach requires to have a motion model for the robot. The motion model projects the what the robot pose would be, if a velocity pair is applied to it for a time period. The robot we used for our implementation is a non-holonomic two wheeled robot. While one can use the general motion equations derived in [26], we used linear approximated motion equations for our robot in Equation 6. Given a robot pose $q^t = (x^t, y^t, \theta^t)$ at time t and an input velocity (v, w) , the projected robot pose at time $t + \Delta t$ is:

$$q^{t+\Delta t} = f_{motion}(q^t, v, w, \Delta t) = \begin{cases} x^t - \frac{v}{w} \sin(\theta^t) + \frac{v}{w} \sin(\theta^t + w\Delta t) \\ y^t + \frac{v}{w} \cos(\theta^t) - \frac{v}{w} \cos(\theta^t + w\Delta t) \\ \theta^t + w\Delta t \end{cases} \quad (6)$$

3.4.3 Results

In this section, we provide qualitative results both in simulation and on the real robot. We used a non-holonomic drive robotic platform, Segway RMP-200, for the real experiments. We used our laser-based torso tracking method presented in Section 4.2.2.

3.4.3.1 Simulation

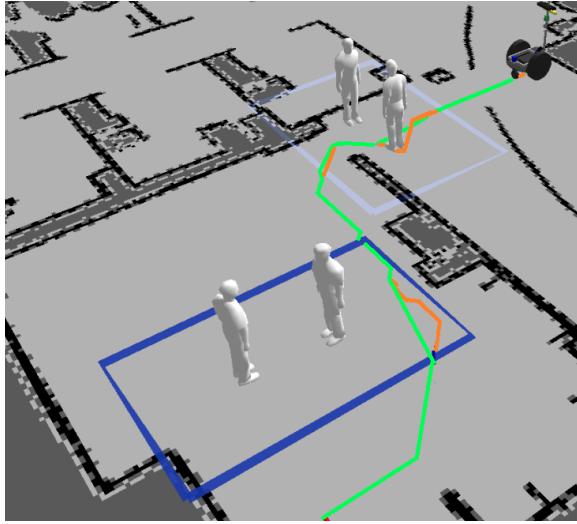
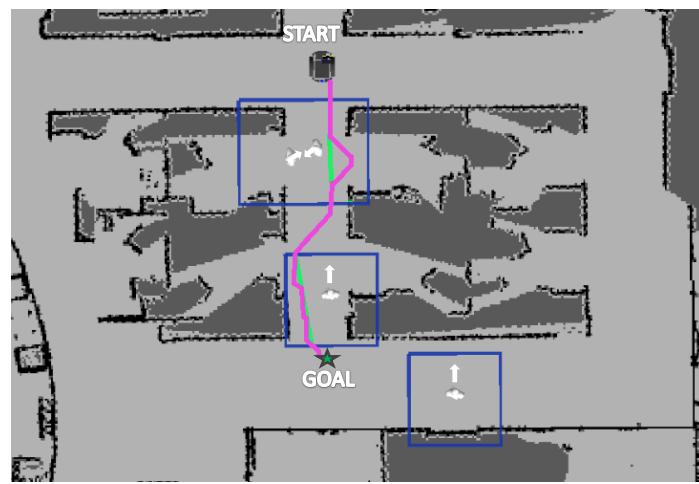
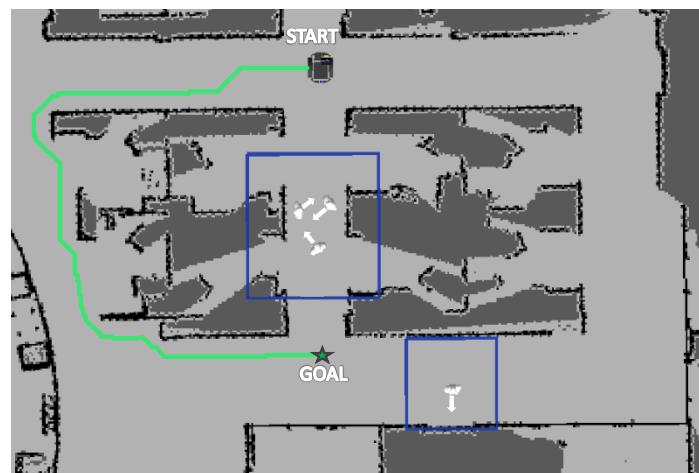


Figure 17: "Room Problem". The robot is outside a room and the goal is inside the room. Traditional planners can not solve the problem because two people are blocking the doorway. Our planner generates a tentative path, with the initial global plan shown in green and the dynamic refinements are shown in orange.

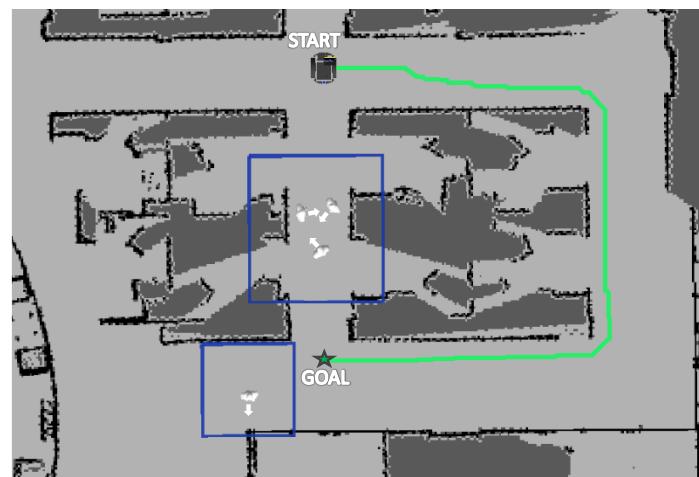
Room Problem: In this scenario, the robot is outside the room and a point inside the room is given as the goal (Figure 17). Traditional planners can not return a solution in this scenario because there is not enough space for the robot to navigate inside. There are two people standing at the doorway and there are two more standing people inside. The static plan and dynamic plans are shown in green and orange, respectively. This path is planned for the current time but makes assumptions about future positions of humans. Note that the dynamic planner modifies only the parts of plan inside group regions (blue rectangles). In the first group region (doorway), the



(a)



(b)



(c)

Figure 18: Paths differ drastically with the poses and grouping of humans. a) The robot takes shortest route, traveling in the vicinity of a group of two and another individual. b) third individual joins the group. Robot takes a longer path that doesn't have humans on path. c) fourth person changes his position, leading the robot to take the longest route.

static plan involves going between the humans. Dynamic simulation suggests that people will get closer to each other if the robot drives towards the side. In the second group region, since two humans are oriented to each other, going between them would add a high disturbance cost, therefore the static plan avoids going between them. Safety costs encourages staying far from the humans, but not too far because a longer path would increase the path length cost. The robot is further led to stay closer to the room boundaries in the dynamic planner due to the repulsive forces from both humans.

Office Environment: Goal of the robot is to navigate to a goal position in an office environment with 4 standing people (Figure 18). In this scenario, we show how the planned path is drastically changing with the poses of humans even though the start and goal position of the robot doesn't change. There are 3 main ways the robot can navigate to its goal: left, center or right corridor.

In the first configuration in Figure 18(a), two people are grouped together as they are looking at each other and likely conversing. The robot decides to take the center corridor, first slightly disturbing the speaking duo, then switches sides in the corridor and reaches its goal. In the figure, the dynamic path (pink line) is overlaid on the static path (green line).

In the second configuration in Figure 18(b), The third person at the center corridor joins the conversation. Now we have 2 group regions (rectangles) in the scene. Since passing through a group of 3 people would introduce a high disturbance cost in addition to the safety cost, the robot decides to take a longer route (left corridor). Since this path does not intersect any group regions, no dynamic simulation is done.

In the third configuration in Figure 18(c), the group of three hasn't moved, but the fourth person has changed its position. In this case, if the left corridor is taken again, an additional safety cost would be incurred. Therefore the robot decides to take the longest route (right corridor). Again, since the robot travels far from

humans, no dynamic simulation is done.

3.4.3.2 Real Robot

We demonstrate our anticipatory navigation planner on the real system in two environments: hallway and kitchen. Each scenario is run twice under different human positions and behaviors in order to show how the planner responds.

Hallway passing: In this scenario (Figure 19), robot’s goal is to navigate to the end of the hallway. In the first run, humans move as the robot anticipates. In the second run, humans do not move as anticipated, and the robot adjusts its path. Each step is described in the caption of the figure. In both cases, the initial plan is to disturb the interaction by going between the two. This is because the safety cost for getting close to one of the humans was more dominant than the disturbance cost.

Narrow corridor: In this scenario (Figure 20), robot’s goal is to drive towards the exit door. There are 3 people nearby the robot. The robot can either take the shorter route that is the direct path, or take a longer path that is to the left of the table. Each important step is described in the caption of the figure. The first run shows that the robot may plan hoping to influence the human. The second run shows that the robot may take a longer route if the disturbance and safety costs are going to be large.

3.5 Speed Limits for Safe Navigation

In this chapter so far, we studied navigation planning that aims to generate human-friendly paths and safety of people must be at the utmost importance. In order to ensure the safety of people, we introduced **Safety Cost** in Section 3.4.1 as well as repulsion forces in Section 3.4.1.1. However, a seemingly safe path generated by our approach can still lead to an accident because it is may not be possible for the robot to track every human in the environment all the time. For example, when the robot is turning a corner, there is a risk that the robot suddenly encounters a person.

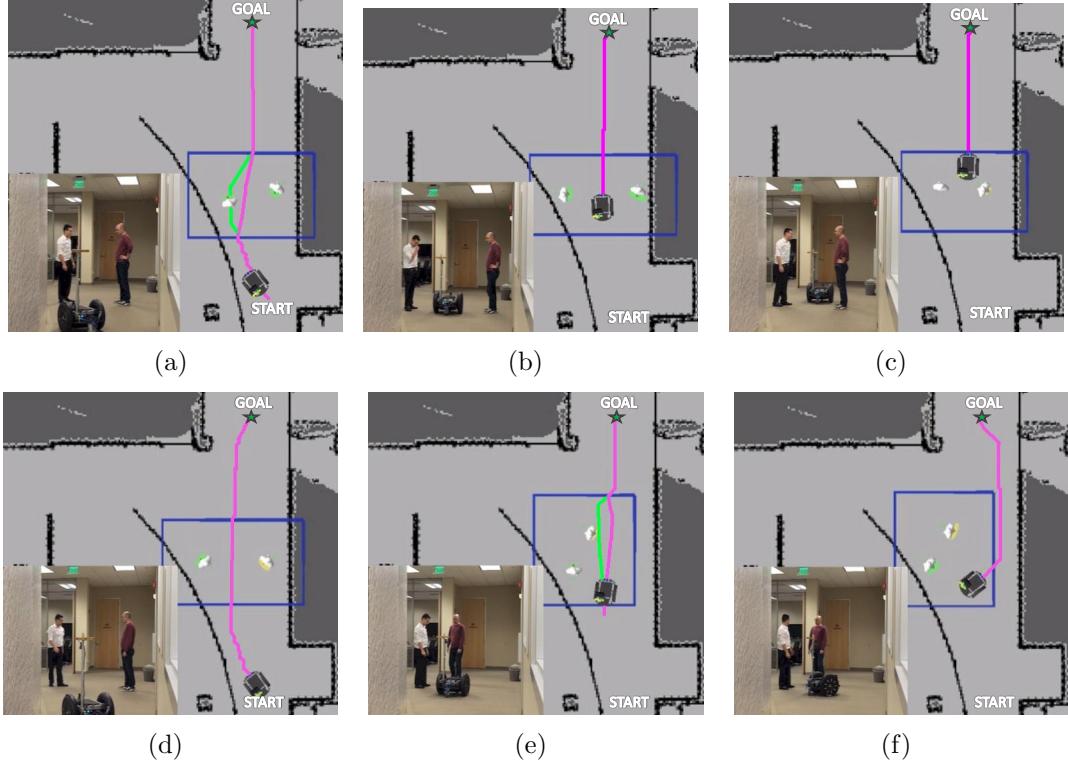


Figure 19: The Hallway scenario. 2 runs are shown in first and second rows. The static plan (green line) and dynamic plan refinement (pink line) are shown. First run: a) Navigation starts. The dynamic planner anticipates that people will give way to the robot when it starts to move towards them. b) Humans notice the robot, and give way by increasing the separation between them. c) The robot continues towards its goal and humans regroup. Second run: d) both the static and dynamic plan involves going in between humans again e) human on the right gets closer to the other person. Since a human made significant movement, dynamic planner re-plans. Plan no longer involves going in between. f) static planner periodic re-plan triggers, leading to robot to stick to the wall to the right.

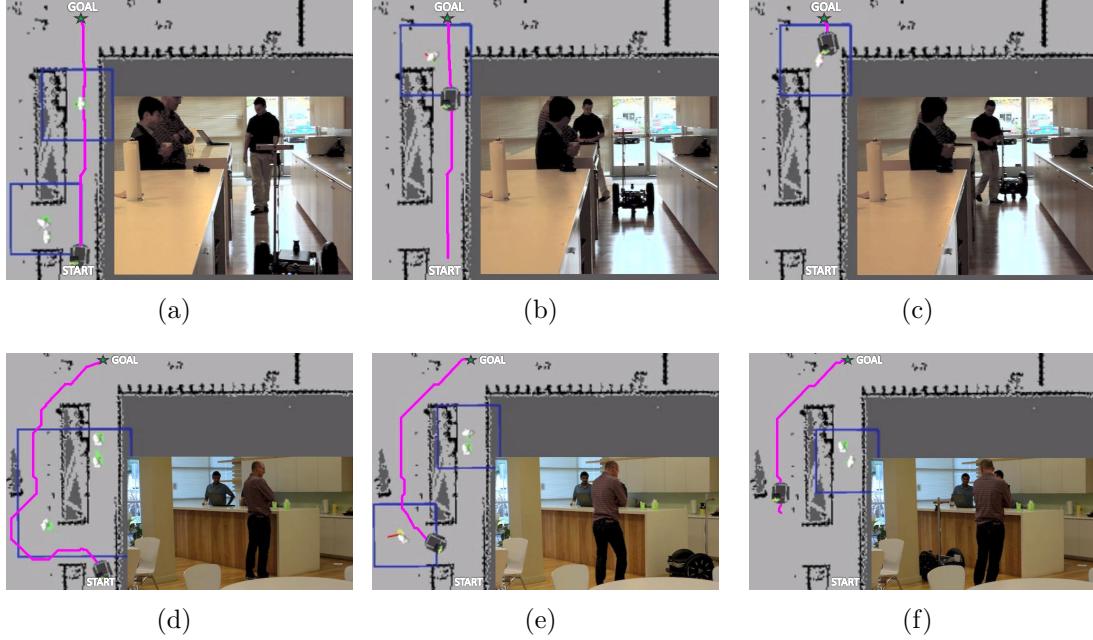


Figure 20: The Kitchen scenario. In the first run, there are two people blocking the path to the left and one person at the narrow corridor. a) robot decides to take the shorter route, because it would disturb one person instead of two. There is not enough space to pass, and dynamic planner assumes the person would get out of the bottleneck to give way. b) human behaves as robot anticipated and gets out of the narrow passage. robot slows down because it enters the human region. c) person gets back to his original position, robot reaches the goal. In the second run: d) there are two people at the narrow corridor and one person on the left. The robot decides to take the longer route and pass the third person from left. The safety cost from the two others would be too high if the robot took the direct route. e) the person steps back as he recognizes the robot. since the person has moved, the dynamic planner re-plans and decides to pass from right. f) after the robot passes the person, it proceeds to its goal.

Since the robot have finite deceleration, a collision may be unavoidable. The speeds that mobile robots should exhibit is also dependent on the context. For instance, a robot should move slowly and carefully in a hospital room. On the other hand, the robot may navigate faster in a long office corridor. The speed limits should better be provided by experts or the building owners, who govern how fast the robots should navigate in a particular environment.

In this section, we introduce speed maps that sets the speed limits for mobile robots in an environment. We claim that usage of such speed maps make the robots safer and potentially more efficient. The speed map designed for the second floor of the College of Computing building at Georgia Tech is shown in Figure 21.

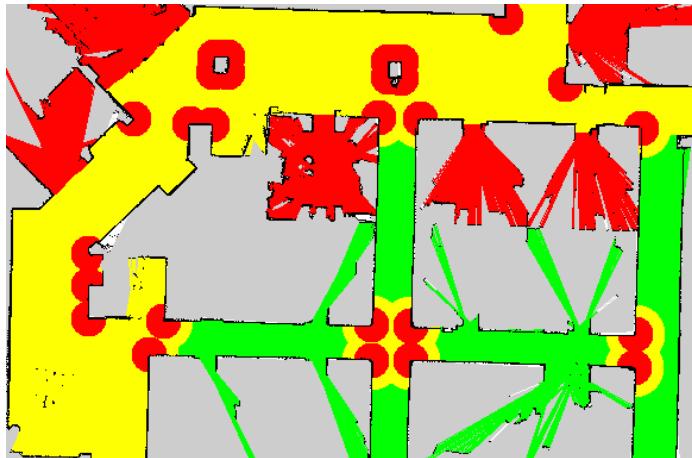


Figure 21: Designed speed limits for. The robot has to be relatively slow in red zones, can have moderate speed in yellow zones and is allowed to move relatively faster in green zones.

The free spaces in this map are divided into three zones:

1. Green zone: The robot is allowed to travel at relatively faster speeds.
2. Yellow zone: Human interaction is possible, top speed is less than the maximum allowed
3. Red zone: Human encounter is probable, top speed is severely limited.

This speed map is designed by hand using the following rules: Spaces corresponding to rooms and cubicles are covered as Red Zones. Blind corners are covered with a Red Zone close to corner and Yellow zone enclosing the Red Zone. Corridors are covered as Green Zones and the rest is covered as Yellow Zones.

The speed map depicted in Figure 21 is designed by hand, however the speed map generation can be automatized using automatic room categorization and additional processing. Room segmentation has been proposed in an interactive fashion by [22], as well as automatically, especially for creating topological maps [79].

