



Principles and Guidelines for Evaluating Social Robot Navigation Algorithms

ANTHONY FRANCIS, Logical Robotics, Simpsonville, South Carolina, USA
CLAUDIA PÉREZ-D'ARPINO, NVIDIA Corp, Santa Clara, California, USA
CHENGSHU LI, Computer Science, Stanford University, Stanford, California, USA
FEI XIA, Google DeepMind, Google Inc, Mountain View, California, USA
ALEXANDRE ALAHI, EPFL, Lausanne, Switzerland
RACHID ALAMI, LAAS-CNRS, University of Toulouse, Toulouse, France
ANIKET BERA, Department of Computer Science, Purdue University, West Lafayette, Indiana, USA
ABHIJAT BISWAS, Robotics Institute, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA
JOYDEEP BISWAS, Computer Science, The University of Texas at Austin, Austin, Texas, USA
ROHAN CHANDRA, University of Virginia, Charlottesville, Virginia, USA
HAO-TIEN LEWIS CHIANG, Google DeepMind, Google Inc, Mountain View, California, USA
MICHAEL EVERETT, Northeastern University—Boston Campus, Boston, Massachusetts, USA
SEHOON HA, Georgia Institute of Technology, Atlanta, Georgia, USA
JUSTIN HART, Computer Science, The University of Texas at Austin, Austin, Texas, USA
JONATHAN P. HOW, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA
HARESH KARNAN, Computer Science, The University of Texas at Austin, Austin, Texas, USA
TSANG-WEI EDWARD LEE, Google DeepMind, Google Inc, Mountain View, California, USA
LUIS J. MANSO, Computer Science, Aston University, Birmingham, United Kingdom of Great Britain and Northern Ireland
REUTH MIRSKY, Computer Science, Bar-Ilan University, Ramat Gan, Israel and Computer Science, Tufts University, Medford, MA, USA
SÖREN PIRK, Computer Science, Kiel University, Kiel, Germany
PHANI TEJA SINGAMANENI, LAAS-CNRS, University of Toulouse, Toulouse, France
PETER STONE, The University of Texas at Austin and Sony AI, Austin, Texas, USA
ADA V. TAYLOR, Robotics Institute, Carnegie Mellon University, Stowe, Vermont, USA
PETER TRAUTMAN, Honda Research Institute USA Inc, Palo Alto, California, USA
NATHAN TSOI and **MARYNEL VÁZQUEZ**, Computer Science, Yale University, New Haven, Connecticut, USA
XUESU XIAO, George Mason University, Fairfax, Virginia, USA
PENG XU, Google DeepMind, Google Inc, Mountain View, California, USA
NAOKI YOKOYAMA, Georgia Institute of Technology, Atlanta, Georgia, USA
ALEXANDER TOSHEV, Apple Inc, Cupertino, California, USA
ROBERTO MARTÍN-MARTÍN, Computer Science, The University of Texas at Austin, Austin, Texas, USA

Anthony Francis and Claudia Pérez-D'Arpino contributed equally to this research.

Alexander Toshev and Roberto Martín-Martín contributed equally as advisors.

Authors' Contact Information: Anthony Francis (corresponding author), Logical Robotics, Simpsonville, South Carolina, United States; e-mail: centaur@logicalrobotics.com; Claudia Pérez-D'Arpino, NVIDIA Corp, Santa Clara, California, United States; e-mail: cdarpino@stanford.edu; Chengshu Li, Computer Science, Stanford University, Stanford, California, United States; e-mail: chengshu@stanford.edu; Fei Xia, Google DeepMind, Google Inc, Mountain View, California, United States; e-mail: xiafei@google.com; Alexandre Alahi, EPFL, Lausanne, Switzerland; e-mail: alexandre.alahi@epfl.ch;

A major challenge to deploying robots widely is navigation in human-populated environments, commonly referred to as *social robot navigation*. While the field of social navigation has advanced tremendously in recent years, the fair evaluation of algorithms that tackle social navigation remains hard because it involves not just robotic agents moving in static environments but also dynamic human agents and their perceptions of the appropriateness of robot behavior. In contrast, clear, repeatable, and accessible benchmarks have accelerated progress in fields like computer vision, natural language processing and traditional robot navigation by enabling researchers to fairly compare algorithms, revealing limitations of existing solutions and illuminating promising new directions. We believe the same approach can benefit social navigation. In this article, we pave the road toward common, widely accessible, and repeatable benchmarking criteria to evaluate social robot navigation. Our contributions include (a) a definition of a socially navigating robot as one that respects the principles of safety, comfort, legibility, politeness, social competency, agent understanding, proactivity, and responsiveness to context, (b) guidelines for the use of metrics, development of scenarios, benchmarks, datasets, and simulators to evaluate social navigation, and (c) a design of a social navigation metrics framework to make it easier to compare results from different simulators, robots, and datasets.

CCS Concepts: • **Computer systems organization** → **Robotics**; • **Human-centered computing** → *HCI design and evaluation methods*; • **Computing methodologies** → *Reinforcement learning*; **Simulation evaluation**;

Additional Key Words and Phrases: social robotics, robot navigation, datasets, benchmarks, simulators

Rachid Alami, LAAS-CNRS, University of Toulouse, Toulouse, France; e-mail: rachid.alami@laas.fr; Aniket Bera, Department of Computer Science, Purdue University, West Lafayette, Indiana, United States; e-mail: bera@umd.edu; Abhijat Biswas, Robotics Institute, Carnegie Mellon University, Pittsburgh, Pennsylvania, United States; e-mail: abhijat-biswas@gmail.com; Joydeep Biswas, Computer Science, The University of Texas at Austin, Austin, Texas, United States; e-mail: joydeepb@cs.utexas.edu; Rohan Chandra, University of Virginia, Charlottesville, Virginia, United States; e-mail: rohanchandra@virginia.edu; Hao-Tien Lewis Chiang, Google DeepMind, Google Inc, Mountain View, California, United States; e-mail: lewispro@google.com; Michael Everett, Northeastern University—Boston Campus, Boston, Massachusetts, United States; e-mail: m.everett@northeastern.edu; Sehoon Ha, Georgia Institute of Technology, Atlanta, Georgia, United States; e-mail: sehoonha@gatech.edu; Justin Hart, Computer Science, The University of Texas at Austin, Austin, Texas, United States; e-mail: hart@cs.utexas.edu; Jonathan P. How, Massachusetts Institute of Technology, Cambridge, Massachusetts, United States; e-mail: jhow@mit.edu; Haresh Karnan, Computer Science, The University of Texas at Austin, Austin, Texas, United States; e-mail: haresh.miriyala@utexas.edu; Tsang-Wei Edward Lee, Google DeepMind, Google Inc, Mountain View, California, United States; e-mail: tsangwei@google.com; Luis J. Manso, Computer Science, Aston University, Birmingham, United Kingdom of Great Britain and Northern Ireland; e-mail: lmanso@aston.ac.uk; Reuth Mirsky, Computer Science, Bar-Ilan University, Ramat Gan, Israel and Computer Science, Tufts University, Medford, MA, USA; e-mail: reuth.mirsky@tufts.edu; Sören Pirk, Computer Science, Kiel University, Kiel, Germany; e-mail: soeren.pirk@gmail.com; Phani Teja Singamaneni, LAAS-CNRS, University of Toulouse, Toulouse, France; e-mail: ptsingaman@laas.fr; Peter Stone, The University of Texas at Austin and Sony AI, Austin, Texas, United States; e-mail: pstone@cs.utexas.edu; Ada V. Taylor, Robotics Institute, Carnegie Mellon University, Stowe, Vermont, United States; e-mail: adat@andrew.cmu.edu; Peter Trautman, Honda Research Institute USA Inc, Palo Alto, California, United States; e-mail: peter.trautman@gmail.com; Nathan Tsoi, Computer Science, Yale University, New Haven, Connecticut, United States; e-mail: nathan.tsoi@yale.edu; Marynel Vázquez, Computer Science, Yale University, New Haven, Connecticut, United States; e-mail: marynel.vazquez@yale.edu; Xuesu Xiao, George Mason University, Fairfax, Virginia, United States; e-mail: xiao@gmu.edu; Peng Xu, Google DeepMind, Google Inc, Mountain View, California, United States; e-mail: pengxu@google.com; Naoki Yokoyama, Georgia Tech, Atlanta, Georgia, United States; e-mail: nyokoyama@gatech.edu; Alexander Toshev, Apple Inc, Cupertino, California, United States; e-mail: toshev@apple.com; Roberto Martín-Martín, Computer Science, The University of Texas at Austin, Austin, Texas, United States; e-mail: robertomm@cs.utexas.edu.