3.5.1 Results

In this section, we evaluate the effect of applying speed limits in a navigation scenario. The robot is in an office environment and it has to turn a corner to reach its goal. There is a person standing right around the corner and the person is not visible to the robot until it gets fairly close to the person. We had two conditions of speed limits:

- Condition 1: Fixed maximum linear speed of $v_{max} = 1.0m/s$
- Condition 2: Variable speed limits are used according to the speed map: $v_{max}(green) = 1.5m/s$, $v_{max}(yellow) = 0.5m/s$, $v_{max}(red) = 0.15m/s$

In both conditions, standard ROS Navigation is used, which found the lowest cost path on a costmap consisting of costs from nearby obstacles. The robot detected the person using our multimodal person detection and tracking method described in Chapter 4. We analyze the distance/velocity of the robot as it encountered the human and measure the time to reach to the goal as evaluation metrics.

The section of the speed map including the corner is shown in Figure 22(a). The trajectory with a fixed maximum speed is shown in Figure 22(b). Note that the robot got dangerously close to the human while it turned the corner in Condition 1 and it was traveling with about $0.3m/s$ when it detected the person. The trajectory using

the proposed speed limits is shown in Figure 22(c). With this approach, the robot slowed down as it approached the corner, therefore allowing earlier detection of the human. The robot gave more space to the human in this case, the speed of the robot when it encountered the robot was $0.15m/s$, half the speed of fixed max velocity case. Considering the speed of the robot and distance to the human at the time of encounter, we can say that the robot was safer.

The second metric we measured was the time to reach the goal. The robot was more efficient with the proposed approach, reaching the goal in $28s$ compared to $29.1s$.

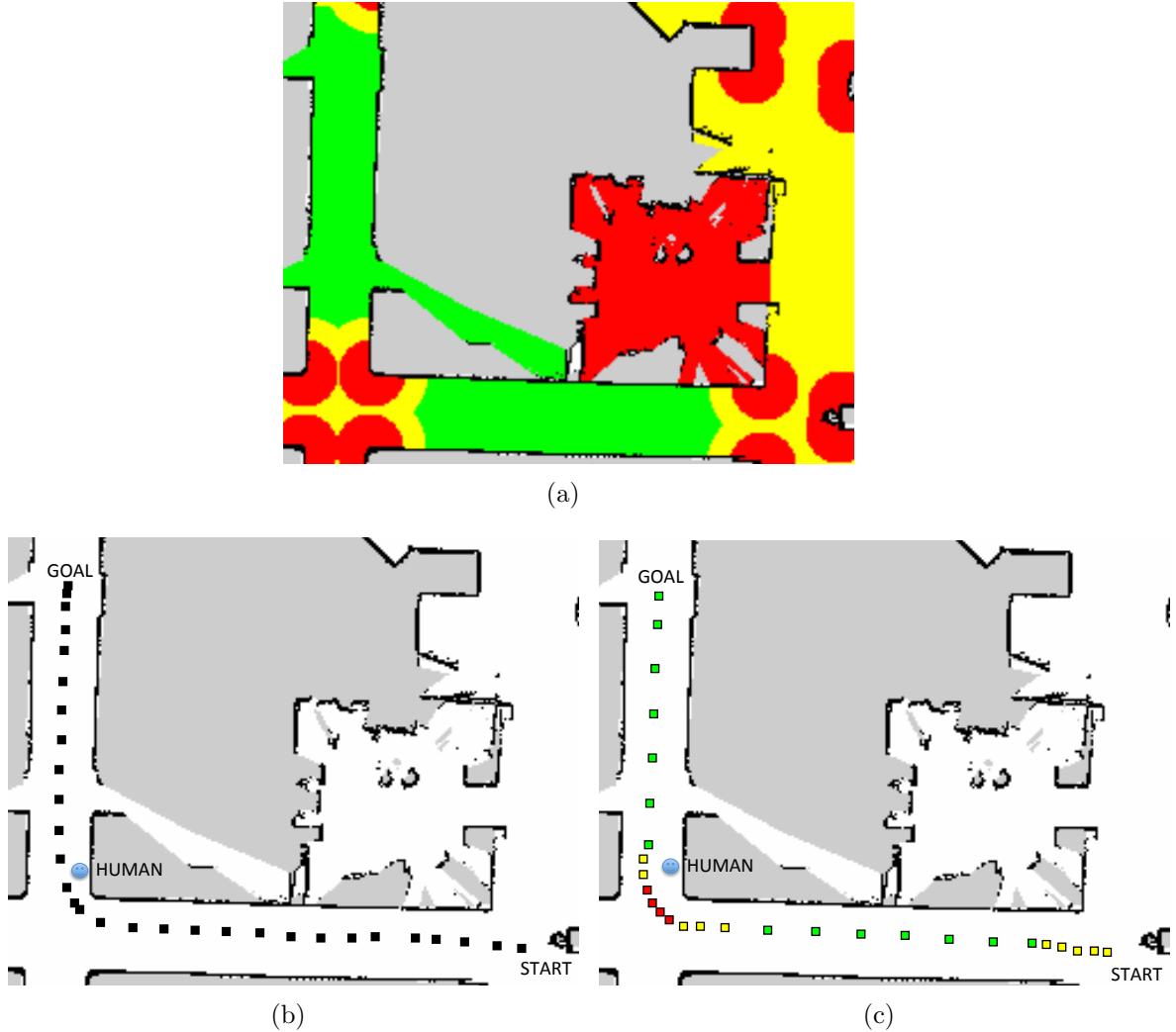


Figure 22: Comparison study of using a maximum top speed versus using location-dependent speed limits. Robot is given a fixed goal location. Right around the corner, there is a bystander human, who is not visible to the robot until the robot makes the turn. Points annotate robot position measured at fixed time intervals. a) Speed map of a corridor intersection at the second floor of College of Computing at Georgia Tech. b) Robot's top speed is fixed at 1.0m/s . Note that the distance between robot positions are mostly constant. The robot gets very close to the bystander because it is moving relatively fast when it turned the corner. c) The robot is allowed to move with 1.5m/s in green, 0.5m/s in yellow and 0.15m/s in red zones. Colors of the sampled points on the path show the associated speed zone. It can be seen by looking at robot's positions that this approach was more gracious turning the corner and respecting human's personal space.

CHAPTER IV

MULTIMODAL PERSON DETECTION AND TRACKING

The ability to robustly track a person is an important prerequisite for human-robot interaction. To realize any task that involves humans, the challenge is the detection and tracking of humans in the vicinity of the robot considering the robot’s movements, sensing capabilities and occlusions. The scope of how much information is needed from the human perception module depends on the objective of the application. First, the robot should determine if there are people nearby. If the robot senses people around, the robot should find out *where* they are. Representing people as points (x,y) in maps is common practice for navigation planning. If the task requires the robot to face a person, then the orientation θ needs be detected. The robot further can determine *who* the detected person is. Identification of humans is necessary for enabling non-generic service. Finally, the robot should interpret *what* the person is doing by analyzing the motion features and through gesture analysis. Tracking body parts of humans over time give significant information about human activity.

We focus on tracking people who are either walking or standing, as these are the two most common human poses around a mobile robot. Many full-body or body part detectors have been developed in the literature, reviewed in Section 4.1. Full-body detectors are not suitable for mobile robot navigation applications because of their inability of capturing the entire body with on-board sensors when people are close to the robot. We aim to robustly track a person 360° around the robot. However, most sensors have a limited field of view and using only a single detector can lead to a system with a single point of failure. Therefore, we think a multimodal detection system is better suited for on-board people tracking for our use cases.

Laser scanners are the natural sensor of choice as state-of-the-art mobile robots are already equipped with an ankle-height laser scanner that is mainly used for navigation. The laser scanners we used on our robot are Hokuyo UTM 30-LX, which has 270° Field of View (FOV), 0.25° angular resolution, 40Hz refresh rate and 30m maximum range. We are only interested in detections in close range (less than 5m). In that range interval, and the accuracy of each laser reading is ±3cm, which is sufficient for our use cases. The relatively higher accuracy and resolution are the two advantages of laser scanners over cameras and RGB-D cameras. Cameras, on the other hand, have the advantage of providing richer information, which can be used to extract body parts. We use a combination of detectors using either a laser scanner and RGB-D camera for robustness and better coverage, described in Section 4.2. Representing people as points in the map is sufficient for mobile robot navigation and each detector produces a point as a person hypothesis. We use a real-time probabilistic tracking framework that relies on the fusion of the multiple person detections, described in Section 4.3. For certain applications, identifying specific users allows the robot to go beyond generic capabilities. We present our face recognition method in Section 4.4.

4.1 Related Work

Person detection was first addressed by the computer vision community as an object detection problem. Early research on person detection using vision is surveyed by Moeslund [77]. Face detection is a common method for detecting people, with the work of Viola and Jones [126] being the most popular one. See Zhang [133] for a survey on contemporary approaches on vision based face detection. Another popular topic has been pedestrian detection in crowded scenes Leibe [64] and Tuzel [124].

In 2000's, laser scanners became the de-facto sensor for localization and mapping. Laser scanners are usually placed slightly above floor for obstacle avoidance, therefore leg detection is common practice. Early works by Montemerlo [78] and Schulz [105]

focused on tracking multiple legs using particle filters. Legs are typically distinguished in laser scans using geometric features such as arcs [129] and boosting can be used to train a classifier on a multitude of features [3]. Topp [119] demonstrates that leg tracking in cluttered environments is prone to false positives. For more robust tracking, some efforts fused information from multiple lasers such as Carballo’s work [15], which uses a second laser scanner at torso level. Glas [31] uses a network of laser sensors at torso height in hall-type environments to track the position and body orientation of multiple people. Several works used different modalities of sensors to further improve the robustness. Kleinehagenbrock [56] and Bellotto [7] combine leg detection and face tracking in a multi-modal tracking framework. Other examples include combining sound localization and vision [10] and combining RFID tracking and vision [30].

Laser-based person methods pertain tracking of humans in 2D, projected to floor plane. Tracking of the body parts has long been a topic of interest in vision [6, 108]. With the recent introduction of 3D sensors such as the Velodyne, Swissranger and Kinect, more robust tracking became possible. Spinello [113] trains geometrical features at different height levels in the 3D point cloud for pedestrian detection. Ganapathi [27] estimates body part locations with a probabilistic model. One of the well-known skeleton tracking algorithms is the Microsoft Kinect SDK by Shotton [107], which trains decision forests using simple depth features in a vast database. This software is not suitable to work on a mobile robot as it is designed to work on a stationary sensor. In the robotics community, there are efforts to develop skeleton trackers that work on mobile robots and in unstructured scenes [14].

Face recognition is a widely used application as surveyed by Phillips [93]. One of the pioneers in face recognition uses a set of patch masks for features that doesn’t necessarily correspond to eyes, ears or noses [123]. [134] combines PCA (Principal

Component Analysis) and LDA (Linear Discriminant Analysis) to improve the generalization capability when only a few samples are available.

There has been some work to identify humans using 3D data, such as the head-to-shoulder signature [55] and body motion characteristics [81]. Biometric person identification techniques, such speaker recognition [53], 3D ear shape [130] and multi-modal cues [28] have potential to be more accurate than face recognition. However, these approaches are better suited to work in controlled environments.

4.2 Person Detection

In this section, we present our person detectors, namely leg detection (Section 4.2.1) and torso detection (Section 4.2.2). We also use an implementation of an upper body detector by Mitzel [75], which uses a template and the depth information of a RGB-D camera to identify upper bodies (shoulders and head), designed to work for close range human detection using head mounted cameras.

4.2.1 Leg Detection

A front-facing laser scanner at ankle height is used for leg detection. The output of a laser scanner at each iteration is an array of range measurements, represented in the polar coordinate system. We first convert the range data to Cartesian coordinate system:

$$x_i = \sum_{\phi=\phi_{start}}^{\phi_{end}} r_i \cos(\phi)$$

$$y_i = \sum_{\phi=\phi_{start}}^{\phi_{end}} r_i \sin(\phi)$$

Then we apply segmentation, Segmentation produces clusters of consecutive scan points, which due to their proximity, have a high likelihood of belonging the same object. Two adjacent distance measurements are considered to be in the same segment if the Euclidean distance between them is below a threshold value. Starting from the

start of the range array, a new segment is started if $|r_i - r_{i+1}| < d_{cluster}$. Although some approaches use a variable segmentation threshold that is a function of the range, we use a fixed clustering threshold $d_{cluster} = 0.1m$. The segmentation process results in a set of segments \mathbf{S} . A set of geometric features are extracted from the laser segment.

In a laser scan, legs can appear in different patterns [119]. We look only single leg and person-wide blob patterns as these two cover all the ways legs can be seen in a laser scan. Depending on the application, we accept either only the single leg pattern or both of the patterns. This is explained in more detail in Section 4.3.

There are a number of geometric features that can be extracted from a laser segment, as delineated by Arras [3]. We use three geometric features that is used to detect a leg: segment width, circularity, and Inscribed Angle Variance (IAV):

1. Segment Width: Measures the Euclidean distance between the first and last point of a segment S_i
2. Segment Circularity: This measure is a simple measure to assess if the segment shape resembles a circle. The circularity criterion we used is the ratio of the perpendicular distance from the middle point to the line segment that connects start and end points, to the segment width. For example, in a perfect half circle in Figure 23, the circularity criterion is $|\overline{P_0P_n}|/d_{mid} = 0.5$. In case of a laser scan, as can be seen in Figure 24, we again consider the ratio of d_{mid} to segment width. For this calculation we only consider the middle point as it provides a simple heuristic on circularity.
3. Inscribed Angle Variance (IAV): This feature is originally proposed by Xavier [129], in order to detect circles. We adopt IAV in order to detect legs, which are not necessarily circle-shaped, especially for the person-wide blob pattern. As an example, inscribed angles on a circle is shown in Figure 25. As a geometric

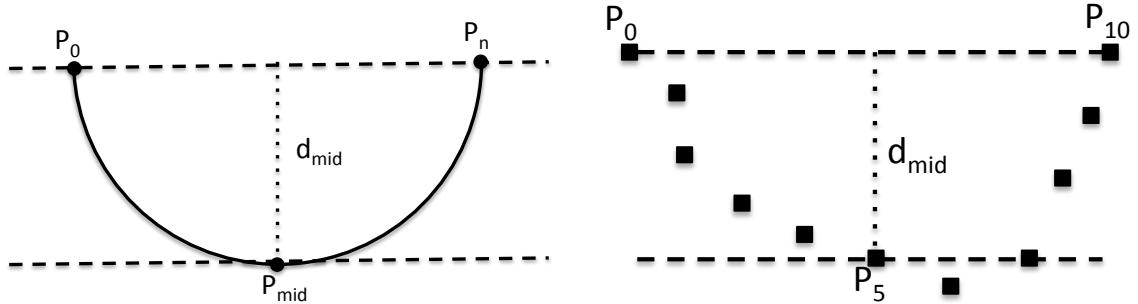


Figure 23: Circularity criterion in a perfect circle is: $|P_0P_n|d_{mid} = 0.5$

Figure 24: Circularity criterion in a this laser segment is: $|P_0P_{10}|/d_{mid}$

property of the circle, $\angle P_0P_1P_4$ and $\angle P_0P_2P_4$ are equal angles. IAV for a given set of points is the average of all inscribed angles:

$$IAV_S = \sum_{P=P_1}^{P_{n-1}} \angle P_0PP_n$$

For a perfect circle, $IAV_S = 90^\circ$. For shapes that are not perfect circles but are similar to circles, IAV feature should be consistent. Laser segments from a leg usually resemble a circle, therefore we use IAV as one of the features for leg detection.

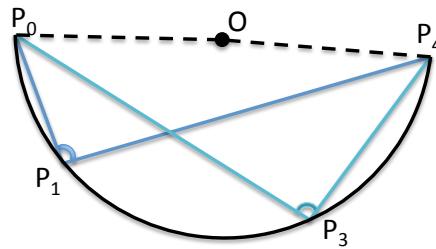


Figure 25: Inscribed angles of an arc are shown in the figure. Inscribed Angle Variance (IAV) is calculated by taking the average of all inscribed angles on a laser segment.

In order to be able to use these values, we first found the nominal feature values for an average human leg. We captured the laser scan data while the robot followed a person through an office environment. The following method used for this experiment will be discussed in detail in Section 5.2. For the training set, two people's legs were

Segment type	Width(<i>m</i>)		Circularity		IAV(<i>radians</i>)	
	μ	σ	μ	σ	μ	σ
Single Leg	0.13	0.03	0.25	0.15	2.23	0.4
Personwide blob	0.33	0.07	0.14	0.09	2.61	0.16
Other	0.22	0.12	0.1	0.11	2.71	0.38

Table 3: Table shows average and standard deviations of geometric leg features calculated in our dataset.

recorded with different clothing (shorts, baggy pants and trousers) to account for variance in the leg parameters. About 17×10^3 Single Leg patterns and 0.6×10^3 person-wide blobs were manually labeled in the data. In addition, 120×10^3 segments were labeled as 'other' or 'not a leg'. The average and variance of the aforementioned geometric features for single leg, personwide blob, as well as other segments are given in Table 3.

For every segment S_i in a test laser scan, we first extract the geometric features f_1^i, f_2^i, f_3^i . We then calculate the weighted Mahalanobis distance to the average leg parameters for each leg pattern:

$$D_{mah}^i = \sum_{j=1}^{n_{features}} w_j \frac{(f_j^i - \mu_j)^2}{\sigma_j^2} \quad (7)$$

where w_j are the weights for each feature, μ_j and σ_j are pulled from Table 3. The resulting Mahalanobis distance is then compared with a detection threshold. If $D_{mah}^i < Threshold_{leg}$, the segment S_i is considered a detection. $Threshold_{leg}$ defines how many standard deviations away from the average features are allowed. In our implementation, we empirically set the feature weights as: $\mathbf{W}_{leg} = (0.35, 0.26, 0.39)$, in the feature order given in Table 3. For normal operation, we set $Threshold_{leg} = 1.5$, which accounts for about %95 of the detections. If only one person is being tracked, we use a higher threshold. The reason behind will be explained in Section 4.3.

4.2.1.1 Associating Leg Segments

After single leg patterns are detected, we try match the leg segments. We extend our leg detection approach to determine which leg segments are connected. Note that this method applies if there is a RGB-D camera pointing to the lower body of the human. For each leg segment pair, if both of them are within the FOV of the RGB-D sensor, we use our algorithm to determine whether there is a connectivity between two candidate leg segments. If a connectivity is found, then the leg segments pair is qualified to be a leg segment pair representing a person. See Figure 26 as an example result. Figure 27 shows the flow chart of the association algorithm.



Figure 26: Two person detections are seen in this figure. Our leg segment association algorithm propagates pixels vertically from candidate leg segments and connects leg pairs.

First, the centroids each of the two candidate leg segments are found. These points are projected onto the depth image acquired from the RGB-D camera. At each iteration, each leg segment, our algorithm first propagates horizontally to both directions in the depth image, then the center pixel is located and it propagates 1 pixel vertically ($+z$ direction). If there are no connectivity after a number of iterations, then we conclude that the candidate leg pair does not represent a person. If there

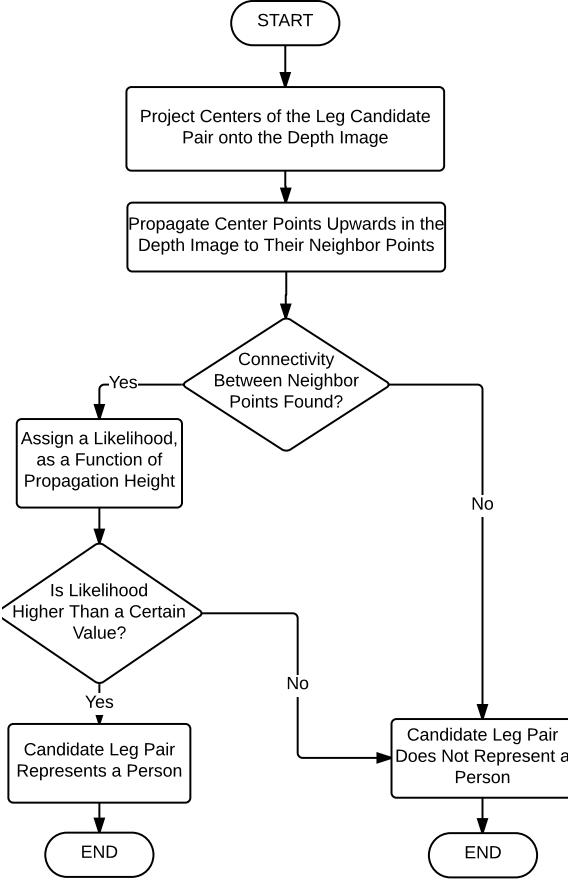


Figure 27: Flow chart for determining if two leg segment candidates belong to a person.

is a connectivity at some point, we then assign a likelihood score to the pair as a function of the vertical propagation height. If this score is higher than a threshold, then the algorithm concludes that the leg candidate segments represent a person. The propagation scoring eliminates most of the false positives due to sensor noise and non-human shapes.

4.2.2 Torso Detection

In this section, we describe our torso detection approach. For this detector, we used another Hokuyo UTM 30-LX laser scanner, placed at torso height ($1.27m$). Our approach relies on fitting an ellipse to laser segments and determining the detection result by interpreting the axis lengths (Figure 28). Our torso detector allows us to

detect the orientation of the person unlike the laser-based leg detectors, therefore this detector is also suitable for applications that relies on extracting the orientation of the person from a single laser scan.

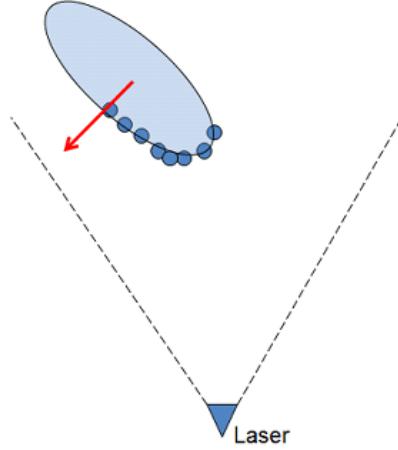


Figure 28: Our torso detector fits an ellipse to the human torso and estimate its position and orientation.

The first step to detect torsos in a laser scan is to segment the laser scan. We use the same segmentation technique used for leg detection, explained in Section 4.2.1. We then fit an ellipse to each laser segment. We use a numerical ellipse fitting method that solves the problem with a generalized eigensystem, introduced Fitzgibbon [25]. This ellipse fitting method is robust, efficient and ellipse-specific, so that even very noisy sensor data will always return an ellipse. Compared to iterative methods, it is computationally very efficient, therefore the speed of the calculations is limited to the laser scan refresh rate.

The ellipse fitting algorithm provides us with the centroid and orientation of the ellipse as well as the minor and major axis lengths. To disambiguate the front/back of a person, we assume that people are facing the sensor when they are first detected. While this is a significant limitation our current system, one can potentially utilize face detection as will be described in Section 4.4 to estimate if the person is facing towards the robot or not.

To detect a torso in a laser segment, we use the minor and major axis lengths,

Torso Features	μ	σ
Width(m)	0.44	0.12
Circularity	0.32	0.18
IAV(radians)	2.57	0.38
Major axis length(m)	0.39	0.08
Minor axis length(m)	0.17	0.06

Table 4: Table shows average and standard deviations of geometric features for a human torso in laser scans.

as well as the three geometric features introduced in Section 4.2.1. We collected 450 laser scans in total while a person stood in front of the sensor and made a one full turn around himself. We calculated the mean and standard deviation of the all five features, which is given in Table 5. For a given laser segment, we find the weighted Mahalanobis distance in Equation 7 to the averaged parameters. If $D_{mah}^i_{torso} < Threshold_{torso}$, the segment is considered a detection. The feature weight constants we used was $\mathbf{W}_{torso} = (0.19, 0.09, 0.35, 0.24, 0.13)$, in respective order given in Table 4. These values were empirically determined, although one can do more sophisticated analysis for optimal weights.

Figure 29 shows how the torso detection rate changes for a given Mahalanobis Distance Threshold in our dataset. What is not displayed in the plot is that higher torso detection rate also means higher rates of false positives. For normal operation, we set $Threshold_{torso} = 1.25$, which accounts for about %90 detection rate. If the tracker is dedicated to track only a single person, then we use a higher threshold: $Threshold_{torso} = 2.5$. The reasoning behind this threshold selection will be discussed in Section 4.3.

4.2.2.1 Evaluation of Torso Detection

In order to evaluate the accuracy of the position and orientation estimations of our torso detection method, we collected torso data from 23 people. Subjects were instructed to stand on 4 targets at different distances with 8 different orientations on

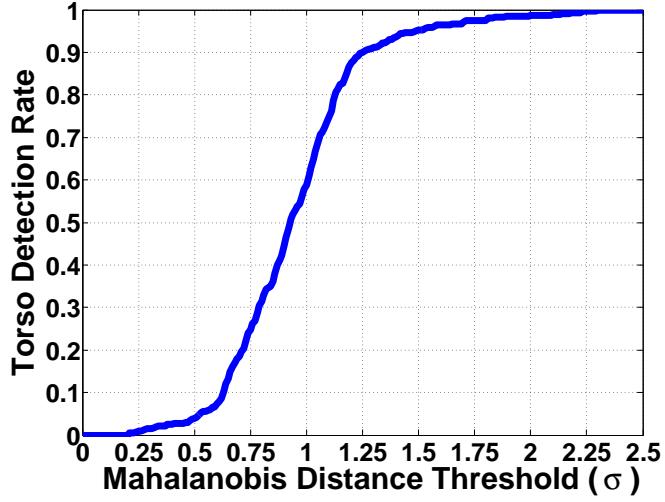


Figure 29: Torso detection rate vs weighed Mahalanobis Distance Threshold in our dataset

each target. Experimental setup from the sensor’s view is shown in Figure 30. For each pose at every target, we logged the position and orientation estimation of the torso detector and compared it with ground truth, which is fixed.

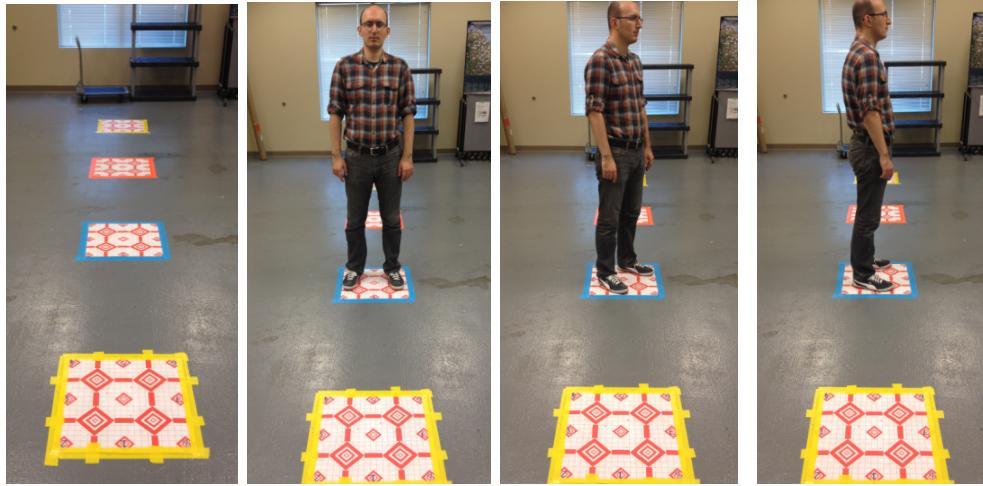


Figure 30: Experimental setup for the evaluation study of the Torso Detector.

Table 5 shows the angular error at every target distance and human orientation with respect to the laser scanner.

The average positional error was about 5cm regardless of the distance and the

Distance To Laser	N	NE	E	SE	S	SW	W	NW	ALL
1.0m	4°	12°	22°	13°	5°	7°	26°	17°	13°
2.5m	5°	16°	19°	10°	3°	6°	14°	17°	11°
4.0m	4°	10°	30°	16°	7°	11°	21°	17°	15°
5.5m	5°	11°	41°	18°	10°	6°	38°	23°	19°
ALL	4°	12°	27°	14°	6°	7°	24°	18°	14.5°

Table 5: Average orientation error of the torso detector with respect to distance from sensor and body pose in a study with 23 people

orientation of the human. The average orientation error throughout all the experiments was 14.5°. Error in orientation, however, varied greatly by pose of the person with respect to the laser scanner. Average error in orientation differed slightly with respect to the distance from the sensor and was the least with 11° when the humans were 2.5m away from the sensor. We attribute to the fact that when humans closer than 2.5m to the laser scanner, it captures more of the arms, which makes the fitted ellipse slightly worse. The orientation of the human with respect to the sensor had a significant effect on orientation error. Least error was achieved when people faced the sensor (4°) or the opposite way (6°). On the other hand, average orientation error was 24° – 27° when humans are perpendicular to the sensor, because a large portion of the torso is not visible to the laser scanner in that configuration.

4.3 Person State Estimation

The position and velocity of the person can not be determined by direct observation due to measurement noise and false detections. Therefore there is a need for a filtering algorithm in order to estimate the state of a person. Using a state predictor for human movement has two advantages. First, the predicted trajectories are smoother than raw detections. Smooth tracking helps the robot maintain consistent trajectories for high-level applications such as Person Following (Section 5). Second, it provides a posterior estimation that can be used for data association when there is a lack of matching detections. This allows the tracker to handle temporary occlusions. We

use a discrete Kalman Filter [46] to predict the position of a person. There are other types of filtering techniques available in the literature, such as Particle Filters [50]. Since the results of the person state estimator is used by time-critical higher level applications, the tracker should come up with an estimate in real time. Therefore the choice of using Kalman Filters was motivated by its computational efficiency. Efficient person state estimation also increases the safety of the robot, as the robot can react faster if there are people in close proximity.

According to Hicheur [39], humans tend to maintain a constant speed when they are walking straight and reduce speed while turning. We used constant velocity model which assumes people will maintain their speed. Even though this assumption is not always true, it provides a simple model without sacrificing too much from tracking performance.

The Kalman filter estimates a process as a predictor-corrector cycle using feedback control. The process has two cycling states: time update and measurement update as shown in Figure. Time update projects the state forward by using the current state and error covariance. Measurement update is responsible for the feedback and corrects the previous estimate.

The Kalman Filter is governed by two linear stochastic difference equations:

$$s_k = As_{k-1} + Bu_{k-1} + w \quad (8)$$

$$z_k = Hs_k + v \quad (9)$$

Where s_k represents the process state at time step k , A is the state propagation matrix, B relates the optional control input u , z_k is a measurement, H is the measurement observation matrix. w and v represent the process and measurement noises, respectively, drawn from normal probability distributions with zero mean $N(0, Q)$ and $N(0, R)$.

We define the state of a person s_k at time step k as:

$$s_k = \begin{bmatrix} x_k \\ y_k \\ \dot{x}_k \\ \dot{y}_k \end{bmatrix} \quad (10)$$

where (x_k, y_k) is the position and (\dot{x}_k, \dot{y}_k) is the velocity of the person in Cartesian Coordinates. With the constant velocity model, the time update equations are:

$$x_k = x_{k-1} + \dot{x}_{k-1}\Delta t_k + w \quad (11)$$

$$y_k = y_{k-1} + \dot{y}_{k-1}\Delta t_k + w \quad (12)$$

$$\dot{x}_k = \dot{x}_{k-1} \quad (13)$$

$$\dot{y}_k = \dot{y}_{k-1} \quad (14)$$

resulting in the following Kalman Filter matrices:

$$A = \begin{bmatrix} 1 & 0 & \Delta t_k & 0 \\ 0 & 1 & 0 & \Delta t_k \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad B = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad H = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \quad (15)$$

where Δt_k is the time difference from the previous detection. A track is lost if there are no detections for a fixed amount of time. At every time update of a filter, if Δt_k is larger than a fixed threshold, the track is killed.

The reason B vector is zero is that we track people in the world frame and robot motion is already accounted for with robot localization. For this reason, we assume there are no control inputs to our system. The noise matrices we used are:

$$Q = qI_4 \quad R = rI_2 \quad (16)$$

where we used $q = 0.02$ and $r = 1.0$ in practice.

Our approach is multimodal in the sense that asynchronous measurements are accepted from different sources as long as they provide a positional estimate in the respective sensor frames. Using the latest localization information, this position is converted to the world frame and then fed as a measurement to the active filters. We apply an additional layer of filtering to every detection before it is considered a measurement. We check if a new detection is in collision with the static map, and if it is in collision, we reject that particular detection. The check against the static map is fast and helps reduce false positives in practice. We use Nearest Neighbor (NN) data association [5], which is a reasonable compromise between performance and computational cost.

Depending on the task, a single person or multiple people must be tracked. We examine each case below:

- **Single target tracking:** For some tasks, such as person following, dedicated tracking of a single specific user is required and tracking bystanders is not required for task success. In this case, our goal is to keep tracking the specific user, so we significantly relax the detection thresholds of the detectors. Even though doing so results in more spurious detections, we do not start more than a single track. This approach improves the tracking performance of a single person.
- **Multi-target tracking:** When the robot is navigating to a goal point with human bystanders, tracking multiple people at the same time is necessary. Moreover, losing track of a bystander would not be very detrimental to task success. We keep a separate Kalman filter for each tracked person. If a detection is matched to multiple filters, only the closest filter is associated with the detection and the other filters are considered to have no detections for that time step.

4.4 Face Recognition

For certain interactive navigation tasks such as finding a specific person, a robot needs to have person recognition capability. Our person recognition approach uses face recognition and optionally shirt color features. We detect faces in RGB images using the popular face detector by Viola and Jones [126]. We use the Eigenface method by Turk and Pentland [123] for face recognition. Our approach allows new faces to be trained on-the-fly.

With the *Eigenface* approach, faces are represented in a lower-dimensional space. Sirovich and Kirby [111] showed that dimension reduction method Principal Component Analysis (PCA) can be used on face images to form a set of basis features. The main idea of PCA for faces is to find vectors that best account for variation of face images in all training images. These vectors are called *eigenvectors*. Then a face space is constructed called *eigenfaces* and the images are projected onto this space. Our approach of face recognition works as follows:

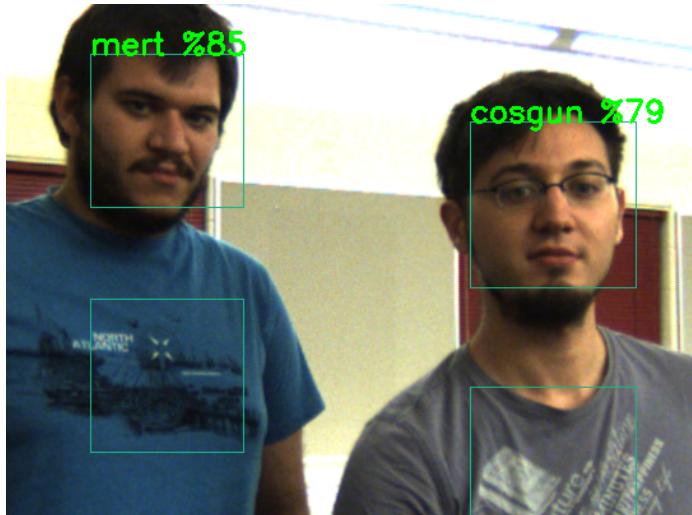


Figure 31: Example results of our person recognition method is shown in the image. We use *Eigenfaces* face recognition method and optionally shirt color recognition.

1. A person unknown to the system comes up to the robot and initiates training.

2. Robot asks the person to turn his face one side to another, and takes M face and shirt images of this person.
3. Eigenfaces from the entire training set is calculated, and every known face is projected to the corresponding M-dimensional weight *facespace*.
4. After training is completed, face recognition is reactivated.
5. A distance value from face recognition and optionally from shirt color recognition is received and it is thresholded for a decision. An example recognition result is in Figure 31.