This work is licensed under a Creative Commons Attribution International 4.0 License.

© 2025 Copyright held by the owner/author(s).

ACM 2573-9522/2025/2-ART34

<https://doi.org/10.1145/3700599>

ACM Reference format:

Anthony Francis, Claudia Pérez-D'Arpino, Chengshu Li, Fei Xia, Alexandre Alahi, Rachid Alami, Aniket Bera, Abhijat Biswas, Joydeep Biswas, Rohan Chandra, Hao-Tien Chiang, Michael Everett, Sehoon Ha, Justin Hart, Jonathan How, Hareesh Karnan, Tsang-Wei Lee, Luis Manso, Reuth Mirksy, Sören Pirk, Phani Teja Singamaneni, Peter Stone, Ada Taylor, Peter Trautman, Nathan Tsoi, Marynel Vázquez, Xuesu Xiao, Peng Xu, Naoki Yokoyama, Alexander Toshev, and Roberto Martín-Martín. 2025. Principles and Guidelines for Evaluating Social Robot Navigation Algorithms. *ACM Trans. Hum.-Robot Interact.* 14, 2, Article 34 (February 2025), 65 pages.

<https://doi.org/10.1145/3700599>

1 Introduction

The study of social robot navigation has a long history, but a crisp definition of what makes navigation “social” remains elusive. Researchers on social robot navigation often have a personal sense of what it is and use that intuition to guide their research into how to make robots move better around people, but the field does not yet have a consensus on a definition of social navigation or how to achieve it. Indeed, at the Social Navigation Symposium,¹ a diverse spectrum of researchers presented a variety of views on what robotic social navigation is and their approaches to solving it, including a range of definitions, variants, problems, and sub-problems.

Ideas presented at the Symposium included a variety of methods to evaluate social navigation performance, involving different experimental setups, evaluation metrics, robot simulators, social datasets, and deployment environments. As the researchers continued their discussions following the symposium, a taxonomy of aspects of social navigation began to emerge, which helped clarify the social robot navigation problem and converged to a set of general recommendations on how to evaluate solutions in ways that were more comparable.

This article summarizes our work to define the social robot navigation problem, identify a taxonomy of its important aspects, create guidelines for its evaluation, and define a common API to make evaluations more comparable. After a review of related work in Section 2, Section 3 proposes a definition of social navigation and a strategy for achieving it by following social navigation principles. Section 4 reviews the different scientific questions asked by social navigation researchers, and Section 5 outlines our taxonomy for analyzing social navigation benchmarks, datasets and simulators. Section 6 discusses the metrics that have been developed for measuring social navigation, including subjective human evaluation metrics, computable analytic metrics, and research toward learned metrics. Section 7 discusses the typical scenarios used in social navigation, and Section 8 describes benchmarks built on these scenarios, while Section 9 reviews datasets collected on social navigation. Section 10 reviews simulators and presents our work to create a unified interface across them.

Figures 1 and 2 illustrate the principles and guidelines we present for the development and evaluation of social navigation. Principles are high-level goals that social navigation methods should try to achieve, as illustrated in Figure 2. Guidelines are concrete, actionable recommendations that practitioners of social navigation research may consider when creating and testing their solutions, as summarized in Figure 1 and unpacked in the rest of the document.

2 Related Work

The field of social robot navigation is vast and we will not attempt to summarize it; instead, we refer to many recent surveys on social navigation [23, 28, 31, 50, 79, 97, 103, 105, 133, 167]. Among these,

¹<https://sites.google.com/view/socialnaviationsymposium/home>

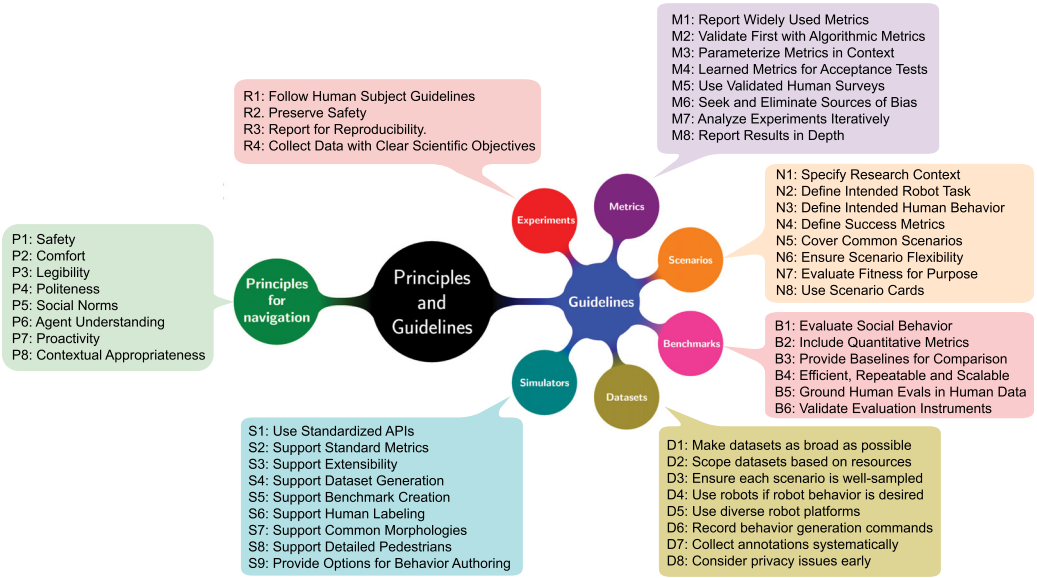


Fig. 1. We identify eight broad principles of social robot navigation—including safety, comfort, legibility, politeness, social competency, agent understanding, proactivity, and contextual appropriateness—which motivate specific guidelines for experiments, metrics, scenarios, benchmarks, datasets, and simulators. Principles and guidelines are labeled with two-letter codes, with P for principles, R for real-world issues, M for metrics, N for scenarios, B for benchmarks, D for datasets, and S for simulators.

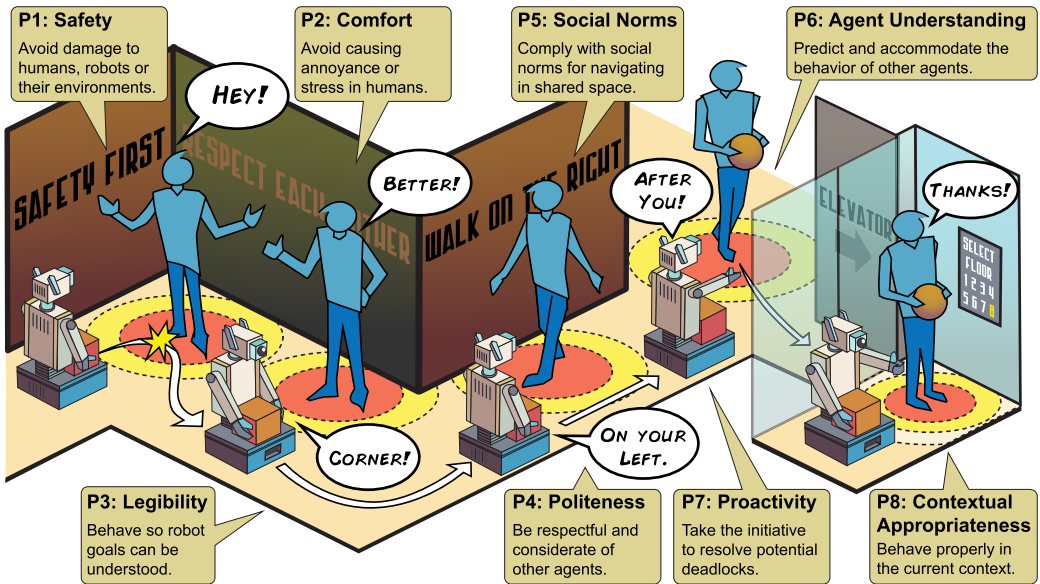


Fig. 2. We define a socially navigating robot as one that interacts with humans and other robots in a way that achieves its navigation goals while enabling other agents to achieve theirs. To make this objective achievable, we propose eight principles for social robot navigation: safety, comfort, legibility, politeness, social competency, agent understanding, proactivity, and contextual appropriateness.