Using the UI of the robot, a user can start training and adjust the information in the person database. The person data is managed by a SQLite database hosted locally on the robot.

Shirt color recognizer can be used when there is little time between the training and recognition. Activating the shirt recognition should improve recognition and reduce false positive detections. We assume a rectangular region below the face captures the shirt (1.5 times below the the face rectangle size). The distribute the histogram into bins using normalized RGB color space because of its relative robustness to lighting. For detection, we calculate the distance between the training histogram to the test histogram using Earth Mover Distance [99]. The color histogram is adaptively updated at every high confidence detection in order to account for illumination changes. The overall person score is calculated by a weighted average of face and shirt distance.

CHAPTER V

PERSON FOLLOWING

In this chapter, we focus on one aspect of human-robot interaction, namely person following in an indoor environment. There are many scenarios in which a person following robot can be useful. For example, a robot can carry luggage of travelers in airports, or groceries in a supermarket. Person following is also the enabling capability for interactive acquisition of the *Home Tour* scenario discussed in Section ???. The robot needs to know how to follow a person before building an environment representation and providing services to the user. There are two properties a person following behavior should achieve: robust following and social awareness.

Typically, a service robot operates in a dynamic and populated environment, therefore the robot must be able to keep track of a single person even when they are temporarily occluded. Multimodal person tracking that is presented helps the robot to have better estimates of a user's position. As discussed in Chapter 4, for the person following task, the robot has to track a designated user, and the detection thresholds of detectors are relaxed for robust tracking at the expense of more false positive detections.

The robot not only has to keep appropriate distance to the user, but it also has to recognize *what* the user is trying to accomplish and move accordingly. For example, during the home tour scenario, when the user stops, the robot should predict that the user is going to annotate a landmark, and it should come beside the user instead of standing behind. Moreover, the robot should be smarter when passing doors or following a person who is cutting a corner. In order to be able to handle these scenarios, the robot should act beyond purely reactive following behavior. It is

desirable for the robot to anticipate what action the user is likely going to take and act accordingly.

In this chapter, after referring to previous studies on person following robots in Section 5.1, we present the elementary person following method in 5.2. After that, in Section 5.3 we demonstrate situation awareness of a person following robot in three commonly encountered scenarios: door passing, user activity awareness and handling corners. In Section 5.4, we present an application of person following to telepresence robots.

5.1 Related Work

A robot that follows a person is a widely studied scenario in robotics. A relevant body of work is pursuit evasion [17], however the target is trying to evade the follower. In person following robots, we assume that the target is cooperating with the follower.

In one of the earliest works in this area by Sidenbladh [109], robot keeps the person centered in the camera image using a P controller. Prassler [94] also offers a fully reactive approach using the *Velocity Obstacles* concept, which uses the velocity of the target to find allowable velocities of the robot that guarantees avoiding collision if both the target and the robot move with constant speed. The approach is applied on a wheelchair, even though social constraints were not considered. The robot does not necessarily follow a person directly from behind. Gockley [32] observed how people walk together. It is reported that partners who were conversing tend to look forward with occasional glances to each other. Ohya [86] presents a following method to escort a target on the side while avoiding obstacles. It was assumed that the target would move with the same acceleration and velocity. Murakami [82] presents a method to first estimate the sub-goal of the leading person and then following as if the robot knows the goal.

Park [91] models the problem as a control problem and offers an algorithm based

on Model Predictive Control. [76] employs randomized tree expansion and biases the calculated paths towards a sub-goal, which is the current position of the person. Hoeller [40] adopts the virtual targets idea and selects a goal position in a circular region around the person. Stein [115] proposes choosing and following a leader to handle navigation in crowds.

Some of the relevant works considered the social side of the interaction. Gockley compared two elementary following methods: direction following, where the robot always attempts to drive towards the tracked person, and path following, where the robot follows the path the person took. It was shown that direction following behavior was perceived as more human-like and natural than path following. Yuan's [131] system switches following behavior between parallel following, direction and path following depending on the layout of the obstacles. Zender [132] emphasizes on situation awareness for following and studies cases of handling doors and corridors. To handle the doors, the robot increases its following distance and that leads the robot to wait for a while. Following in a corridor is handled with an approach similar to Pacchierotti [89], and the robot's speed is increased. [67] presents a system that is capable of responding to verbal and non-verbal gestures and follow a person. When following a group of people, a common method is to choose and follow a leader. Granata [34] presents behaviors such as going towards, following and searching a user. Ota [87] touches upon the recovery functions whenever the robot loses tracks of the leading person.

5.2 Basic Person Following

In this section, we describe our basic person following method. To keep track of the person

When the following behavior is initiated from a higher level process, first the person to be followed must be tracked. The robot looks for the closest person in the

vicinity of the robot (within 2m). If no person is detected for some time, then the command is invalidated. We use the person detection and tracking system described in Section 4.

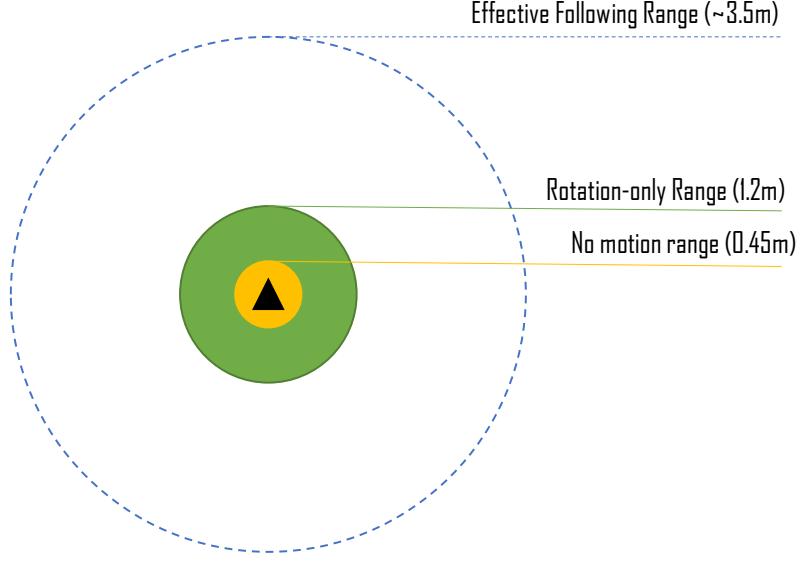


Figure 32: Overhead view of relevant ranges for person following. Robot is represented as the triangle in the middle.

In the basic person following mode, the robot has three different strategies depending on the distance towards the followed person. The distance to the user is calculated as the distance from the center of the robot base to the person's current location estimation. We used Hall's characterization of personal spaces in order to determine the distance limits. See Section 1.3 for a review of Hall's work. The three distinct zones and corresponding robot behavior are given as follows:

Intimate Zone [0 – 0.45m]: In this short interval, the person is very close to the robot, therefore any motion of the robot may be potentially unsafe. Therefore the robot comes to a complete halt in this zone ($v = 0, w = 0$). This behavior also allows the user to safely interact with the on-board User Interface.

Personal Zone [0.45 – 1.2m]: When the robot is in the personal zone of the user, the robot stops and only rotates towards the followed person ($v = 0, w = w_{rotation}$). The rotational velocity is determined with a P controller, with the error term defined

as the difference between the current orientation and the tracked person’s orientation. The rotational velocity is capped at a fixed value, so that the person feels comfortable with the rotation of the robot.

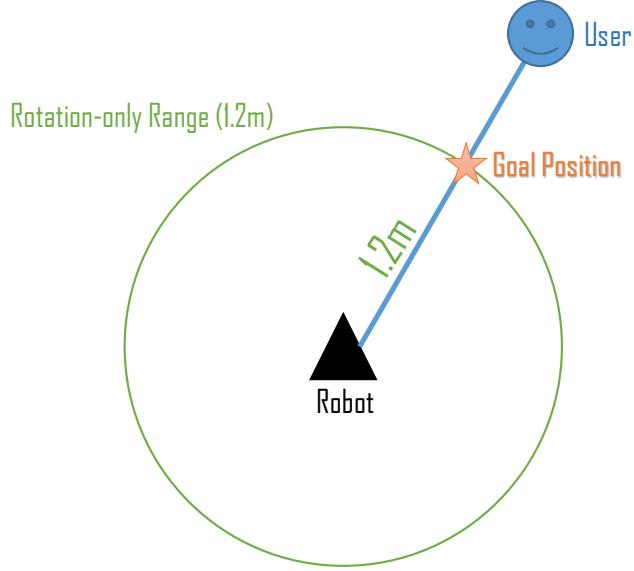


Figure 33: An illustration of how the goal position is calculated when the user is in the social space [1.2m – 3.5m].

Social Zone [1.2m – 3.5m]: In this range interval, the robot executes the main following behavior. At every time step, a goal pose that is $1.2m$ away from the user and headed towards the user is calculated, see Figure 33 for an illustration. A collision-free path is found and the robot executes this path until a new measurement from the tracked person is received. The path is found using Dynamic Window Approach (DWA) local planning method. We use the ROS implementation of DWA for the basic following behavior.

Sometimes it is inevitable that the person tracking system loses the target, particularly when the person is consistently faster than the robot or the person goes outside the range of the sensors (further than $\sim 3.5m$ in our case). When this happens, the robot will attempt to go to its last calculated goal position and look for the person. By this means, the robot attempts to keep up with the lost person as far as possible

with the hope that the person will re-appear in the vicinity of the last seen position. After the robot reaches this goal, it stops and waits for an amount of time. If the user is saved to the database, or the robot already knows that he/she is in the database, then the face recognition system described in Section 4.4 is activated. Otherwise, the robot continues following the closest person that appears in this position. If no person is detected within a fixed amount of time, 5 seconds in our implementation, then the robot declares that the person is lost.

5.3 Situation Aware Person Following

5.3.1 Door Passing

5.3.2 User Activity Awareness

5.3.3 Corners

5.4 Application To Telepresence Robots

A telepresence robot can be described as *Skype on wheels*, where a remote user teleconferences while having the control of the movement of a robotic system in a physical environment. Telepresence robots constitute a promising area in the consumer robotics industry as evidenced by multiple start-up companies working on telepresence products. However, currently all the telepresence robots that are available in the market are controlled by manual driving - usually via the keyboard or a joystick. In this section, we present an implementation of person following on a telepresence robot and a user study that evaluates effect of having person following capability on a telepresence robot.

Telepresence robots are a level above video conferencing since the robot is used as the communication medium and the remote user can now control the movement. Therefore, the spatial interaction between people and a telepresence robot in social situations is worth investigating. One of those situations is moving with a group of people. In an effort to analyze the spatial and verbal interaction, we focus on

engagement with one person where the remote user interacts with the person while following him/her in a corridor. This situation is very likely to happen in office environments, for example when the remote user is having a discussion with a co-worker while walking to his office after a meeting. As telepresence robots become more common, there will be need to have the functionality of autonomous following of a person so that the remote user doesn't have to worry about controlling the robot.

We evaluate our system by conducting a user study, where there are two following conditions:

1. Manual Person Following: Robot is controlled with an Xbox controller
2. Autonomous Person Following: Initiated by clicking on a user in RGB-D image

The aim of the user study is to measure how remote users like using the autonomous following feature compared to the manual. For the study, the remote user has a task that consists of listening to a passage the followed person reads and answering related questions after the interaction. We also observe subjects' experiences using the system, get useful feedback and pinpoint future challenges that can be helpful designing new applications for telepresence robots.

5.4.1 Robot Platform

The system described in this paper is implemented on an experimental telepresence robot shown in Figure 34. The robot has a differential drive base and can be used for about 8 hours with full charge. For the experiments in this paper, the speed of the robot was limited to 0.55 m/s. A laser scanner with 360° field of view, which was taken from Neato XV-11 vacuum cleaning robot, was mounted horizontally at 0.3m height. The system runs on Windows 7 and Microsoft Robotics Developer Studio (MRDS) as its distributed computing environment. The remote user connects to the robot via wireless internet and communicates with others using Skype. On the



Figure 34: The telepresence robot platform we used for our experiments.

remote end, . The robot is also equipped with an omni-directional microphone and a high-end speakerphone.

There are two operation modes for the robot: Teleoperation and Autonomous Person Following. A Xbox 360 Wireless Controller is used to remotely teleoperate the robot. A wide-angle camera is placed on top of the monitor and tilted slightly downward to help the remote user to see the floor, robot base and people's faces at the same time. A Kinect Sensor is also placed above the monitor. Person following is initiated through the user interface shown in Figure 35.

A modified version of the local planner used in Section 3.4.2 is used for person following. A utility function consisting multiple factors, including the respective position to the person, is optimized over multiple steps using Breadth-First Search. Details of the planning method can be found in [21].

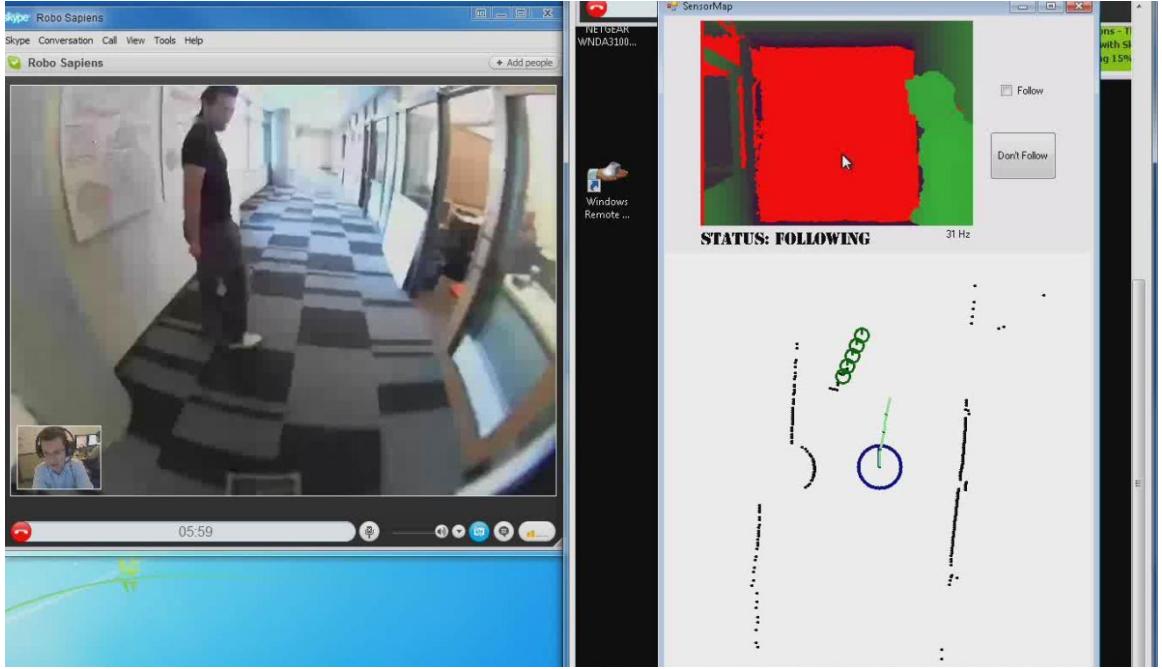


Figure 35: User Interface of the robot for the remote user.

5.4.2 User Study

In this study, remote user is the subject and the followed person is the experimenter. To investigate the effectiveness of using autonomous person following for an interaction task, we ran a controlled experiment and varied manual vs. autonomous following within subjects.

Design: The experiments were conducted in working hours and bypassers were allowed to walk across the experiment area or talk. The subjects were given the task of following the experimenter through the course for a lap and listen to the passage he is reading. In the first run, the subject used the autonomous following or teleoperation method to follow a person and complete the lap. In the second run, the subject used the other method. At the end of each run, the subject was asked to complete a 4-question quiz about the passage. The passages and quiz questions were taken from Test of English as a Foreign Language (TOEFL) listening section examples. One passage was about “behaviorism” and the other one was about “manila

hemps”, and passages were chosen so that they are at a similar difficulty level. The time it takes to read a passage corresponded approximately to the same time a lap is completed. We also asked numbered 7 point Likert scale questions, administered after each run, about how *Understandable* the experimenter was, *Easiness of UI*, if the robot exhibited *Natural Motions*, how *Safe* the remote user felt, if the subject was able to *Pay Attention* to the passage, how *Fast* the robot was and how much *Fun* the subject had. At the end of both runs, the user was asked which method he/she will prefer over the other for this type of a scenario. The exact questionnaire was show in Appendix C.

Participants: 10 volunteers participated in the study (6 male and 4 female between the ages of 25-48). Participants consisted of 4 researchers and 6 interns at Microsoft Research. 5 of the participants had little knowledge, 4 had average knowledge and 1 had above average knowledge on robotics. The participants weren’t gamers: 4 participants never played console games, 4 played rarely, 1 sometimes played and 1 often played. 6 of the participants often used video conferencing software, while 2 sometimes and 2 rarely used. 9 of the participants were not native English speakers and all of them had taken the TOEFL before. Participants were recruited through personal relations and were given a small gift (valued at approximately US\$10) for their help.

Procedure: The participants were first greeted by the experimenter and instructed to complete a pre-task questionnaire regarding their background. The robot was shown to the participant and basic information about its capabilities was told. The experimenter explained the task while walking with the participant in the corridor and showing the course to be followed. Participants were told that they should stay close to the experimenter while he is walking and there will be a quiz regarding the passage afterwards. The participant was informed that there are 2 operation modes:

manual and autonomous person following.

Before the experiment started, the participant went through training for about 15 minutes. First, the participant learned the basic controls for the Xbox controller when he/she was nearby the robot. Then the participant was taken to the remote station, which was in a room about 20 meters away from the corridor area. The participant was informed about the UI and was shown how the autonomous following can be activated. Then a test run was executed, where the remote user followed the experimenter via teleoperation and had a conversation.

After the training, the actual run was executed using either the manual or autonomous method. When the lap was completed, first the passage quiz, then the survey questions were answered by the subject. Then the second experiment using the other method was executed, and the second passage quiz and survey questions were given to the subject. As the last question, the subject was asked to state his/her method of preference. Lastly, the participants were debriefed about the study and engaged in a discussion. We switched the starting method for every other experiment in order not to bias the subjects' opinions about one particular method.