[50] focuses on evaluating social robot navigation algorithms, reviewing 177 recent papers to gather evaluation methods, scenarios, datasets, and metrics, using their findings to discuss shortcomings of existing research and to make recommendations for future research directions. Another recent survey by Mavrogiannis et al. [97] focuses on the core challenges of social navigation with respect to navigation algorithms, human behavior models, and evaluation. Our work builds on the works of [50] and [97] and similar surveys to map the field. We contribute a crisp definition of social robot navigation based on discussions held at the 2022 Social Navigation Symposium, an overview of methodologies for research, and a taxonomy of the field which we use to examine existing metrics, scenarios, benchmarks, datasets, and simulators, and shared principles to make social navigation evaluations comparable across the community.

Wang et al. [167] propose new metrics evaluating the principles defined in [79], comfort, naturalness, and sociability. We expand the principles in [79] and propose a lifecycle of social navigation with recommendations for metrics, scenarios, benchmarks, datasets, and simulators, along with guidelines for metric usage.

Beyond social navigation, clear, repeatable, and accessible benchmarks have accelerated progress in fields like computer vision [139] and natural language processing [15, 130, 166] enabling researchers to compare algorithms, revealing limitations of existing solutions, and illuminating promising new directions. Our effort builds on benchmark challenges in traditional robot navigation [2, 35, 36, 141, 169, 170], social navigation benchmarks and challenges [11, 12, 39, 45, 72, 83, 109, 122, 145], and social navigation scenario development [34, 126, 171]. We review social navigation scenarios and benchmarks in Sections 7 and 8. We contribute guidelines for scenario development, a review of scenarios in the literature, a social navigation scenario card, as well as guidelines for social navigation benchmarking and dataset development.

Simulators are a key component in social navigation, though many simulators exist with diverse APIs which are largely not compatible. [72], discussed in more detail in Section 8 is a benchmark that provides an API for easily generating new worlds and tasks for two different simulators. This article proposes guidelines for simulator development and usage, as well as a common API design to unify simulator outputs to facilitate common evaluations using shared metrics.

3 Toward a Definition of Social Navigation

Social navigation refers to a range of behaviors from simple navigation around dynamic obstacles, to complying with complex social norms, up to navigating with communicative intent. As such, it risks becoming a “suitcase word,” defined by Minsky [101] as words that carry other concepts inside them, like memory, emotions, or consciousness; these terms must be unpacked to understand their meanings fully.

In this section, we unpack the term “social robot navigation.” First, we examine social robotics and what it means. Then we examine the sub-problems of social navigation and how context can affect what behaviors are considered social. To guide research, we formulate these problems as social navigation principles.

3.1 What Is a Social Robot?

Intuitively, we expect social robots to be able to recognize social cues, norms, and expectations, to have the understanding to interpret them correctly, and to have the capabilities to respond appropriately. This raises the question of what “social” is, and what kinds of sensing, interpretation, and capabilities social robots need to effectively navigate social interactions.

In their review of **Human–Robot Interaction (HRI)** for social robotics, Kanda and Ishiguro [68] argue that in addition to navigation (moving robots from place to place) and manipulation (changing objects in the environment) capable robots must also leverage social interactions, i.e., be

able to interplay with humans or other robots to perform tasks. Further, they distinguish robots that simply encounter humans from those that have socially interactive features, such as voices, expressive faces, or the ability to gesture.

But simply having socially interactive features in a robot does not mean that the quality of its interactions would be acceptable to humans or efficient for others; additional principles are needed to apply these features in a positive way. Developing solutions that create high-quality social interactions autonomously is difficult; many social interactions that Kanda and Ishiguro studied were beyond the technology of the time and required a human to teleoperate the robot. What distinguishes “social” robotics from pure interactivity?

To define “social” more precisely, we examined the terms social and antisocial for humans. Social sometimes means participating in society, i.e., participating in an interacting group whose individuals modify their behavior to accommodate the needs of others while achieving their own. But social has a second meaning: a social individual has outstanding skills to work with others, based on an understanding of their feelings and needs and adapting to them. Antisocial individuals, in contrast, fail to follow the customs of society or live without consideration for others. Inspired by these terms when applied to humans, we generalize this notion to other agents, and offer this definition of social robot navigation:

A socially navigating robot acts and interacts with humans or other robots, achieving its navigation goals while modifying its behavior so the experience of agents around the robot is not degraded or is even enhanced.

This social quality may be reflected through overt behavior changes, such as respecting social norms, or through understanding other agents’ needs, feelings, and capabilities.

3.2 Principles of Social Navigation

It is often difficult for an agent to know exactly what other agents, especially humans, need to achieve, or what they feel and like, and social norms that could guide us are often not verbalized. To operationalize these concerns, we identified *principles of social navigation* that can be used to evaluate the quality of social behavior, including (1) *safety*, (2) *comfort*, (3) *legibility*, (4) *politeness*, (5) *social competency*, (6) *understanding other agents*, (7) *proactivity*, and (8) *responding appropriately to context*, as illustrated in Figure 2.

Seen from the lens of optimization, the first seven principles of social navigation can be formalized as additional objectives that the robot needs to optimize for while still achieving its main objective, and the eighth principle, context, can be seen as weighting which principles are most important at any given time, as shown in Figure 3. These principles are not completely orthogonal: improving legibility might improve safety and even comfort, whereas non-verbal politeness depends on understanding other agents’ trajectories. In addition, what is considered appropriate or polite behavior depends on both the cultural context [132] and the robot’s main objective: for example, a delivery robot arguably should maintain a greater distance from humans than one functioning as a guide.

The principles mentioned above guide the development of metrics to evaluate social robots, discussed in more depth in Section 6. Properly studying these principles of social navigation directly impacts which metrics to measure [50], what datasets to collect, how to build simulators, and how to structure benchmarks. In the following sections, we unpack these principles as often used in social robotics research:

Principle P1: Safety—Protect Humans, Other Robots, and Their Environments. A minimal requirement for robots and human sociality is not harming others in the course of business [14, 80], as the