Measures: We had three measurement criteria to compare manual vs autonomous following: 1) Number of correct answers to passage quizzes: Assuming the standardized TOEFL exercises were of same difficulty, we ran a paired *t – test* on two groups of autonomous and manual. 2) Survey questions: We ran a paired *t – test* using 7-point Likert Scale on each of the seven questions. 3) Preferred Method: We looked at which method subjects chose over the other one.

Results: Out of 4 quiz questions, the correct answers for autonomous group ($\mu = 2.9$, $\sigma = 0.9$) were more than the manual group ($\mu = 2.2$, $\sigma = 1.2$) but the statistical difference was not statistically significant ($t(9) = 1.48$, $p = 0.17$ on *t – test*).

Table 6 summarizes the survey results. For *Understandable* and *Fun*, the scores

Question	Autonomous		Manual Drive		$t - test$	
	μ	σ	μ	σ	p	t
1. Understandable	4.0	1.5	3.6	1.7	0.47	0.73
2. Easy UI	6.5	0.9	5.0	2.2	0.06	2.13
3. Natural Motion	5.4	1.0	3.5	1.9	0.03	2.52
4. Safe	5.1	1.7	2.3	1.4	0.01	3.09
5. Pay Attention	5.3	1.8	3.4	1.5	0.02	2.63
6. Fast	3.9	0.3	4.3	0.8	0.10	-1.8
7. Fun	5.3	1.5	5.1	1.7	0.66	0.45

Table 6: Survey results of the user study for person following for telepresence robots. Table displays survey question average and standard deviations for the two conditions: Autonomous Person Following and Manual Person Following.

slightly favored autonomous method but the difference wasn't statistically significant. Manual method User Interface (gaming controller) was found to be easy to use ($\mu = 5.0, \sigma = 2.2$), but the UI for autonomous method (clicking) was found to be marginally easier ($\mu = 6.5, \sigma = 0.9$), ($t(9) = 2.13, p = 0.06$). The motions of the robot was found to be significantly more *Natural* to have a conversation for autonomous ($\mu = 5.4, \sigma = 1.0$) than manual ($\mu = 3.5, \sigma = 1.9$), ($t(9) = 2.52, p = 0.03$). Participants thought they were able to *Pay more Attention* to the passage the experimenter is reading when the robot was following the him autonomously ($\mu = 5.3, \sigma = 1.8$) compared to manual control ($\mu = 3.4, \sigma = 1.5$) and the statistical difference was significant ($t(9) = 2.63, p = 0.02$). Participants have found the autonomous method ($\mu = 5.1, \sigma = 1.7$) much safer than manual method ($\mu = 2.3, \sigma = 1.4$) and there was a significant difference between two groups ($t(9) = 3.09, p = 0.01$). The speed of the robot was found to be neither fast nor slow for both methods ($\mu = 3.9, \sigma = 0.3$) and ($\mu = 4.3, \sigma = 0.8$).

All 10 subjects chose autonomous person following over teleoperation for this task.

5.4.3 Design Implications

Our user study showed that a person following behavior is desirable for telepresence robots when there is interaction. The follow-up discussions also agreed with the survey

results, as one subject (R10) stated: “*It just gives me more focus and concentration.*” Below, we list our observations and implications for future research and design for telepresence robots:

Motor Noise: Even though the motors on the robot were relatively quiet, 8 out of 10 participants expressed that the motor noise made communication harder. This justifies the close scores we collected in the survey question asking if the subject was able to understand what the experimenter was saying. (R8) was disturbed by the noise: “*When I was driving, it was always this constant sound. It was worse for the autonomous one. It was constantly adjusting and compensating for the movement.*” On the other hand, (R5) found the motor noise useful: “*I actually like it because it gives me the feedback whether I’m driving faster or slower. It also gives me a little bit feeling of life.*” Thus, although excessive motor noise should be avoided, some noise might be useful.

Wireless Connection: Second most cited problem for video conferencing was the video quality and time lags. (R8) clearly expressed why it was hard to walk with the experimenter using the manual method: “*The frame rate drops all of a sudden and you have no choice but to stop.*” Another subject (R9) made use of the displayed sensor data when the video conferencing quality went bad: “*Because of the lag, I just switched to the Kinect (depth image) and the overhead view (laser).*” This was possible because the wide angle camera image was coming from Skype whereas sensor displays were received from the Windows Remote Assistance. Clearly, a big challenge for telepresence systems is to deal with wireless connection problems.

Natural Interaction: Even though the participants thought the motions of the robot were natural to have a conversation ($\mu = 5.4$, $\sigma = 1.0$), some didn’t feel it was a natural way to communicate. As seen in Figure 34, the screen displaying the remote user’s face is flat and it introduced problems when the robot was traveling on the side of the person. (R5), when asked about walking side by side: “*..we don’t*

have face-to-face. It is not really a conversation.” This raises design considerations on how the remote user’s face is brought out. One of the subjects (R5) discovered that the microphone characteristics are different than human hearing: “*I don’t have a distance sense if the experimenter is further away or close. If you have the fading audio, then I’ll immediately notice.*” Whether a telepresence robot should exhibit the same characteristics of human perception or not is an open question and needs further investigation.

Assisted Teleoperation: Telepresence robots should possess a layer to assist the remote user to avoid obstacles and collisions. *Safety* ratings for the manual method were very low ($\mu = 2.3$, $\sigma = 1.4$) and (R8) expressed the concern: “*I was especially worried about running into the experimenter.*” This suggests that scenarios involving interaction would demand more attention of the remote users. The teleoperation should also be intuitive and be similar to driving modalities that people are already used to. (R4) stated: “*I was thinking about Manual mode compared to driving a car.*” before suggesting “.. maybe something like a cruise control might be good.”

Gaming Experience: Since the robot was controlled by a gaming console controller, some participants likened the manual mode to gaming. (R9) said: “*Manual is like playing video games.*” and (R5) said: “*I don’t play video games so controlling those consoles is not natural to me.*” Thus, it is possible that gamers are less likely to have trouble driving the robot. This observation is also made by Takayama [117].

Long Term Interaction: None of the subjects participated in our study had used a telepresence robot before. (R6) justified the inability to use the manual method: “*Maybe if I have some more practice for about several hours of driving the robot, I can use manual as well as autonomous.*” (R8) on having fun using teleoperation: “*It was fun because it was the first time I did it but I can imagine that over time, I’ll get bored of it.*” The *Fun* question in the survey received similar scores for autonomous

and manual, possibly because using a telepresence robot was a new experience for the subjects. Studies regarding long term interaction for telepresence robots can yield interesting results, as in [63].

Error recovery: When the person was lost during following, the UI displayed a text that the person was lost so that the remote user can re-initiate the following by clicking on the person. None of the subjects complained about the robot losing the person. When asked explicitly about the robot losing the experimenter, (R10) answered: “*That’s not a big deal in comparison to me driving the robot.*” Therefore, applications developed for telepresence robots can take advantage of the human being in the loop and does not have to be error-free for deployment.

5.4.4 Discussion

User studies showed that autonomous person following is a desired capability for a telepresence robot and it was favored over direct teleoperation for an accompanying task. Autonomous following was found to be safer, easier to use and helped the remote users to pay more attention to the conversation instead of the robot control. From the experience earned from user studies, there are still interesting challenges to explore in terms of human-robot interaction for telepresence robots.

CHAPTER VI

PERSON GUIDANCE

The application of guiding a person to a location with a robot has many uses, such as giving tours in museums, showing a visitor locations or individuals of interest and helping the elderly or the visually impaired go to places. We think the guidance behavior is one of the most fundamental capabilities a socially interactive robot should have.

A straightforward approach to guide a person would be the stop-and-wait method: robot plans a path and executes it normally as long as the guided person is nearby, and stops then the person is outside the radius. However, this method of guidance may lead to sudden stops and would not be socially accepted. A guide robot should consider the distance to the human and incorporate this information in its control strategy.

In this chapter, after reviewing the literature on guide robots in Section 6.1, in Section 6.2 we present our guidance method, which adjusts the speed of the robot as a function of the distance to the person. Then we present a special scenario for guidance in Section 6.3, specifically tailored to be used by blind people.

6.1 Related Work

The earliest works in guide robots focused on the implementation and long-term deployment of tour guide robots in public places such as museums. Burgard [13] presents the robot Rhino, that was deployed to a museum for 47 hours. The Minerva robot was an improved model over Rhino [118], and was deployed to a museum with an order of magnitude larger floor space. This robot was in operation for two weeks, and it was able to have short-term interaction with people , i.e. head motion, facial

expressions. Siegwart [110] presents a robot that was deployed to an exhibition for 6 months. Nourbakhsh [84] presents a project where four guide robots were deployed to museums for a period of five years. The authors remark that it is indeed possible to deploy guide robots public places, unsupervised. All these tour-guide robots had various degrees of autonomy and was received with enthusiasm. However, in all of these works it was apparent that there is a need for research in the area of Human-Robot Interaction.

Pacchierotti [88] demonstrates an office guide robot, but the main focus is on passing people in corridors. Clodic [20] presents another robot deployed in a museum. It was reported that a continuous interaction all along the guiding mission is fundamental to keep visitor's interest. Martin [70] studies the scenario of guiding a visitor in an office environment and focuses on robust person tracking. Pandey [90] focuses on the leave-taking of the guided person. The robot predicts the intent of the discontinuation of the task and either breaks the mission or searches for the user depending on the waiting time. Martinez-Garcia [72] focuses on the scenario of guiding a group of people with multiple robots at the same time. Garrell [29] works on a similar problem, where the task of the two robots is to control group of people and guide them. Another relevant scenario is the evacuation scenario, in which there is a danger and robots guide people to the safe a location [52, 98].

6.2 Guide Robot

In this section, we describe our method of person guidance. After a request to guide a person is received from a higher level process, the robot first plans a path using our navigation planner presented in Chapter 3, or standard ROS Navigation. The robot continues on its path as long while constantly monitoring the distance between itself and the guided person. The robot adjusts its speed according to this distance.

We use a variable speed profile so that the robot can better keep up with the

person and the motions of the robot is smoother. We define a speed function that is a function of the distance between the robot and the user. This speed profile is shown in Figure 36. The robot moves at a low speed v_{safe} if the human is dangerously close. The speed is peaked at distance d_{peak} and the robot stops if the distance is larger than d_{guide} , which may indicate that the human is not interested in being guided. Note that v_{peak} is capped by the speed limits in the environment, provided by the speed map approach presented in Section 3.5.

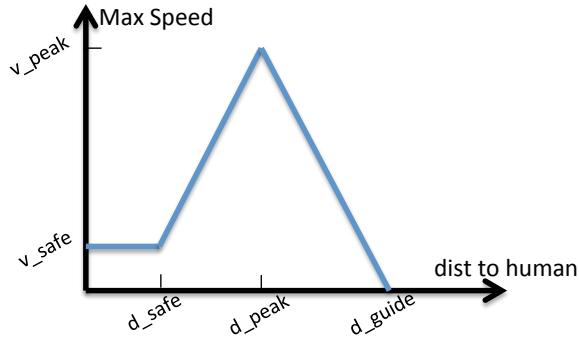
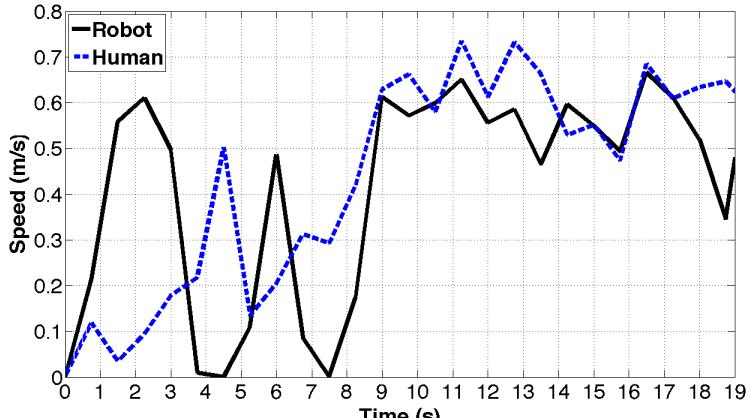


Figure 36: Speed profile of a person guiding robot as a function of the distance to the user.

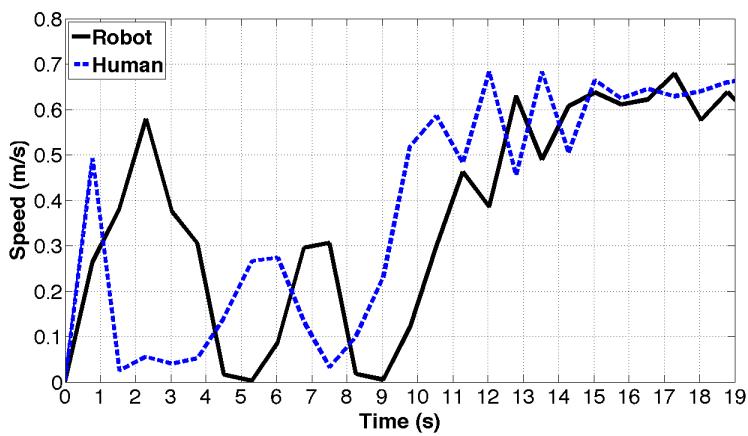
In the second scenario, the robot has the same goal point and guiding a person. In the first condition, we use ROS Navigation but robot stops if the distance to the human is over a threshold. In second condition, we use our method of dynamic speed adjustment for guidance. In this scenario, there is no person standing around the corner. We measured the instantaneous speeds of the robot and the human.

6.2.1 Results

The robot is given a fixed goal to guide a person, who is tracked with a torso-level laser scanner, by fitting an ellipse to the torso. We compared the velocity profile in Figure 37 with $d_{guide} = 1.7\text{m}$, $d_{peak} = 0.9\text{m}$, $d_{safe} = 0.1\text{m}$, $v_{safe} = 0.1\text{m/s}$, $v_{peak} = 1.0\text{m/s}$. We compared our approach with the simple strategy: If the human is closer than d_{guide} , then the navigation continues with a fixed max speed. Otherwise robot stops and waits.



(a)



(b)

Figure 37: Comparison of robot and human speeds with respect to time. a) Standard ROS Navigation b) Our approach: accelerations are less steeper than a), which employs the dynamic speed adjustment for guidance.

In the experiment, when guiding was enabled, the human first waited until the robot stopped at d_{guide} . Then he took a step and waited for a second time, and then started following the robot. The comparison of robot speeds is given in Figure 37(a) for fixed max speed and Figure 37(b) with the speed profile. Between $t = 0$ and $t = 9s$, the accelerations are much more rapid for the fixed max speed case. Robots that exhibit high accelerations will likely be perceived as unsafe, therefore our approach exhibits a more socially acceptable behavior. Moreover, after the person started following ($t > 9s$), our approach is better at mimicking the speed of the

human.

6.3 Application To Blind Users

In this section, we present our person guidance system specifically tailored for guiding blind users. Our approach consists of planning a path for the user and applying vibrations via a haptic belt to keep the user on the path.

6.3.1 Tactile Belt

In the previous section, we assumed that the guided person can detect where the robot is and follow him/her. With a blind user, this assumption does not hold. Therefore we need a mechanism to give directions to the user. Readily available options for assistive interfaces are limited to Braille or devices that presents content with speech synthesis. These ways of presenting information have difficulty dealing with representing spatial information. We also think visually impaired individuals would prefer a non-speech interface because they mostly rely on their sense of hearing in daily life. We therefore use a tactile belt for navigation guidance, because it can represent directions and rotations, be worn discreetly and does not occupy the hearing sense.

The belt has 8 pancake vibration motors, linearly spaced around the waist, and the motors can be controlled asynchronously via an Arduino board. We used two distinct vibration patterns to control the person's movements:

1. Directional Movement: When the guided person should move in a direction
2. Rotational Movement: When the guided person should turn around self

The vibration patterns that induce directional and rotational movements in the human can be specified in many ways. We evaluated four patterns in each category with a user study. The details of this user study as well as a survey about the usability of the Tactile Belt is provided in Appendix D. The user study showed that the

directional motion pattern with the highest recognition rate and least reaction time was the **TWO TAPS** pattern, which is illustrated in Figure 38. For the rotational motion pattern, we used the continuous rotation with a single motor, illustrated in Figure 39.

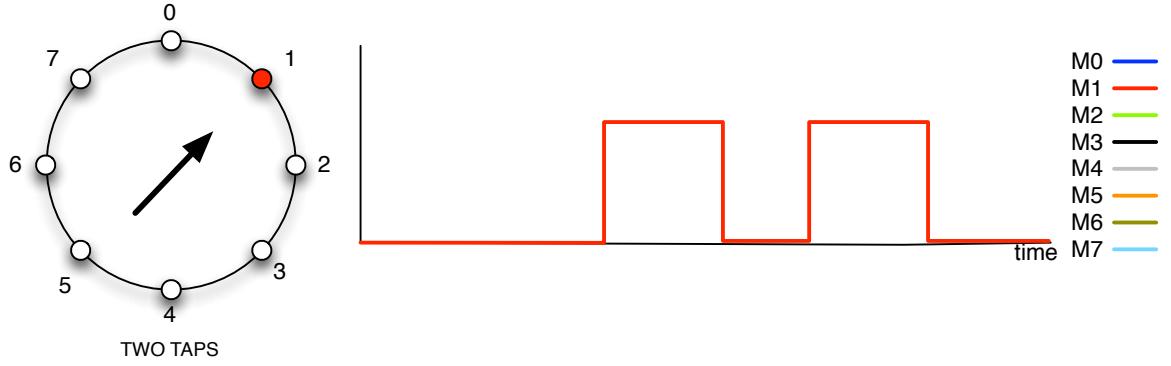


Figure 38: The vibration pattern applied by the Tactile Belt to induce directional movements in the blind user. A motor is fired for a duration of $250ms$, inactivated for $250ms$ and fired again for $250ms$.

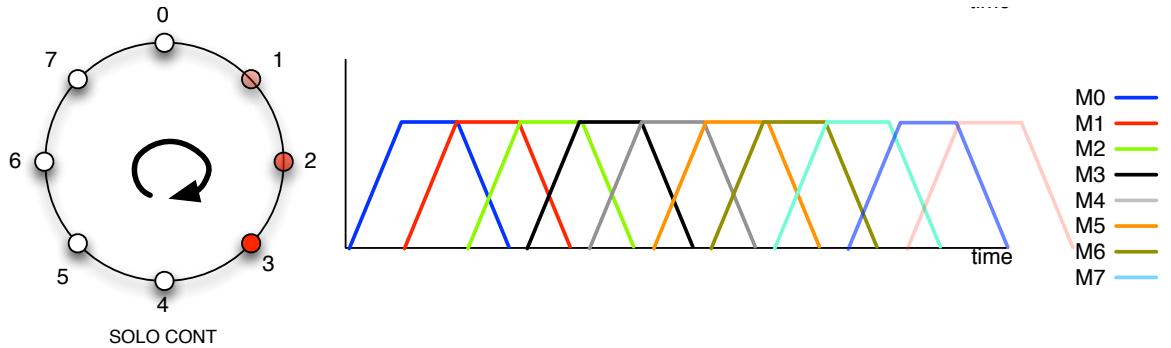


Figure 39: The vibration pattern applied by the Tactile Belt to induce rotational movements in the blind user. The consequent vibrations motors are fired consecutively, starting from left for CW and right for CCW rotation.

6.3.2 Planning the Path of the User

ROS provides an easy-to-use navigation stack for mobile robots. The input to the navigation stack is a map and a goal point and the output is a path and linear and angular velocities (v, w) necessary to keep the robot on the course of the path. We assumed that the human is a non-holonomic robot with a circular footprint. The

obstacle information is acquired from the laser scanner and the goal is provided in the sensor frame. Coupled with the Human Tracker, the 'robot' stays localized in the map and with respect to the path. The path is re-planned every second to deal with possible deviations. Next section is concerned with how the linear and angular velocities are converted to the vibration patterns.

6.3.3 Velocity to Vibration Mapping

Given a desired velocity that the 'robot' should execute, we first determine if a directional or rotational vibration pattern should be applied by the belt. If the linear velocity is dominant, then the human should walk towards that direction. If the angular velocity is dominant, the human should rotate around self. If both the linear and angular velocity is close to zero, the human should not move. To calculate which motion is appropriate, the 'robot' is simulated using Equation 6. If the distance the 'robot' took is larger than a threshold, then a directional vibration pattern is used. If it is less than the threshold, a rotational pattern is used. If both of the velocities are small enough, the no vibration is applied.

When the human gets to the vicinity of the goal point, a special stop signal is applied to inform the person that the destination is reached. Stop signal is implemented similar to **TWO TAPS** pattern except all the motors are activated instead of one.

6.3.4 Demonstration

We demonstrated that our system can successfully guide a blindfolded person to a goal location in a room. Based on our evaluation results of vibration patterns, we used **CONT** for directional motions and **SOLO CONT** for rotational motions. The experimenter manually provided several goal poses using the GUI. Note that since the system is re-planning frequently, the planner is able to accommodate dynamic obstacles and compensate unpredictable motions of the person. The demonstration steps as well as their explanations are shown in Figure 40.

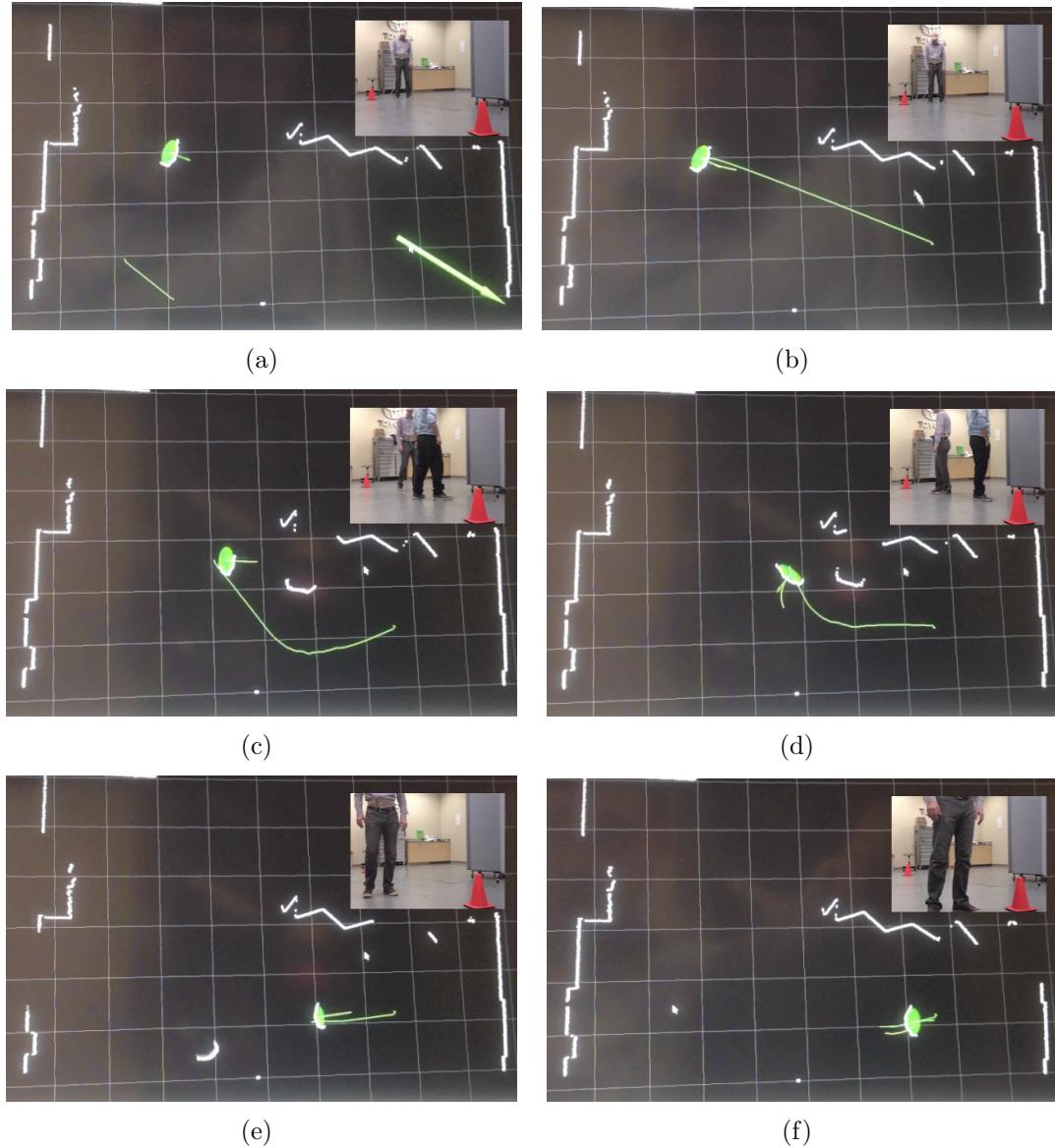


Figure 40: Autonomous guiding of a blindfolded person using the tactile belt. a) The guidance starts. The user is blindfolded and is standing at the left of the screen. The human detection system detects him and places an ellipse marker with an arrow depicting his orientation. The operator gives a goal point by clicking on the screen. The goal point is the right traffic cone, and given by the big arrow. b) The system autonomously generates a path for the user. As seen in the picture the path is collision free. At this stage the belt begin to vibrate towards the front of the user. c) An unexpected obstacle (another person) appears and stops in front of the user. The system detects the other person as an obstacle, and reevaluates the path. A new path going around the obstacle is immediately calculated and sent to the user by the belt. d) The user receives a rotation vibration modality, and begins to turn towards the new path. And follows this path from now on. e) The obstacle leaves. The path is then reevaluated and changed. The user receives forward directional belt signal, and advances towards the goal. f) The person reaches to the vicinity of the goal and stop signal is applied.

CHAPTER VII

CONCLUSION

Conclusion

APPENDIX A

QR CODE BASED LOCATION INITIALIZATION

QR Code Based Location Initialization

APPENDIX B

ASSISTED REMOTE CONTROL

Assisted Remote Control

APPENDIX C

TELEPRESENCE STUDY SURVEY QUESTIONS

APPENDIX D

VIBRATION PATTERN ANALYSIS FOR HAPTIC BELTS

In this Chapter, we provide the user study for vibration patterns to be used on the Tactile Belt described in Chapter 6.

D.1 Vibration Patterns

We define two main classes of vibration patterns depending on the type of intended human motion: directional and rotational. Directional patterns induce a motion towards a direction. Rotational patterns induce a rotation around self, which is intended to control the orientation of the human.

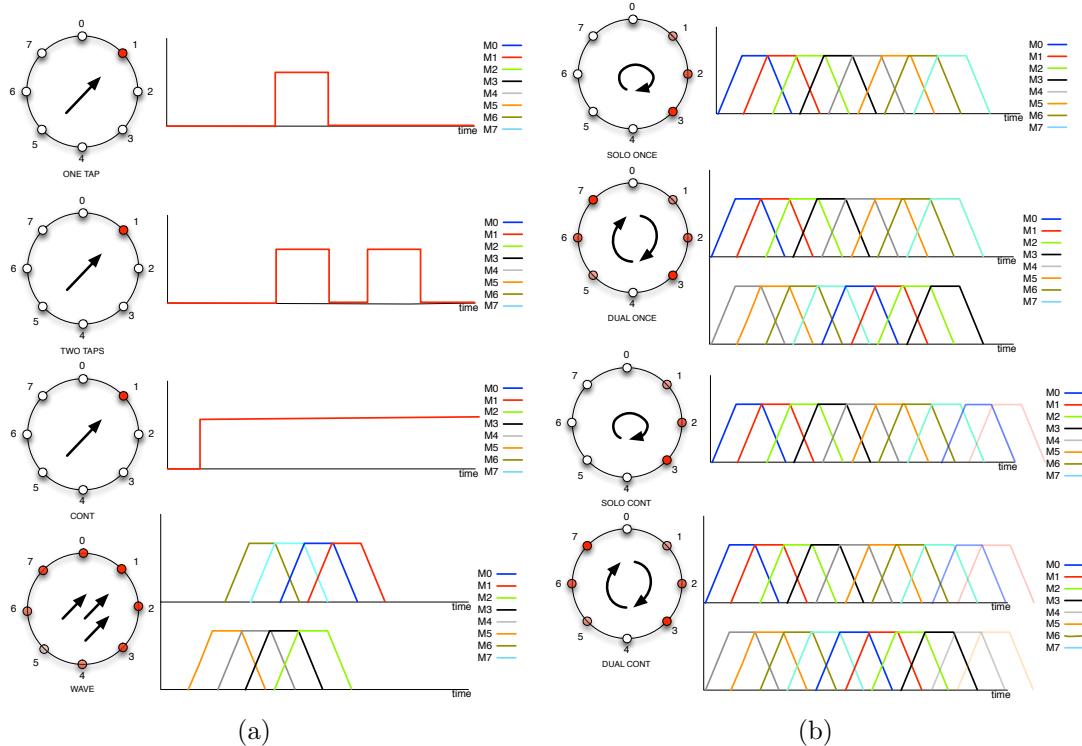


Figure 41: Evaluated vibration patterns. a) Directional b) Rotational

- **ONE TAP:** a motor is active for 250ms

- **TWO TAPS:** a motor is active 250ms, inactive for 250ms and active for 250ms again
- **CONT:** a motor is active until a new pattern is applied
- **WAVE:** a feeling that starts from the opposite end of desired direction and ends in desired direction

Figure 41(a) illustrates all 4 directional patterns for towards northeast cardinal direction (Vibration Motor 1).

- **SOLO ONCE:** activates all 8 motors consecutively, starting from left motor for clockwise, right for counter-clockwise.
- **SOLO CONT:** repeats **SOLO ONCE** pattern
- **DUAL ONCE:** circle motion is executed for one full circle with two opposing motors instead of one
- **DUAL CONT:** repeats **DUAL ONCE** pattern

Figure 41(b) illustrates all 4 rotational patterns for a clockwise rotation motion.

D.2 Procedure

The subject first went through a training for 5 minutes where experimenter applied directional and rotational patterns and told the correct direction/rotation. The subject is instructed to walk randomly in a confined area so that the patterns are tested while the human is in motion. The subject was asked to press a button on a Xbox controller whenever he/she decides on the direction/rotation and then tell the direction (coded 1-8) or rotation (left or right). The experiments consisted of 4 studies:

1. 5 samples from each of the **ONE TAP**, **TWO TAPS**, **WAVE** directional patterns are applied in random order.

2. 8 samples from **CONT** directional pattern are applied.
3. 2 samples from each of the 4 rotational patterns are applied.
4. 2 samples from each of the 4 rotational patterns while a random directional motion pattern is applied in between samples.

Upon the completion of the experiment, the experimenter applied all vibration patterns one by one and asked which directional and rotational pattern the he/she would prefer. Then, a post-study survey is conducted.

We hypothesize that subjects will prefer a one-time signal pattern in Experiment 1 to continuous pattern in Experiment 2 because a long lasting vibration may be annoying to the users. Our second hypothesis is that applying intermediate directional patterns will reduce the recognition rate of rotational patterns. This would be tested by comparing Study Experiment 3 and Experiment 4 results.

D.3 Measures

3 metrics were used for evaluation:

- **Recognition error (RE):** Angle difference between applied and perceived direction
- **Recognition accuracy (RA):** Percentage of correct recognition
- **Reaction time (RT):** Time between the start of the pattern to the instant the subject decides on an answer. A timeout occurs if subject can not give an answer in 5 seconds.

The post-survey consisted of 4 usability questions on a 10-point Likert scale. Participants were also asked for their preferred directional and rotational patterns.

D.4 Results

D.4.1 Directional Patterns

Recognition error and RT w.r.t. pattern type: A total of 344 directional pattern samples were sampled in Experiments 1 and 2. Since the directions are discretized, RE from any single test ranges from 0 to 180° with 45° increments. Results with respect to pattern type is in Table 7.

	ONE TAP	TWO TAPS	CONT	WAVE
Directional Error	12.4°	10.6°	8.4°	23.1°
Reaction Time (s)	1.32	1.13	1.26	1.92

Table 7: Average recognition error and reaction times of directional patterns

With a mean RE of 8.4°, **CONT** pattern was the most accurate directional pattern, followed by **TWO TAPS** pattern with 10.6°. Subjects recognized **TWO TAPS** the fastest by an average of 1.13 seconds. **WAVE** pattern performed significantly worse than others on both measures.

Recognition accuracy w.r.t. applied direction: The confusion matrix for RA of the applied directions is given in Table 42.

		Perceived Direction							
		1	2	3	4	5	6	7	8
Applied Direction	1	0.97	0.03	0	0	0	0	0	0
	2	0	0.59	0.38	0	0.03	0	0	0
	3	0	0.1	0.55	0.3	0.05	0	0	0
	4	0	0.04	0.04	0.76	0.16	0	0	0
	5	0	0.02	0	0.09	0.81	0.02	0	0
	6	0	0.02	0	0.02	0.04	0.84	0.08	0
	7	0.02	0	0	0	0	0.26	0.65	0.06
	8	0.03	0	0	0	0	0.03	0.14	0.80

Figure 42: Confusion matrix of recognition accuracy of directional vibration patterns. Our results show that the recognition accuracy is highly dependent on the applied direction. The subjects recognised the front direction (Direction 1) with the highest accuracy (%97) whereas Direction 3 had the least accuracy (%55)

Our results show that the recognition accuracy is highly dependent on the applied

direction. The subjects recognized the front direction with highest RA (%97) whereas Direction 3 (right) had the least RA (%55).

We expected the directional RA to be symmetrical around the belt, meaning that the accuracy of right and left directions should be equally sensitive to haptic feedback. However, there is at least %10 accuracy difference for all 3 such pairs. This could be due to imperfect alignment of motors in the belt prototype.

D.4.2 Rotational Patterns

	SOLO ONCE	DUAL ONCE	SOLO CONT	DUAL CONT
Recog. Accuracy	%100	%92	%100	%98
Reaction Time (s)	1.32	1.84	1.16	1.68

Table 8: Average recognition accuracy and reaction times of rotational patterns

A total of 256 rotational pattern samples were tested Experiments 1 and 2. Results are in Figure 8.

Subjects recognized **SOLO ONCE** and **SOLO CONT** patterns with %100 accuracy, whereas had **DUAL ONCE** had %92 RA. **SOLO CONT** had the least RT by 1.16 seconds. It is interesting to note that this reaction time is almost identical to the directional pattern with the least RT (1.13s).

D.4.3 Usability of Tactile Belt

Subjects thought that it was fairly easy to move while wearing the belt ($M=9.2$). Most subjects thought that the belt was comfortable ($M=8.3$) and it fit their waist well ($M=8.4$). The vibration motors were not found very silent ($M=6.8$). Actual questions in the post-study survey results are shown in Figure 43.