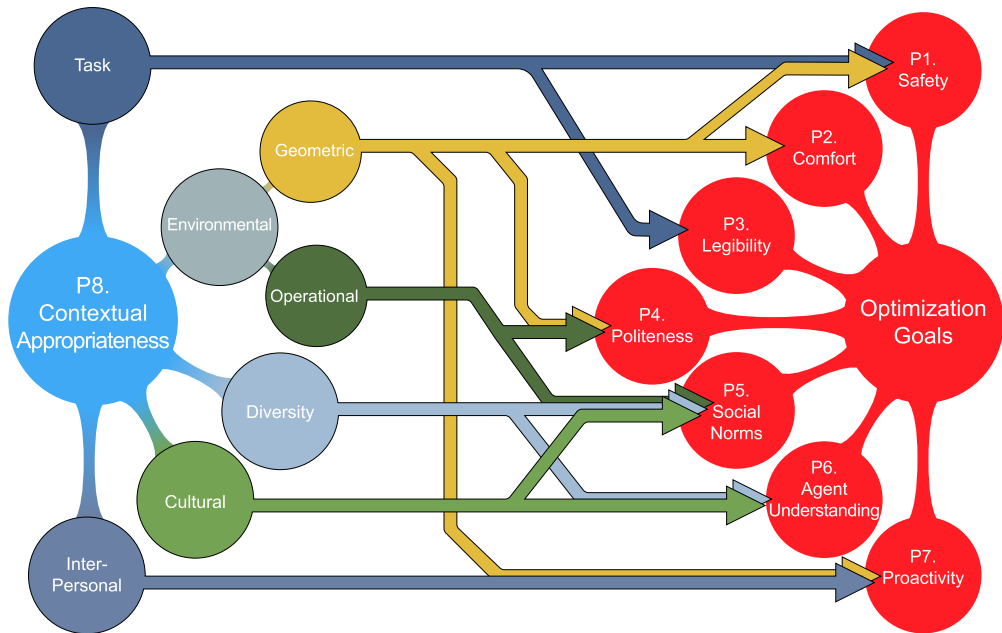


Fig. 3. Contextual factors of social navigation. While the first seven principles represent factors to optimize, the eighth principle, contextual appropriateness, calls out that the weighting of these factors can be affected by many features, including cultural, diversity, environmental, task, and interpersonal context. Lines in the diagram are representative of common interactions but are not exclusive.

robot fails to do in the first scenario in Figure 2 when it collides with a human’s toe. Avoiding collisions with humans is important but is not the only safety concern [9, 97, 110]; robots can damage each other or their environments. While it might be acceptable for a factory robot to bump a guardrail defining the edge of its workspace, social robots should generally avoid damaging human environments, which often contain important objects that can be damaged or wall coverings whose visual appearance is important. Robots should also avoid damaging each other or behaving in a way that induces humans or other robots to injure themselves.

Principle P2: Comfort—Do Not Create Annoyance or Stress. Humans should also feel comfortable around robots, defined in [79] as *the absence of annoyance and stress for humans in interaction with robots*. Many features contribute to comfort, including maintaining human–robot distance, not cutting humans off, and naturalness of motion. Unacceptable robot speed, navigation jitter, and unexpected head movements are factors that degrade humans’ perception of comfort. Additionally, social robots should arguably not behave in a way that triggers the safety layers of other robots. Kruse et al. [79] further argue that annoyance can be triggered by a failure to respect proxemics, the virtual personal space around a human that other humans instinctively respect [57]. Figure 2 illustrates proxemics with the “intimate” distance of 0.45 m shown in red and the “personal” distance of 1.2 m; after initially violating a human’s personal distance, the robot is shown attempting to stay in “personal” or more distant “social” spaces, except as required by the geometric context. Proxemics is a rich and controversial field; for an in-depth survey see [133].

Principle P3: Legibility—Behave so Goals Can Be Understood. Legibility refers to the property of an agent’s behavior that makes it possible for other agents to infer their goals [37]. This includes not

only the robot's goal but also incidental interactions when performing other tasks, e.g., moving to the right or left when passing in a hallway. Dragan et al. [37] define legibility as relaxing constraints such as predictability of trajectories (in the sense of an agent's own predictable style) in favor of more clearly understood behaviors (in the sense of moving to make goals explicit). While [37] focused on changes to robot paths to make them legible, a robot capable of communicating could explicitly announce its intentions, the way restaurant staff are trained to call "corner" when entering a blind corner, as the robot does in the middle of Figure 2.

Principle P4: Politeness—Be Respectful and Considerate. Politeness refers to behavior that is respectful and considerate of people. There are at least two dimensions: physical politeness (how robots navigate around people, such as not cutting people off) and communicative politeness (gestures or verbal signals, such as saying "excuse me," or "on your left," as the robot does in Figure 2 when a narrow hallway forces it to transgress on a human's personal space). Politeness can have a strong effect on people's perception of robots [66, 131]. Social robots should also be considerate of each other, so they do not prevent other robots from accomplishing their tasks.

Principle P5: Social Competency—Comply with Social Norms. Robots should comply with social, political, and legal norms for sharing space. Many social competencies are matters of following conventions rather than optimizing performance [26, 29, 131]. For example, in the absence of norms there is no optimization preference for driving on the left or right, but identifying and following the local norms helps prevent conflicts [103] in the third hallway interaction in Figure 2. Some social competencies, like turn-taking, can emerge naturally [77], whereas others must be specifically engineered. Social norms may apply to more than just humans; conventions of behaviors may make it easier for robots to interact.

Principle P6: Understanding Other Agents—Predict and Accommodate the Behavior of Other Agents. Understanding, accommodating, and even facilitating other agents' activities is a key element of social behavior. Accommodating other agent's goals and comfort requires an understanding of what they are perceiving, doing, and trying to accomplish. For example, to pass between two agents politely, it is important to understand whether they are conversing [125]. Understanding when the *interaction potential*—the likelihood of robots entering human personal space—should be minimized [8, 156] or maximized [99] depends on the task [102]. Further, understanding how agents move can reduce the potential for *conflicts* (short-term encounters in which humans and robots would collide without intervention [20, 103]), as in the right side of Figure 2, where a robot recognizes the human's path will cross theirs and stops to prevent a conflict.

Principle P7: Proactivity—Taking the Initiative to Prevent and Resolve Issues. Simply understanding other agents is not enough, however: in some circumstances, being social involves taking the initiative or even interrupting other agents in their navigational task [93]. For instance, at a four-way intersection, delays occur if all drivers act conservatively, necessitating one to take the lead or propose a solution, such as signaling others to proceed [21, 151]. Research on self-driving vehicles shows that non-conservative (or "aggressive") behavior can be effective or even desirable when expected by others [19, 22]. Although robots are less likely to be mistaken for humans as in the self-driving domain, similar deadlock scenarios, like two pedestrians dodging in the same direction, can arise. In such cases, a robot that takes the initiative to avoid a human or to proactively suggest a solution is arguably more socially adept, as shown on the right side of Figure 2, where the robot proactively proposes that the human enter the elevator first to prevent this kind of deadlock. Section 7 discusses measuring proactivity through scenarios designed to elicit this behavior.

Principle P8: Contextual Appropriateness—Behave Properly in the Current Context. Social navigation should be evaluated within the context that it is to be deployed. Context helps us understand the relative importance of the previous objectives and is a complex construct in its own right, as shown in Figure 3. An example shared in the symposium was a CRASH CART robot in a hospital bringing an emergency drug to a doctor: politeness is less important than task success. Also, when navigating a narrow corridor, we may be “less polite” and get closer to other agents. We identified the following forms of context, all of which can change which response is right in a given situation:

- *Cultural Context:* Different cultures have different social norms, as notably documented in [57]; more recently [7, 9] and [132] examined cultural norms in social robotics but concluded more work needs to be done.
- *Diversity Context:* Different individuals with different abilities or different background histories may need different accommodations [132].
- *Environmental Context:* The environment may affect the social navigation problem [112] and includes both geometric factors—the shape of the space—and operational factors—how that shape is to be used.
 - *Geometric Context:* The geometry of the environment may affect the social navigation problem. For example, the more crowded the space is, the smaller the acceptable distance is between the robot and other agents.
 - *Operational Context:* The operational domain the robot is intended to work in affects what behaviors are considered good: for example, a robot may drive slower in a daycare than in an office, even if the two settings had identical geometric layouts.
- *Task Context:* In turn, the task the robot operates in affects what behaviors are appropriate: for example, even in a single environment like a hospital, whether a robot is performing mail delivery or is a crash cart changes its weighting of politeness against speed.
- *Interpersonal Context:* While there are many different areas of context that are appropriate, interpersonal context (e.g., whether humans are independent pedestrians, are traveling in a group, or stopped and conversing is critical to knowing how to navigate among them).

As an illustration of context, the robot in Figure 2 is first shown violating a person’s intimate space distance in red, then attempting to avoid proxemics violations going forward. However, the corridor around the bend is too narrow to prevent the robot from passing through a person’s personal space distance in yellow, prompting the robot to politely call out its presence. Then, in the relatively small elevator, the standard interpersonal distances are no longer easy to achieve, and both the robot and human adjust their perceived proxemics radii based on the current context, shown as a contraction of the proxemics circles from their original size.

Social navigation should be evaluated within the context that it is intended to be deployed. While defining the context in a sufficiently precise way for a robot to identify or respond to it is a challenging problem, at the least, the intended context should be defined well enough in terms of cultural, diversity, environmental, operational, task, and interpersonal context for other researchers to gauge the applicability of the ideas and findings conveyed by research.

4 Research Methodologies of Social Navigation

Benchmarks require measures and an evaluation methodology for comparing social robot navigation systems. They consider different phenomena including human perceptions of robot behavior as well as objective properties, such as behavior around dynamic obstacles, which can affect social principles like safety and comfort. The scientific questions that benchmarks ask determine the phenomena they study and the data they collect, which in turn guide the development of social navigation methods, creating a lifecycle of social navigation research.

4.1 Research Questions of Social Navigation

The overarching research question of social robot navigation is developing a scientific understanding of the problem sufficient to build computational models that enable robots to perform acceptably in human environments. This involves understanding the factors that influence social navigation, developing models of those factors, and implementing algorithms that take them into account. To fairly evaluate how algorithms compare to each other, we need standardized benchmarks that help us understand their differences and identify which ones are better as opposed to evaluations crafted for each algorithm. Given the complexity of social navigation, different benchmarks often focus on different aspects of the problem and thus different, more specific scientific questions. Some of these questions arise from traditional robot navigation research and can arguably be evaluated using traditional methods, with adjustments for human participants:

- *How do methods compare with each other against baselines?* Some aspects of method evaluation involve quantitative metrics measurable in simulation, such as revealing problems in a robot’s safety layer as it faces increasing obstacle densities. However, when human evaluations are required, these are typically conducted in the real world, though toolkits are now coming into use that enable labeling simulated trajectories as well [6].
- *How do components of a method affect its overall performance?* These are generally conducted by turning method components off, often called “ablation studies” in analogy to ablation studies in neuroscience [100]. While in theory ablation studies could be conducted on-robot, in practice these studies are often only conducted in simulation, as real human participant time can be wasted on variants of the algorithm expected to perform poorly (or known to perform poorly in simulation).
- *How do behaviors generalize to different environments?* Benchmarks can test methods under different conditions to evaluate this, a task that is easier (though less realistic) in simulation.

Success at task performance is often measurable and quantitative, but determining whether a robot is satisfying the principles of social navigation is trickier. While the physical aspects of Principle P1, Safety, may arguably be measurable quantitatively (at least in the sense that the lack of safety can be measured through damage and collisions), others like Principle P4, Politeness, often require human evaluation, and still, others like Principle P2, Comfort, are often explicitly defined in terms of human reactions. Scientific questions involving these subjective aspects therefore generally require measuring human perceptions and reactions to robot behaviors and are best investigated through HRI studies:

- *Human Ratings:* How do humans rate the socialness of social navigation methods, either intrinsically or in comparison to a baseline? For some researchers, human ratings of policy behavior in real contexts are the gold standard for policy performance, but for these ratings to be effective, studies must follow proper HRI protocols and use validated survey instruments [64, 121].
- *Behavior Analysis:* How does human behavior change when exposed to different robot navigation policies? While ratings are explicit, behavior change is implicit or even unconscious. Studies should be conducted according to HRI guidelines to ensure conditions are appropriately blinded so participant and rater reactions are valid.
- *Issue Discovery:* Benchmarks can also be used to conduct exploratory analyses. For example, these analyses could find out the frequency of encounter types between humans and robots as well as the frequency of problems that affect a given policy. This can guide research in the direction of problems that occur in the wild. These studies must be conducted with a robot in a live deployment.

Many benchmarks focus on a subset of these questions because different researchers have different aims and different groups have different needs. As a result, social navigation evaluation methodologies have become fragmented and a comprehensive evaluation methodology does not yet exist.

Because different lines of research have different needs, we do not aim to provide one evaluation protocol for all social navigation methods, but a methodology by which researchers can make principled decisions to guide their own evaluations. Such a protocol will allow researchers to compare social navigation methods along the dimension relevant to a specific subdomain.

In Section 8 we argue that because social navigation involves understanding both how robots affect other agents and which methods are effective, most benchmarks will benefit from incorporating both HRI components that evaluate human reactions in the real world as well as ablation studies, even if those are constrained to simulation.

4.2 Types of Social Navigation Studies

We advocate viewing in-the-wild studies and controlled scenarios as part of a lifecycle of study of social navigation phenomena. To define terms, we can distinguish several different major classes of social navigation studies:

- (1) *Field Studies*: Field studies involve pedestrians who are not confederates of the experimenters, such as a mall, campus, or boardwalk. Such studies are often called “in the wild” as they are conducted in uncontrolled environments. Field studies provide an opportunity to collect natural data about robot-human interactions outside the influence of experiments or instructions, but individual encounters are not directly reproducible. Very-large-scale studies offer a proxy of reproducibility when rare events re-occur with enough statistical frequency to be analyzed; however, large-scale field studies are the most resource-intensive, complex, and potentially dangerous to conduct.
- (2) *Robot Deployments*: Robot deployments are conducted in environments partially under the control of the experimenters, such as an office, a classroom building, or a factory. In this case, robot deployments necessarily involve experimenters informing participants about the robots, which may change their responses compared to someone encountering a robot in the wild; furthermore, participants necessarily develop experiences about the robots that can distort HRIs. Symposium participants reported that users unfamiliar with robots were less accepting of errors than robot researchers, who in turn were less accepting than experienced “robot wranglers” responsible for managing the deployment; these anecdotal reports mirror studies that found evidence that both general computer user skill level [69] and familiarity with particular robots [75, 115, 140] could improve assessments of robot capabilities and behavior. Semi-controlled robot deployments are similar to, but less naturalistic than field studies, but because robot deployment environments are more controlled than true in-the-wild studies, a larger scale is often more practicable by conducting experiments over a longer period of time. For example, [13] collected 1,000 kilometers of indoor navigation on a college campus, and the system described in [171] was part of a deployment at Google that collected over 3,000 kilometers of data.
- (3) *Laboratory Experiments*: Laboratory experiments are sometimes considered the gold standard in science but may have distorting effects on human behavior due to the controlled environment and experiment instructions. While A/B testing in field studies or robot deployments can compare some algorithms, laboratory experiments are often necessary to answer scientific questions about human reactions to changes in robot behavior or to evaluate algorithmic changes prior to large-scale deployments. However, we also need to ensure that laboratory

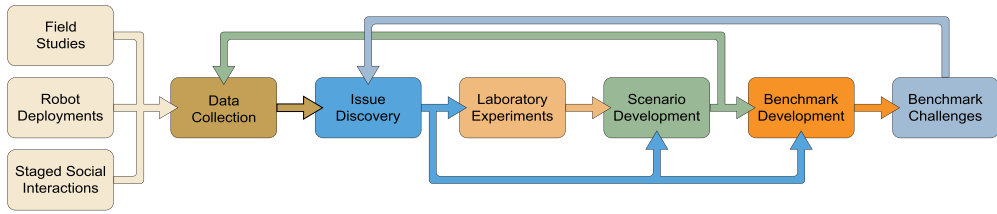


Fig. 4. Lifecycle of social navigation research. Field studies, robot deployments, and staged social interactions can be used to collect data, which helps identify issues and their prevalence. Issues discovered guide laboratory experiments and the development of social navigation scenarios, which in turn can inform data collection. Issue discovery also helps guide the development of benchmarks that test these issues, along with public benchmark challenges; attempts at solutions of these challenges can also help identify issues.

experiments have good *ecological validity*, defined as the degree to which laboratory results generalize to the real world [74, 114, 142]. For social navigation experiments, the ecological validity of an experiment in turn depends on whether the scenario it tests has been properly validated. We discuss a methodology for scenario design in Section 7.1, but validating scientific instruments to determine whether they correctly evaluate the variables they are designed to test, often called *construct validity*, can take several iterations of experiment and analysis [88, 106].