D.4.4 Discussion

Among directional patterns, **TWO TAPS** had the least RT and was found the most intuitive by our usability study. When asked which directional pattern they would

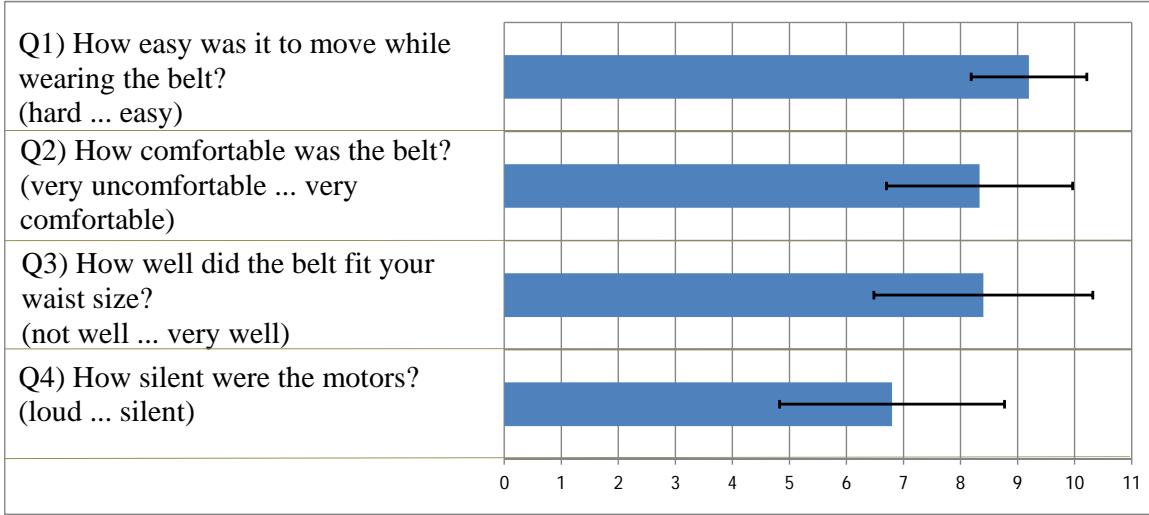


Figure 43: Results of post-study survey.

prefer, 10 out of 15 subjects chose **TWO TAPS**, 3 subjects chose **CONT** and 2 subjects selected **ONE TAP** and no one chose **WAVE**. Although **CONT** have less recognition error, most subjects preferred **TWO TAPS** probably because people didn't feel comfortable when the vibration lasted for a long time. This supports our first hypothesis that continuous patterns would not be found preferable.

7 out of 15 subjects preferred **SOLO ONCE**, 7 subjects preferred **SOLO CONT** and 1 subject preferred **DUAL CONT** in the post-study survey. Our results show that subjects rarely made mistakes in recognizing rotations and using a single motor for rotation patterns is preferable over using two motors. Therefore an application may use either of the **SOLO** patterns.

Our second hypothesis did not hold, as we found that application of other pattern between rotational patterns did not deteriorate the recognition rate.

REFERENCES

- [1] ALTHAUS, P., ISHIGURO, H., KANDA, T., MIYASHITA, T., and CHRISTENSEN, H. I., "Navigation for human-robot interaction tasks," in *Robotics and Automation, 2004 IEEE International Conference on*, vol. 2, pp. 1894–1900, IEEE, 2004.
- [2] ALY, A. and TAPUS, A., "An integrated model of speech to arm gestures mapping in human-robot interaction," in *Information Control Problems in Manufacturing*, vol. 14, pp. 817–822, 2012.
- [3] ARRAS, K. O., MOZOS, O. M., and BURGARD, W., "Using boosted features for the detection of people in 2d range data," in *Robotics and Automation, 2007 IEEE International Conference on*, pp. 3402–3407, IEEE, 2007.
- [4] AVRUNIN, E. and SIMMONS, R., "Using human approach paths to improve social navigation," in *8th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pp. 73–74, IEEE, 2013.
- [5] BAR-SHALOM, Y. and LI, X.-R., *Multitarget-multisensor tracking: principles and techniques*, vol. 19. YBS Storrs, Conn., 1995.
- [6] BAUMBERG, A. and HOGG, D., "Learning deformable models for tracking the human body," in *Motion-Based Recognition*, pp. 39–60, Springer, 1997.
- [7] BELLOTTO, N. and HU, H., "Multisensor-based human detection and tracking for mobile service robots," *Systems, Man, and Cybernetics, Part B: Cybernetics, IEEE Transactions on*, vol. 39, no. 1, pp. 167–181, 2009.
- [8] BENNEWITZ, M., AXENBECK, T., BEHNKE, S., and BURGARD, W., "Robust recognition of complex gestures for natural human-robot interaction," in *Proc. of the Workshop on Interactive Robot Learning at Robotics: Science and Systems Conference (RSS)*, 2008.
- [9] BENNEWITZ, M., BURGARD, W., CIELNIAK, G., and THRUN, S., "Learning motion patterns of people for compliant robot motion," *The International Journal of Robotics Research*, vol. 24, no. 1, pp. 31–48, 2005.
- [10] BERNARDIN, K. and STIEFELHAGEN, R., "Audio-visual multi-person tracking and identification for smart environments," in *Proceedings of the 15th international conference on Multimedia*, pp. 661–670, ACM, 2007.

- [11] BLODOW, N., MARTON, Z.-C., PANGERJCIC, D., RÜHR, T., TENORTH, M., and BEETZ, M., “Inferring generalized pick-and-place tasks from pointing gestures,” in *IEEE International Conference on Robotics and Automation (ICRA), Workshop on Semantic Perception, Mapping and Exploration*, 2011.
- [12] BROOKS, A. G. and BREAZEAL, C., “Working with robots and objects: Revisiting deictic reference for achieving spatial common ground,” in *Proceedings of the 1st ACM SIGCHI/SIGART conference on Human-robot interaction*, pp. 297–304, ACM, 2006.
- [13] BURGARD, W., CREMERS, A. B., FOX, D., HÄHNEL, D., LAKEMEYER, G., SCHULZ, D., STEINER, W., and THRUN, S., “The interactive museum tour-guide robot,” in *AAAI/IAAI*, pp. 11–18, 1998.
- [14] BUYS, K., CAGNIART, C., BAKSHEEV, A., DE LAET, T., DE SCHUTTER, J., and PANTOFARU, C., “An adaptable system for rgb-d based human body detection and pose estimation,” *Journal of Visual Communication and Image Representation*, 2013.
- [15] CARBALLO, A., OHYA, A., and YUTA, S., “Fusion of double layered multiple laser range finders for people detection from a mobile robot,” in *Multisensor Fusion and Integration for Intelligent Systems, 2008. MFI 2008. IEEE International Conference on*, pp. 677–682, IEEE, 2008.
- [16] CHENG, K. and TAKATSUKA, M., “Hand pointing accuracy for vision-based interactive systems,” in *Human-Computer Interaction-INTERACT 2009*, pp. 13–16, Springer, 2009.
- [17] CHUNG, T. H., HOLLINGER, G. A., and ISLER, V., “Search and pursuit-evasion in mobile robotics,” *Autonomous robots*, vol. 31, no. 4, pp. 299–316, 2011.
- [18] CIOPOLLA, R. and HOLLINGHURST, N. J., “Human-robot interface by pointing with uncalibrated stereo vision,” *Image and Vision Computing*, vol. 14, no. 3, pp. 171–178, 1996.
- [19] CLARK, H. and BRENNAN, S., “Grounding in communication,” *Perspectives on socially shared cognition*, vol. 13, no. 1991, pp. 127–149, 1991.
- [20] CLODIC, A., FLEURY, S., ALAMI, R., CHATILA, R., BAILLY, G., BRETHES, L., COTTRET, M., DANES, P., DOLLAT, X., ELISEI, F., and OTHERS, “Rackham: An interactive robot-guide,” in *Robot and Human Interactive Communication, 2006. ROMAN 2006. The 15th IEEE International Symposium on*, pp. 502–509, IEEE, 2006.
- [21] COSGUN, A., FLORENCIO, D. A., and CHRISTENSEN, H. I., “Autonomous person following for telepresence robots,” in *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 4335–4342, IEEE, 2013.

- [22] DIOSI, A., TAYLOR, G., and KLEEMAN, L., “Interactive slam using laser and advanced sonar,” in *Robotics and Automation, 2005. ICRA 2005. Proceedings of the 2005 IEEE International Conference on*, pp. 1103–1108, IEEE, 2005.
- [23] DROESCHEL, D., STUCKLER, J., and BEHNKE, S., “Learning to interpret pointing gestures with a time-of-flight camera,” in *Human-Robot Interaction (HRI), 2011 6th ACM/IEEE International Conference on*, pp. 481–488, IEEE, 2011.
- [24] DROESCHEL, D., STUCKLER, J., HOLZ, D., and BEHNKE, S., “Towards joint attention for a domestic service robot-person awareness and gesture recognition using time-of-flight cameras,” in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, pp. 1205–1210, IEEE, 2011.
- [25] FITZGIBBON, A., PILU, M., and FISHER, R. B., “Direct least square fitting of ellipses,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 21, no. 5, pp. 476–480, 1999.
- [26] FOX, D., BURGARD, W., and THRUN, S., “The dynamic window approach to collision avoidance,” *Robotics & Automation Magazine, IEEE*, vol. 4, no. 1, pp. 23–33, 1997.
- [27] GANAPATHI, V., PLAGEMANN, C., KOLLER, D., and THRUN, S., “Real time motion capture using a single time-of-flight camera,” in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pp. 755–762, IEEE, 2010.
- [28] GARCIA-SALICETTI, S., BEUMIER, C., CHOLLET, G., DORIZZI, B., LES JARDINS, J. L., LUNTER, J., NI, Y., and PETROVSKA-DELACRÉTAZ, D., “Biomet: a multimodal person authentication database including face, voice, fingerprint, hand and signature modalities,” in *Audio-and Video-Based Biometric Person Authentication*, pp. 845–853, Springer, 2003.
- [29] GARRELL, A. and SANFELIU, A., “Local optimization of cooperative robot movements for guiding and regrouping people in a guiding mission,” in *Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on*, pp. 3294–3299, IEEE, 2010.
- [30] GERMA, T., LERASLE, F., OUADAH, N., and CADENAT, V., “Vision and rfid data fusion for tracking people in crowds by a mobile robot,” *Computer Vision and Image Understanding*, vol. 114, no. 6, pp. 641–651, 2010.
- [31] GLAS, D. F., MIYASHITA, T., ISHIGURO, H., and HAGITA, N., “Laser-based tracking of human position and orientation using parametric shape modeling,” *Advanced robotics*, vol. 23, no. 4, pp. 405–428, 2009.
- [32] GOCKLEY, R., “Developing spatial skills for social robots.,” in *AAAI Spring Symposium: Multidisciplinary Collaboration for Socially Assistive Robotics*, pp. 15–17, 2007.

- [33] GOCKLEY, R., FORLIZZI, J., and SIMMONS, R., “Natural person-following behavior for social robots,” in *Proceedings of the ACM/IEEE international conference on Human-robot interaction*, pp. 17–24, ACM, 2007.
- [34] GRANATA, C. and BIDAUD, P., “A framework for the design of person following behaviors for social mobile robots,” in *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pp. 4652–4659, IEEE, 2012.
- [35] HALL, E. T., *The hidden dimension*. Anchor Books, 1966.
- [36] HARRIS, C. and STEPHENS, M., “A combined corner and edge detector,” in *Alvey vision conference*, vol. 15, p. 50, Manchester, UK, 1988.
- [37] HATO, Y., SATAKE, S., KANDA, T., IMAI, M., and HAGITA, N., “Pointing to space: modeling of deictic interaction referring to regions,” in *Proceedings of the 5th ACM/IEEE international conference on Human-robot interaction*, pp. 301–308, IEEE Press, 2010.
- [38] HELBING, D. and MOLNAR, P., “Social force model for pedestrian dynamics,” *Physical review E*, vol. 51, no. 5, p. 4282, 1995.
- [39] HICHEUR, H., VIEILLENT, S., RICHARDSON, M., FLASH, T., and BERTHOZ, A., “Velocity and curvature in human locomotion along complex curved paths: a comparison with hand movements,” *Experimental brain research*, vol. 162, no. 2, pp. 145–154, 2005.
- [40] HOELLER, F., SCHULZ, D., MOORS, M., and SCHNEIDER, F. E., “Accompanying persons with a mobile robot using motion prediction and probabilistic roadmaps,” in *Intelligent Robots and Systems, 2007. IROS 2007. IEEE/RSJ International Conference on*, pp. 1260–1265, IEEE, 2007.
- [41] HOSOYA, E., SATO, H., KITABATA, M., HARADA, I., NOJIMA, H., and ONOZAWA, A., “Arm-pointer: 3d pointing interface for real-world interaction,” in *Computer Vision in Human-Computer Interaction*, pp. 72–82, Springer, 2004.
- [42] HU, K., CANAVAN, S., and YIN, L., “Hand pointing estimation for human computer interaction based on two orthogonal-views,” in *Pattern Recognition (ICPR), 2010 20th International Conference on*, pp. 3760–3763, IEEE, 2010.
- [43] HÜTTENRAUCH, H., EKLUNDH, K. S., GREEN, A., and TOPP, E. A., “Investigating spatial relationships in human-robot interaction,” in *Intelligent Robots and Systems, 2006 IEEE/RSJ International Conference on*, pp. 5052–5059, IEEE, 2006.
- [44] JOJIC, N., BRUMITT, B., MEYERS, B., HARRIS, S., and HUANG, T., “Detection and estimation of pointing gestures in dense disparity maps,” in *Automatic Face and Gesture Recognition, 2000. Proceedings. Fourth IEEE International Conference on*, pp. 468–475, IEEE, 2000.

- [45] KAHN, R. E. and SWAIN, M. J., “Understanding people pointing: The perseus system,” in *Computer Vision, 1995. Proceedings., International Symposium on*, pp. 569–574, IEEE, 1995.
- [46] KALMAN, R. E., “A new approach to linear filtering and prediction problems,” *Journal of Fluids Engineering*, vol. 82, no. 1, pp. 35–45, 1960.
- [47] KEHL, R. and VAN GOOL, L., “Real-time pointing gesture recognition for an immersive environment,” in *Automatic Face and Gesture Recognition, 2004. Proceedings. Sixth IEEE International Conference on*, pp. 577–582, IEEE, 2004.
- [48] KEMP, C. C., ANDERSON, C. D., NGUYEN, H., TREVOR, A. J., and XU, Z., “A point-and-click interface for the real world: laser designation of objects for mobile manipulation,” in *Human-Robot Interaction (HRI), 2008 3rd ACM/IEEE International Conference on*, pp. 241–248, IEEE, 2008.
- [49] KENDON, A., *Conducting interaction: Patterns of behavior in focused encounters*, vol. 7. CUP Archive, 1990.
- [50] KHAN, Z., BALCH, T., and DELLAERT, F., “An mcmc-based particle filter for tracking multiple interacting targets,” in *Computer Vision-ECCV 2004*, pp. 279–290, Springer, 2004.
- [51] KHATIB, O., “Real-time obstacle avoidance for manipulators and mobile robots,” *The international journal of robotics research*, vol. 5, no. 1, pp. 90–98, 1986.
- [52] KIM, Y.-D., KIM, Y.-G., LEE, S. H., KANG, J. H., and AN, J., “Portable fire evacuation guide robot system,” in *Intelligent Robots and Systems, 2009. IROS 2009. IEEE/RSJ International Conference on*, pp. 2789–2794, IEEE, 2009.
- [53] KINNUNEN, T. and LI, H., “An overview of text-independent speaker recognition: From features to supervectors,” *Speech communication*, vol. 52, no. 1, pp. 12–40, 2010.
- [54] KIRBY, R., SIMMONS, R., and FORLIZZI, J., “Companion: A constraint-optimizing method for person-acceptable navigation,” in *Robot and Human Interactive Communication, 2009. RO-MAN 2009. The 18th IEEE International Symposium on*, pp. 607–612, IEEE, 2009.
- [55] KIRCHNER, N., ALEMPIJEVIC, A., and VIRGONA, A., “Head-to-shoulder signature for person recognition,” in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*, pp. 1226–1231, IEEE, 2012.
- [56] KLEINEHAGENBROCK, M., LANG, S., FRITSCH, J., LOMKER, F., FINK, G. A., and SAGERER, G., “Person tracking with a mobile robot based on multi-modal anchoring,” in *Robot and Human Interactive Communication, 2002. Proceedings. 11th IEEE International Workshop on*, pp. 423–429, IEEE, 2002.

- [57] KOREN, Y. and BORENSTEIN, J., “Potential field methods and their inherent limitations for mobile robot navigation,” in *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1398–1404, IEEE, 1991.
- [58] KOWADLO, G., YE, P., and ZUKERMAN, I., “Influence of gestural salience on the interpretation of spoken requests.,” in *INTERSPEECH*, pp. 2034–2037, 2010.
- [59] KRUSE, T., BASILI, P., GLASAUER, S., and KIRSCH, A., “Legible robot navigation in the proximity of moving humans,” in *Advanced Robotics and its Social Impacts (ARSO), 2012 IEEE Workshop on*, pp. 83–88, IEEE, 2012.
- [60] KRUSE, T., KIRSCH, A., SISBOT, E. A., and ALAMI, R., “Exploiting human cooperation in human-centered robot navigation,” in *RO-MAN, 2010 IEEE*, pp. 192–197, IEEE, 2010.
- [61] KUDERER, M., KRETZSCHMAR, H., and BURGARD, W., “Teaching mobile robots to cooperatively navigate in populated environments,” in *Intelligent Robots and Systems (IROS), 2013 IEEE/RSJ International Conference on*, pp. 3138–3143, IEEE, 2013.
- [62] LAM, C.-P., CHOU, C.-T., CHIANG, K.-H., and FU, L.-C., “Human-centered robot navigation towards a harmoniously human–robot coexisting environment,” *IEEE Transactions on Robotics*, vol. 27, no. 1, pp. 99–112, 2011.
- [63] LEE, M. K. and TAKAYAMA, L., “Now, i have a body: Uses and social norms for mobile remote presence in the workplace,” in *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pp. 33–42, ACM, 2011.
- [64] LEIBE, B., SEEMANN, E., and SCHIELE, B., “Pedestrian detection in crowded scenes,” in *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, vol. 1, pp. 878–885, IEEE, 2005.
- [65] LI, Z., HOFEMANN, N., FRITSCH, J., and SAGERER, G., “Hierarchical modeling and recognition of manipulative gesture,” in *Proc. of the Workshop on Modeling People and Human Interaction at the IEEE Int. Conf. on Computer Vision*, vol. 77, 2005.
- [66] LICHTENTHÄLER, C. and KIRSCH, A., “Towards Legible Robot Navigation - How to Increase the Intent Expressiveness of Robot Navigation Behavior,” in *International Conference on Social Robotics - Workshop Embodied Communication of Goals and Intentions*, 2013.
- [67] LOPER, M. M., KOENIG, N. P., CHERNOVA, S. H., JONES, C. V., and JENKINS, O. C., “Mobile human-robot teaming with environmental tolerance,” in *Proceedings of the 4th ACM/IEEE international conference on Human robot interaction*, pp. 157–164, ACM, 2009.