- (4) *Social Navigation Scenarios*: Social navigation scenarios, such as FRONTAL APPROACH, PEDESTRIAN OVERTAKING, and INTERSECTION, can be viewed as a subset of laboratory experiments that test specific scenarios discovered through field studies or robot deployments, with well-defined configurations validated through theoretical analysis, pilot studies, or social navigation issue discovery in existing datasets. The social navigation community is collecting a growing set of scenarios to guide experiments, enable data collection for imitation learning, and serve as regression tests for behavior.
- (5) *Staged Social Interactions*: Due to the excessive costs of field studies and the lack of rare, naturally occurring human–robot encounters in robot deployments and laboratory experiments, researchers developed staged social interactions to evaluate robot social navigation. In staged social interactions, participants are recruited to act in a structured but free-form fashion; this can be an explicit set of scripts (so-called “Guided Crowd Scenarios”) or a less structured activity such as a “Robot Happy Hour” where participants are recruited to perform a social activity around where robots are operating. These studies are less controllable than social navigation scenarios, and their “staged” nature makes them closer to robot deployments or laboratory experiments rather than true field studies. However, they can create higher-density free-form interactions than may otherwise be available.

4.3 Lifecycle of Social Navigation Research

Arguably, social navigation research should be driven by data collected from field studies or robot deployments, but these can be prohibitively expensive; conversely, validated social navigation scenarios enable analysis of known problems, but may not cover novel experimental conditions or detect problems that show up in the wild. Rather than focus on one or the other, it is more useful to think of the following lifecycle of social navigation benchmarking:

- (1) *Data Collection*: Field studies, robot deployments, and staged interactions can be used for the first step of the scientific process: data collection. Ideally, these should be used for more than just A/B testing; they should be used to generate datasets that can be shared to extend the power of the social navigation research community to collect data at scale. Data that can be

collected includes but is not limited to robot and human behavior, surveys (e.g., subject's opinions on safety or comfort), or biometric data (e.g., heart rate, skin impedance).

- (2) *Issue Discovery*: The foundation of social robot navigation is humans interacting with robots. Issue discovery refers to mining human–robot encounter datasets for repeating problematic scenarios that can be reliably detected, enabling the statistical analysis of their frequency and properties. Ideally, the focus should be on high-frequency issues (challenging scenarios that often occur, like a frontal approach in a narrow hallway or the freezing robot problem) and high-risk issues (challenging scenarios where there is a high risk, like compromising a person's safety). Robot deployments in desired target environments are often the best way to collect these data, but large-scale field studies can serve as a proxy.
- (3) *Laboratory Experiments*: Many scientific questions about HRI can be conducted even if large-scale datasets or issue statistics do not exist. Research groups not able to conduct large-scale studies or deployments can nevertheless formulate scientific questions and answer them. Where feasible, these experiments should use benchmarking procedures and metrics validated by the research community, such as those discussed in Section 6. Ideally, these should use scenarios identified as frequent issues in the target domain.
- (4) *Scenario Development*: One outcome of data collection, issue discovery, and laboratory experiments should be the identification of social navigation scenarios that can be reliably detected in datasets, occur frequently in target environments, and can be replicated in controlled laboratory settings. While social navigation scenarios are not a substitute for in-the-wild data collection, using validated social navigation scenarios in laboratory experiments can ensure that experiments are ecologically valid and can ensure that A/B tests are backed up by regression tests of known social navigation issues. Scenarios also aid targeted data collection for both analysis of human behavior and generation of datasets for imitation learning.
- (5) *Social Benchmarking*: Social navigation scenarios can be composed to create benchmarks for social navigation. Most social navigation benchmarks consist, at least implicitly, of a set of social navigation scenarios, real or simulated, that are used to test robot social navigation behavior, along with metrics to gauge performance; many also define datasets of social navigation behavior for comparisons and may also provide simulation environments where the scenarios are defined. From a lifecycle perspective, using reliable, validated scenarios frequently occurring in target environments would make a social navigation benchmark more valuable.
- (6) *Benchmark Challenges*: Finally, benchmarks can be publicly released as “challenges” which include success criteria, a call for solutions, methods for collecting and evaluating solutions, and a leaderboard of ranked solutions. Benchmark challenges have been used for a wide variety of embodied AI tasks and have proved useful for promoting improvements in the field, sometimes leading to challenges being solved and retired (see discussion in [35]). The iGibson Challenge [164] was one of the first publicly available social navigation challenges.

4.4 Guidelines for Real-World Studies

Real-world social navigation studies have aspects that do not come up in simulated experiments or even traditional navigation experiments. Robots controlled by untested policies can damage themselves, other robots, humans, or their environments; human participants captured by robot sensors have privacy and consent concerns. Here we present guidelines for conducting social navigation experiments in the real:

Guideline R1: Follow Guidelines for Human Subject Research. User studies for behavioral and social research should follow guidelines involving respecting the privacy, safety, and well-being of human participants, as well as informed consent. Researchers should follow the guidelines provided by their institution, including study protocol approval for those with an institutional review board or

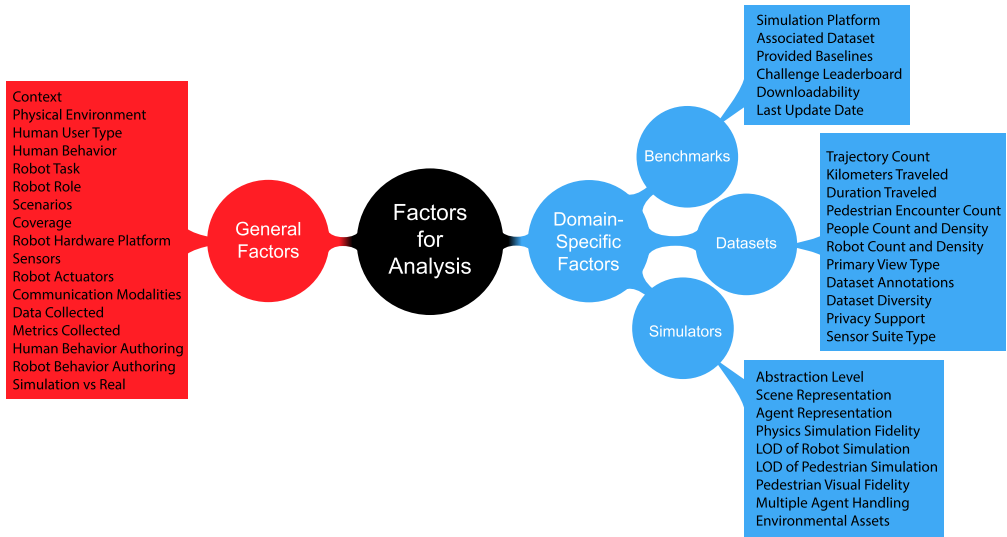


Fig. 5. A taxonomy of social navigation. Most social navigation instruments share common factors like overall context, physical environments, human user type, robot role and task, and so on. However, datasets, benchmarks, and simulators have additional factors particular to them.

independent ethics committees in universities and research institutions, or formal internal ethics review processes in industrial settings.

Guideline R2: Preserve Safety. Real-world benchmarks should preserve the safety of humans, robots, and the environment through active measures such as experiment monitors and safety layers. In particular, policies that have been ablated to illuminate sources of power may have unintuitive behavior; if a safety layer is not available to prevent unsafe actions, either these policies should be tested in simulation or an experiment monitor should be ready to stop the robot in case of issues.

Guideline R3: Report for Reproducibility. As the context of a social navigation experiment can strongly affect its outcome, it is important to report the experiment set up, data collected and metrics analysis clearly and comprehensively enough for other researchers to reproduce the study.

Guideline R4: Collect Data with Clear Scientific Objectives. Real-world experiments are expensive and expose humans, robots, and their environments to risk and should be justified with a clear notion of what is to be learned from conducting the experiment. Furthermore, the default data collected by navigating robots may lack information needed to answer scientific questions about their sociability. Therefore, the purpose of real-world benchmark studies should be clearly defined to ensure enough data are collected to make the experiment scientifically useful.

Clarifying the scientific objectives of a research program can help guide researchers wanting to set up a social navigation study. If the social navigation phenomenon is not yet well understood, a field study might be the most appropriate. Developing datasets for social navigation can be done with robot deployments, which are also useful for experimenting with social robot navigation methods prior to full laboratory studies. When testing algorithms prior to deployment or comparing algorithms to each other, clearly enumerating the important social navigation scenarios and crafting benchmarks or staging social interactions where these scenarios occur frequently can help ensure the comparisons have statistical power.

5 A Taxonomy of Social Navigation

Creating principles and guidelines that are broadly useful to the social navigation community requires understanding the research efforts already underway. Therefore, we have developed a taxonomy for social navigation research (Figure 5) in terms of the metrics, datasets, simulators, and benchmarks used, analyzed with a common vocabulary for factors of analysis.

5.1 A Taxonomy for Analysis

We propose that social navigation research instruments can be analyzed along a set of formal axes which include the metrics they collect, the datasets that they use, if any, the simulator platforms they use, if any, and any formalized scenarios or benchmarks they use for comparison.

- *Metrics*: Recent surveys of social navigation metrics have uncovered close to a hundred different metrics in use (see [50, 103] for recent reviews). Some metrics are algorithmically computed, while others are gathered by surveying humans, either explicitly via questionnaires or implicitly through sensors measuring affect. Algorithmic metrics in turn can be hand-crafted or learned from data gathered from surveys. Other axes of metrics include the type of variable(s) being modeled, and whether metrics cover the behavior of the robot at a specific point in time (step-wise) or during a whole navigation task (episode-wise). The nuances of metrics are discussed in depth Section 6.
- *Scenarios*: Social navigation studies include field studies of behavior in the wild, long-term robot deployments at particular sites, controlled laboratory experiments, social navigation scenarios that aim to create a particular in-the-wild behavior, and “staged” scenarios that attempt to recreate the chaos of crowd scenarios. We have developed a “scenario card” which enables us to compare scenarios, discussed in further depth in Section 7.
- *Benchmarks*: Social navigation benchmarks involve an evaluation protocol for collecting metrics for social robot navigation methods in social navigation scenarios. Current benchmarks are discussed in further depth in Section 8. “Challenges” are benchmarks that are publicly available, include success criteria, and provide evaluation mechanisms along with leaderboards to compare solutions; challenges have shown success in other fields in promoting the improvement of the state of the art [35].
- *Datasets*: We have used these factors to analyze social navigation datasets, discussed in further depth in Section 9. Note that datasets require additional parameters for analysis such as coverage, sampling distribution, annotations, and privacy and fairness handling.
- *Simulators*: Social navigation simulators enable the evaluation of policies controlling agents around other agents in simulation, discussed in further depth in Section 10. Note most simulators have different APIs and metrics.

We next unpack factors common to metrics, scenarios, benchmarks, datasets, and simulators before drilling into these topics in more detail in Sections 6, 7, 8, 9, and 10.

5.2 Factors Common to Social Navigation

Benchmarks, datasets, and simulators for social navigation all face similar challenges: characterizing contexts, representing environments, defining robot roles, tasks, and embodiments, and so on. Rather than analyzing benchmarks, datasets, and simulators separately, we argue that many factors are shared among all three, and here present common factors in attempt to create a common vocabulary for analysis.

- (1) *Context*: Broadly speaking, the context of a social navigation endeavor could refer to its scope, objective, and intended application. As discussed in Section 3.2, context is a complex

construct, and symposium participants did not come to an agreement on a crisp definition, often preferring to use more specific terms when available. However, when it is used, the context of a social navigation research tool often refers to factors implicit in its definition, e.g., a “pedestrian outdoor dataset” or a “benchmark for indoor environments.” Often, the generic concept of “context” can be unpacked into more specific statements about the expected environment, human behavior, or robot tasks; simulators may have these aspects of their context embedded into their design. Aspects of context include the scope of a dataset, what a benchmark tests and what it doesn’t, and what the focus of experiments are: perception, trajectory forecasting, collision avoidance, algorithm benchmarking, human simulation testing, gesture and gaze interaction, body language and affect sensing, human–robot collaboration, indoors vs. outdoors, and individuals vs. crowds.

—Synonyms: Application, Scope

—Related Factors: Robot Role

- (2) *Physical Environment*: Although it could be considered part of the “Context,” the physical space in which the robot(s) and humans operate is particularly relevant. The description of the physical environment includes high-level descriptions such as indoor or outdoor and can be as detailed as one desires. For example, “nearby a water cooler in an office space crowded with cubicles.” Simulator environment definitions may be scanned from the real or authored. Environment definitions also include constraints such as the layouts and traversability of areas of the scene, as robot objectives and constraints are conditioned on the scene layout. The representation of this may be an explicit scene or map or may be implicit in the physical layout of the experiment.

—Synonyms: Location, Scene Type, Context

—Related Factors: Environmental Constraints

- (3) *Type of Human User*: Specifying who the intended or expected human users are is also important. The key is gauging whether the humans are familiar with robots. Humans behave quite differently when they see the robot the first few times, then they get used to it. This type of behavior shift should be noted in a benchmark or dataset since benchmark results obtained by interacting with a group of roboticists may not be representative of when the robot interacts with the public.

—Synonyms: Human Role

—Related Factors: Human Behavior, Robot Role

- (4) *Human Behavior*: A description of the actions taken by specific humans or groups of humans as they relate to the robot. In benchmarks, the desired agent behavior needs to be specified. In simulation, this means the algorithms and scripts that guide the movements of simulated pedestrians. In the real world, this means the instructions to human participants; these could range from a scripted setting, where humans are instructed to perform a specific task or attempt to navigate to a specific location, to unscripted scenarios, where the humans are not explicitly instructed on how to move. Examples of behavior descriptions include humans navigating to specific waypoints, humans blocking the robot or passing. These range from in-the-wild behaviors to carefully specified tasks and everything in between. Simulated human behavior is currently far more constrained than behaviors in the wild.

—Synonyms: Pedestrian Behaviors, Human Tasks

—Related Factors: Robot Task, Robot Role

- (5) *Robot Task*: The piece of work assigned to the robot. The typical robot task is navigation from the robot’s current location to a goal location. Further, higher-level tasks could be specified, such as the delivery of an object, or guarding of an area in the physical environment.