- [68] LU, D. V. and SMART, W. D., “Towards more efficient navigation for robots and humans,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 1707–1713, IEEE, 2013.
- [69] LUBER, M., SPINELLO, L., SILVA, J., and ARRAS, K. O., “Socially-aware robot navigation: A learning approach,” in *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pp. 902–907, IEEE, 2012.
- [70] MARTIN, C., BÖHME, H.-J., and GROSS, H.-M., “Conception and realization of a multi-sensory interactive mobile office guide,” in *Systems, Man and Cybernetics, 2004 IEEE International Conference on*, vol. 6, pp. 5368–5373, IEEE, 2004.
- [71] MARTIN, C., STEEGE, F.-F., and GROSS, H.-M., “Estimation of pointing poses for visually instructing mobile robots under real world conditions,” *Robotics and Autonomous Systems*, vol. 58, no. 2, pp. 174–185, 2010.
- [72] MARTÍNEZ-GARCÍA, E., AKIHISA, O., YUTA, S., and OTHERS, “Crowding and guiding groups of humans by teams of mobile robots,” in *Advanced Robotics and its Social Impacts, 2005. IEEE Workshop on*, pp. 91–96, IEEE, 2005.
- [73] MATIKAINEN, P., PILLAI, P., MUMMERT, L., SUKTHANKAR, R., and HEBERT, M., “Prop-free pointing detection in dynamic cluttered environments,” in *Automatic Face & Gesture Recognition and Workshops (FG 2011), 2011 IEEE International Conference on*, pp. 374–381, IEEE, 2011.
- [74] MATUSZEK, C., BO, L., ZETTLEMOYER, L., and FOX, D., “Learning from unscripted deictic gesture and language for human-robot interactions,” 2014.
- [75] MITZEL, D. and LEIBE, B., “Close-range human detection for head-mounted cameras,” in *British Machine Vision Conference (BMVC)*, 2012.
- [76] MIURA, J., SATAKE, J., CHIBA, M., ISHIKAWA, Y., KITAJIMA, K., and MASUZAWA, H., “Development of a person following robot and its experimental evaluation,” in *Proceedings of the 11th International Conference on Intelligent Autonomous Systems, Ottawa, Canada*, pp. 89–98, 2010.
- [77] MOESLUND, T. B. and GRANUM, E., “A survey of computer vision-based human motion capture,” *Computer Vision and Image Understanding*, vol. 81, no. 3, pp. 231–268, 2001.
- [78] MONTEMERLO, M., THRUN, S., and WHITTAKER, W., “Conditional particle filters for simultaneous mobile robot localization and people-tracking,” in *Robotics and Automation, 2002. Proceedings. ICRA’02. IEEE International Conference on*, vol. 1, pp. 695–701, IEEE, 2002.

- [79] MOZOS, O. M., TRIEBEL, R., JENSFELT, P., ROTTMANN, A., and BURGARD, W., “Supervised semantic labeling of places using information extracted from sensor data,” *Robotics and Autonomous Systems*, vol. 55, no. 5, pp. 391–402, 2007.
- [80] MÜLLER, J., STACHNISS, C., ARRAS, K., and BURGARD, W., “Socially inspired motion planning for mobile robots in populated environments,” in *Proc. of International Conference on Cognitive Systems*, 2008.
- [81] MUNSELL, B. C., TEMLYAKOV, A., QU, C., and WANG, S., “Person identification using full-body motion and anthropometric biometrics from kinect videos,” in *Computer Vision–ECCV 2012. Workshops and Demonstrations*, pp. 91–100, Springer, 2012.
- [82] MURAKAMI, R., MORALES SAIKI, L. Y., SATAKE, S., KANDA, T., and ISHIGURO, H., “Destination unknown: walking side-by-side without knowing the goal,” in *Proceedings of the 2014 ACM/IEEE international conference on Human-robot interaction*, pp. 471–478, ACM, 2014.
- [83] NICKEL, K. and STIEFELHAGEN, R., “Pointing gesture recognition based on 3d-tracking of face, hands and head orientation,” in *Proceedings of the 5th international conference on Multimodal interfaces*, pp. 140–146, ACM, 2003.
- [84] NOURBAKHSH, I. R., KUNZ, C., and WILLEKE, T., “The mobot museum robot installations: A five year experiment,” in *Intelligent Robots and Systems, 2003.(IROS 2003). Proceedings. 2003 IEEE/RSJ International Conference on*, vol. 4, pp. 3636–3641, IEEE, 2003.
- [85] OHKI, T., NAGATANI, K., and YOSHIDA, K., “Collision avoidance method for mobile robot considering motion and personal spaces of evacuees,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 1819–1824, IEEE, 2010.
- [86] OHYA, A. and MUNEKATA, T., “Intelligent escort robot moving together with human-interaction in accompanying behavior,” in *Proceedings 2002 FIRA Robot World Congress*, pp. 31–35, 2002.
- [87] OTA, M., OGITSU, T., HISAHARA, H., TAKEMURA, H., ISHII, Y., and MIZOGUCHI, H., “Recovery function for human following robot losing target,” in *Industrial Electronics Society, IECON 2013-39th Annual Conference of the IEEE*, pp. 4253–4257, IEEE, 2013.
- [88] PACCHIEROTTI, E., CHRISTENSEN, H., JENSFELT, P., and OTHERS, “Design of an office-guide robot for social interaction studies,” in *Intelligent Robots and Systems, 2006 IEEE/RSJ International Conference on*, pp. 4965–4970, IEEE, 2006.

- [89] PACCHIEROTTI, E., CHRISTENSEN, H. I., and JENSFELT, P., “Human-robot embodied interaction in hallway settings: a pilot user study,” in *Robot and Human Interactive Communication, 2005. ROMAN 2005. IEEE International Workshop on*, pp. 164–171, IEEE, 2005.
- [90] PANDEY, A. K. and ALAMI, R., “A step towards a sociable robot guide which monitors and adapts to the person’s activities,” in *Advanced Robotics, 2009. ICAR 2009. International Conference on*, pp. 1–8, IEEE, 2009.
- [91] PARK, J. J. and KUIPERS, B., “Autonomous person pacing and following with model predictive equilibrium point control,” in *Robotics and Automation (ICRA), 2013 IEEE International Conference on*, pp. 1060–1067, IEEE, 2013.
- [92] PELLEGRINI, S., ESS, A., SCHINDLER, K., and VAN GOOL, L., “You’ll never walk alone: Modeling social behavior for multi-target tracking,” in *Computer Vision, 2009 IEEE 12th International Conference on*, pp. 261–268, IEEE, 2009.
- [93] PHILLIPS, P. J., FLYNN, P. J., SCRUGGS, T., BOWYER, K. W., CHANG, J., HOFFMAN, K., MARQUES, J., MIN, J., and WOREK, W., “Overview of the face recognition grand challenge,” in *Computer vision and pattern recognition, 2005. CVPR 2005. IEEE computer society conference on*, vol. 1, pp. 947–954, IEEE, 2005.
- [94] PRASSLER, E., BANK, D., and KLUGE, B., “Motion coordination between a human and a robotic wheelchair,” in *Robot and Human Interactive Communication, 2001. Proceedings. 10th IEEE International Workshop on*, pp. 412–417, IEEE, 2001.
- [95] QUINTERO, C. P., FOMENA, R. T., SHADEMAN, A., WOLLEB, N., DICK, T., and JAGERSAND, M., “Sepo: Selecting by pointing as an intuitive human-robot command interface,” in *IEEE Int. Conference of Robotics and Automation, Karlsruhe, Germany*, 2013.
- [96] RAZA ABIDI, S. S., WILLIAMS, M., and JOHNSTON, B., “Human pointing as a robot directive,” in *Proceedings of the 8th ACM/IEEE international conference on Human-robot interaction*, pp. 67–68, IEEE Press, 2013.
- [97] RIOS-MARTINEZ, J., SPALANZANI, A., and LAUGIER, C., “Understanding human interaction for probabilistic autonomous navigation using risk-rrt approach,” in *Intelligent Robots and Systems (IROS), 2011 IEEE/RSJ International Conference on*, pp. 2014–2019, IEEE, 2011.
- [98] ROBINETTE, P. and HOWARD, A. M., “Incorporating a model of human panic behavior for robotic-based emergency evacuation,” in *RO-MAN, 2011 IEEE*, pp. 47–52, IEEE, 2011.
- [99] RUBNER, Y., TOMASI, C., and GUIBAS, L. J., “A metric for distributions with applications to image databases,” in *Computer Vision, 1998. Sixth International Conference on*, pp. 59–66, IEEE, 1998.

- [100] RUSU, R. B. and COUSINS, S., “3d is here: Point cloud library (pcl),” in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, pp. 1–4, IEEE, 2011.
- [101] SASAKI, T. and HASHIMOTO, H., “Human observation based mobile robot navigation in intelligent space,” in *Intelligent Robots and Systems, 2006 IEEE/RSJ International Conference on*, pp. 1044–1049, IEEE, 2006.
- [102] SATAKE, S., KANDA, T., GLAS, D. F., IMAI, M., ISHIGURO, H., and HAGITA, N., “How to approach humans?-strategies for social robots to initiate interaction,” in *Human-Robot Interaction (HRI), 2009 4th ACM/IEEE International Conference on*, pp. 109–116, IEEE, 2009.
- [103] SCANDOLO, L. and FRAICHARD, T., “An anthropomorphic navigation scheme for dynamic scenarios,” in *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, pp. 809–814, IEEE, 2011.
- [104] SCHMIDT, J., HOFEMANN, N., HAASCH, A., FRITSCH, J., and SAGERER, G., “Interacting with a mobile robot: Evaluating gestural object references,” in *Intelligent Robots and Systems, 2008. IROS 2008. IEEE/RSJ International Conference on*, pp. 3804–3809, IEEE, 2008.
- [105] SCHULZ, D., BURGARD, W., FOX, D., and CREMERS, A. B., “Tracking multiple moving targets with a mobile robot using particle filters and statistical data association,” in *Robotics and Automation, 2001. Proceedings 2001 ICRA. IEEE International Conference on*, vol. 2, pp. 1665–1670, IEEE, 2001.
- [106] SHOTTON, J., FITZGIBBON, A., COOK, M., SHARP, T., FINOCCHIO, M., MOORE, R., KIPMAN, A., and BLAKE, A., “Real-time human pose recognition in parts from single depth images,” in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pp. 1297–1304, IEEE, 2011.
- [107] SHOTTON, J., SHARP, T., KIPMAN, A., FITZGIBBON, A., FINOCCHIO, M., BLAKE, A., COOK, M., and MOORE, R., “Real-time human pose recognition in parts from single depth images,” *Communications of the ACM*, vol. 56, no. 1, pp. 116–124, 2013.
- [108] SIDENBLADH, H., BLACK, M. J., and FLEET, D. J., “Stochastic tracking of 3d human figures using 2d image motion,” in *Computer Vision?ECCV 2000*, pp. 702–718, Springer, 2000.
- [109] SIDENBLADH, H., KRAGIC, D., and CHRISTENSEN, H. I., “A person following behaviour for a mobile robot,” in *Robotics and Automation, 1999. Proceedings. 1999 IEEE International Conference on*, vol. 1, pp. 670–675, IEEE, 1999.
- [110] SIEGWART, R., ARRAS, K. O., BOUABDALLAH, S., BURNIER, D., FROIDEVAUX, G., GREPPIN, X., JENSEN, B., LOROTTE, A., MAYOR, L., MEISSER, M., and OTHERS, “Robox at expo. 02: A large-scale installation of personal robots,” *Robotics and Autonomous Systems*, vol. 42, no. 3, pp. 203–222, 2003.

- [111] SIROVICH, L. and KIRBY, M., “Low-dimensional procedure for the characterization of human faces,” *JOSA A*, vol. 4, no. 3, pp. 519–524, 1987.
- [112] SISBOT, E. A., MARIN-URIAS, L. F., ALAMI, R., and SIMEON, T., “A human aware mobile robot motion planner,” *Robotics, IEEE Transactions on*, vol. 23, no. 5, pp. 874–883, 2007.
- [113] SPINELLO, L., ARRAS, K. O., TRIEBEL, R., and SIEGWART, R., “A layered approach to people detection in 3d range data,” in *AAAI Conf. on Artif. Intell.(AAAI)*, 2010.
- [114] STEIN, P., SPALANZANI, A., SANTOS, V., LAUGIER, C., and OTHERS, “Robot navigation taking advantage of moving agents,” in *IROS Workshop on Assistance and Service robotics in a human environment*, 2012.
- [115] STEIN, P. S., SANTOS, V., SPALANZANI, A., LAUGIER, C., and OTHERS, “Navigating in populated environments by following a leader,” in *Ro-man 2013-International Symposium on Robot and Human Interactive Communication*, 2013.
- [116] SVENSTRUP, M., BAK, T., and ANDERSEN, H. J., “Trajectory planning for robots in dynamic human environments,” in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 4293–4298, IEEE, 2010.
- [117] TAKAYAMA, L., MARDER-EPPSTEIN, E., HARRIS, H., and BEER, J., “Assisted driving of a mobile remote presence system: System design and controlled user evaluation,” in *IEEE Int. Conf. on Robotics and Automation (ICRA)*, 2011.
- [118] THRUN, S., BENNEWITZ, M., BURGARD, W., CREMERS, A. B., DELLAERT, F., FOX, D., HAHNEL, D., ROSENBERG, C., ROY, N., SCHULTE, J., and OTHERS, “Minerva: A second-generation museum tour-guide robot,” in *Robotics and automation, 1999. Proceedings. 1999 IEEE international conference on*, vol. 3, IEEE, 1999.
- [119] TOPP, E. A. and CHRISTENSEN, H. I., “Tracking for following and passing persons,” in *Intelligent Robots and Systems, 2005.(IROS 2005). 2005 IEEE/RSJ International Conference on*, pp. 2321–2327, IEEE, 2005.
- [120] TRAUTMAN, P. and KRAUSE, A., “Unfreezing the robot: Navigation in dense, interacting crowds,” in *Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on*, pp. 797–803, IEEE, 2010.
- [121] TREVOR, A. J. B., GEDIKLI, S., RUSU, R., and CHRISTENSEN, H. I., “Efficient organized point cloud segmentation with connected components,” in *ICRA Workshop on Semantic Perception Mapping and Exploration*, 2013.
- [122] TREVOR, A. J., ROGERS III, J. G., CHRISTENSEN, H., and OTHERS, “Planar surface slam with 3d and 2d sensors,” in *Robotics and Automation (ICRA), 2012 IEEE International Conference on*, pp. 3041–3048, IEEE, 2012.

- [123] TURK, M. A. and PENTLAND, A. P., “Face recognition using eigenfaces,” in *Computer Vision and Pattern Recognition, 1991. Proceedings CVPR’91., IEEE Computer Society Conference on*, pp. 586–591, IEEE, 1991.
- [124] TUZEL, O., PORIKLI, F., and MEER, P., “Human detection via classification on riemannian manifolds,” in *Computer Vision and Pattern Recognition, 2007. CVPR’07. IEEE Conference on*, pp. 1–8, IEEE, 2007.
- [125] VASQUEZ, D., STEIN, P., RIOS-MARTINEZ, J., ESCOBEDO, A., SPALANZANI, A., LAUGIER, C., and OTHERS, “Human aware navigation for assistive robotics,” in *ISER-13th International Symposium on Experimental Robotics-2012*, 2012.
- [126] VIOLA, P. and JONES, M. J., “Robust real-time face detection,” *International journal of computer vision*, vol. 57, no. 2, pp. 137–154, 2004.
- [127] WAINER, J., FEIL-SEIFER, D. J., SHELL, D. A., and MATARIC, M. J., “The role of physical embodiment in human-robot interaction,” in *15th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*, pp. 117–122, IEEE, 2006.
- [128] WILSON, A. D. and BOBICK, A. F., “Parametric hidden markov models for gesture recognition,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 21, no. 9, pp. 884–900, 1999.
- [129] XAVIER, J., PACHECO, M., CASTRO, D., RUANO, A., and NUNES, U., “Fast line, arc/circle and leg detection from laser scan data in a player driver,” in *Robotics and Automation, 2005. ICRA 2005. Proceedings of the 2005 IEEE International Conference on*, pp. 3930–3935, IEEE, 2005.
- [130] YAN, P. and BOWYER, K. W., “Biometric recognition using 3d ear shape,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 29, no. 8, pp. 1297–1308, 2007.
- [131] YUAN, F., HANHEIDE, M., SAGERER, G., and OTHERS, “Spatial context-aware person-following for a domestic robot,” 2008.
- [132] ZENDER, H., JENSFELT, P., and KRUIJFF, G.-J. M., “Human-and situation-aware people following,” in *Robot and Human interactive Communication, 2007. RO-MAN 2007. The 16th IEEE International Symposium on*, pp. 1131–1136, IEEE, 2007.
- [133] ZHANG, C. and ZHANG, Z., “A survey of recent advances in face detection,” tech. rep., Tech. rep., Microsoft Research, 2010.
- [134] ZHAO, W., KRISHNASWAMY, A., CHELLAPPA, R., SWETS, D. L., and WENG, J., “Discriminant analysis of principal components for face recognition,” in *Face Recognition*, pp. 73–85, Springer, 1998.

- [135] ZUKERMAN, I., KOWADLO, G., and YE, P., “Interpreting pointing gestures and spoken requests: a probabilistic, salience-based approach,” in *Proceedings of the 23rd International Conference on Computational Linguistics: Posters*, pp. 1558–1566, Association for Computational Linguistics, 2010.
- [136] ZUKERMAN, I., MANI, A., LI, Z., and JARVI, R., “Speaking and pointing?from simulations to the laborator,” *Knowledge and Reasoning in Practical Dialogue Systems*, p. 58, 2011.