- Synonyms: Robot Behaviors
- Related Factors: Robot Role, Human Behavior
- (6) *Robot Role*: The relationship intended between the robot and the humans, e.g., servant, companion, or fellow pedestrian in a space.
 - Related Factors: Robot Task, Type of Human User
- (7) *Scenarios*: A specific configuration of physical environment, human behavior, and robot task. Scenarios combine three other factors into a package to enable specific configurations of environment, behavior, and task of research interest to be shared in the community. A scenario can be as detailed as a scripted interaction, although free-form scenarios, which are unscripted, are also possible. A robot's role may be specified as part of a scenario, or it may be a variable that is changed and tested.
 - Scenario Classifiers* and *Behavior Graphs* are methods to automatically extract scenarios from data and/or to provide an unambiguous way of labeling
- (8) *Coverage*: The breadth and frequencies of scenarios are also important. Datasets, benchmarks, and simulators can focus on narrowly specified scenarios, a suite of scenarios, or a broad range of cases. Even if the coverage is broad, the distribution of the tests is important, as is explicit coverage of *corner cases*, such as tests or data collection that include erratic, non-cooperative humans.
 - Synonyms: Edge Cases, Regression Tests
- (9) *Robot Hardware Platform*: The specific robot morphology, including its shape, sensors, effectors, displays, and communications modalities. Robot hardware platforms can be instantiated in the real world, in simulation, or both; while many robots have associated simulators, not every robot is represented in every simulator. Unifying robot embodiments is unnecessary and likely impossible, as different robot embodiments are used in different contexts. For this reason, while some benchmarks specify robot embodiments, others are embodiment agnostic.
 - Synonyms: Form Factor, Platform, Embodiment
- (10) *Sensors*: The devices that detect or measure physical properties and record, indicate, or otherwise respond to them. Sensing can include on-board sensors only or may include external sensors or trackers.
 - Synonyms: Inputs, Observation Space
 - Major Divisions: *Robot Sensors* on the robot and *Third-Person Sensors* in the environment.
- (11) *Robot Actuators*: What is the action space of the robot? Conceivably, this may also include third-party actuators such as automatic doors, but this usage is rare.
 - Synonyms: Effectors, Action Space
- (12) *Communications Modalities*: How can humans and robots communicate? Not at all, the robot speaking but not hearing, the robot hearing but not speaking, or two-way? For example, possible communication modalities include visual and audio signals, body and head motion, or no communication at all.
- (13) *Data Collected*: In addition to any robot sensation, actuation, and communication, benchmarks, datasets, or simulators may collect other data such as people tracks, maps of the spaces, and so on. This can include information about pedestrians, such as access to explicit pedestrian states (e.g., position, velocity) or just sensor data; sensor data itself can include third-person sensors like external cameras, or be restricted to the robot's observation space. Pedestrian and robot data can be ground truth (either from a simulator or from motion capture in the real world) or noisily extracted with detection and tracking. The range of visibility of pedestrian is also important, as is whether the visibility is restricted to that of

robot sensors (including range, occlusions, directionality, and sensing delay) or ground truth (again, from the simulator or non-robot sensors). This is further discussed in Section 9.

- (14) *Metrics Collected*: Metrics transform raw collected data into standardized measures with shared definitions that can be compared across different algorithms, robots, and scenarios. Having shared metrics is important for communicating benchmarks, datasets, and simulators and is being looked at by several research groups; we present a view of this field in Section 6.
- (15) *Human Behavior Authoring Methods*: How are the human behaviors generated for the dataset or benchmark? E.g., real pedestrians, confederates of the experimenters, recordings, simulated via a standard social model, or generated by a policy. For simulated environments, these behaviors may include non-reactive (pedestrians driven by pre-recorded data), reactive (ORCA, social force, or generative models), and animated (character animations including static moving shapes and animated walking); for real environments, these may include natural behaviors, scripted behaviors, or randomized behaviors. For both simulated and real environments, the goals of the movement may be random, goal-directed, and potentially customized depending on the context.
 - Synonyms: Pedestrian Simulation, Crowd Simulation, Microscopic Crowd Simulation
- (16) *Robot Behavior Authoring Methods*: These are similar to the human behavior authoring methods, except there is no “real robot” class corresponding to “real pedestrians,” just the robot policies under test.
 - Synonyms: Agent Behaviors, Baseline Policies
- (17) *Simulation vs. Real*: Whether the dataset or benchmark is in simulation, in the real, or some combination of both. Sim and/or real: Is the benchmark operated in the simulation or in the real world? The participants noted that the simulation can be effectively used for issue discovery but cannot replace real-world testing.
 - *Subfactors*: Simulation Fidelity, which ranges from dots in an abstract geometrical space to fully rendered simulations. This includes both *Human Simulation Fidelity* and *Robot Simulation Fidelity*, as robots are simulated more often than full humans.

While the above factors are common across many social navigation instruments, one size does not fit all: there are additional factors particular to benchmarks, datasets, or simulators:

- *Dataset Properties* include trajectory count, kilometers traveled, duration traveled, number of pedestrian encounters, people count and density, robot count and density, primary view type (robot POV, pedestrian POV, third-person POV), dataset annotations, dataset diversity, privacy support, and sensor suite type (moving robot/stationary robot/third-person sensor suite).
- *Benchmark Properties* include the simulation platform, associated dataset, provided baselines, challenge leaderboard, downloadability, and the most recent update.
- *Simulator Properties* include abstraction level, scene representation, agent representation, physics simulation fidelity, level of detail of robot simulation (points, cylinders, robot morphologies), level of detail of pedestrian simulation (with or without gait), pedestrian visual fidelity (basic meshes, movements, photorealistic), handling of multiple agents (flow-based crowd, agent-based individuals), and environmental assets (realistic scenes or simulated layouts; indoors, outdoors, or abstract scenes).

When providing datasets, benchmarks, and simulators, researchers should identify the key features of the social environment and agent behavior that would enable other researchers to properly evaluate and replicate their work. While social robot navigation is highly varied, key features of the social environment might include the physical layout of the indoor or outdoor

environment and the expected human and robot density within it, as well as other contextual features such as intended task or human and robot roles. Similarly, key features of agent behavior should include both traditional navigation metrics (such as **Success Rate (SR)**, navigation time and distance, and environment collisions) and social metrics (such as safety, comfort, and agent collisions) so other researchers can understand how well robots perform at their tasks and how their behavior impacts others.

6 Social Navigation Metrics

Unlike traditional navigation, where the community largely agrees on a few evaluation metrics, such as **Success weighted by Path Length (SPL)** [2], finding a consensus for social navigation metrics is challenging. One reason for this difficulty is that we care about multiple aspects of human–robot encounters in social robot navigation, e.g., how safe a robot’s behavior is near people and how well the robot communicates its intent in order to facilitate motion coordination. Measuring any one of these factors from a human perspective is difficult, let alone deciding how to combine them into a single metric.

For example, while safety is a generally agreed-upon factor that drives the implementation and evaluation of social navigation systems, safety is a complex construct [80]. While one can think of physical safety in terms of collisions, as is often the case in the broader robot navigation literature (and is captured in Principle P1, Safety), safety also can be viewed from a psychological standpoint [67] (which might be captured in Principle P2, Comfort), or even in terms of not disrupting social and moral values [14] (which might be captured in Principle P5, Social Norms). Careful thought must be put into even obvious terms as the context of their usage may change their meaning (Principle P7, Contextual Appropriateness).

The next section provides a taxonomy of social navigation metrics, followed by a discussion of the challenges of measuring social navigation. We then present recommendations on metrics for social navigation, along with guidelines for using metrics to evaluate the success of social navigation systems.

6.1 Taxonomy of Existing Social Navigation Metrics

In the past years, a wide range of metrics have been proposed to quantitatively measure key aspects of social robot navigation and allow for fair comparison among social navigation solutions (see [50, 103] for recent reviews). We describe three ways social navigation metrics can be classified according to, (a) their nature, (b) the variable being modeled, and (c) their temporal scope. To fully classify a metric, it should therefore be classified according to the three taxonomies.

6.1.1 Taxonomy Based on Their Nature. We can distinguish two main groups of metrics, those that are algorithmic, and those that are not computed but surveyed (see Figure 6).

Surveyed metrics are usually human ratings of desired properties of social robot navigation, e.g., safety, comfort, or legibility. They can be classified into *questionnaire-based* (*in situ* or *ex situ*), where the ratings are explicitly requested, or *sensor-based*, where the ratings are transduced from sensor data. Although surveyed metrics are (arguably) the best way to measure social navigation success, they are expensive, difficult to scale, and time-consuming. While small-scale human studies are commonly conducted, results can have high variance and can be non-reproducible. To address this shortcoming, researchers have also created a variety of *algorithmic* metrics that serve as proxies for surveyed metrics. These algorithmic metrics are cheap to compute and reproducible, properties that are key for benchmarking. Unlike traditional navigation, where SPL [2] is a commonly accepted metric, social navigation has no single metric of reference. Instead, method comparison is usually performed using multiple metrics.

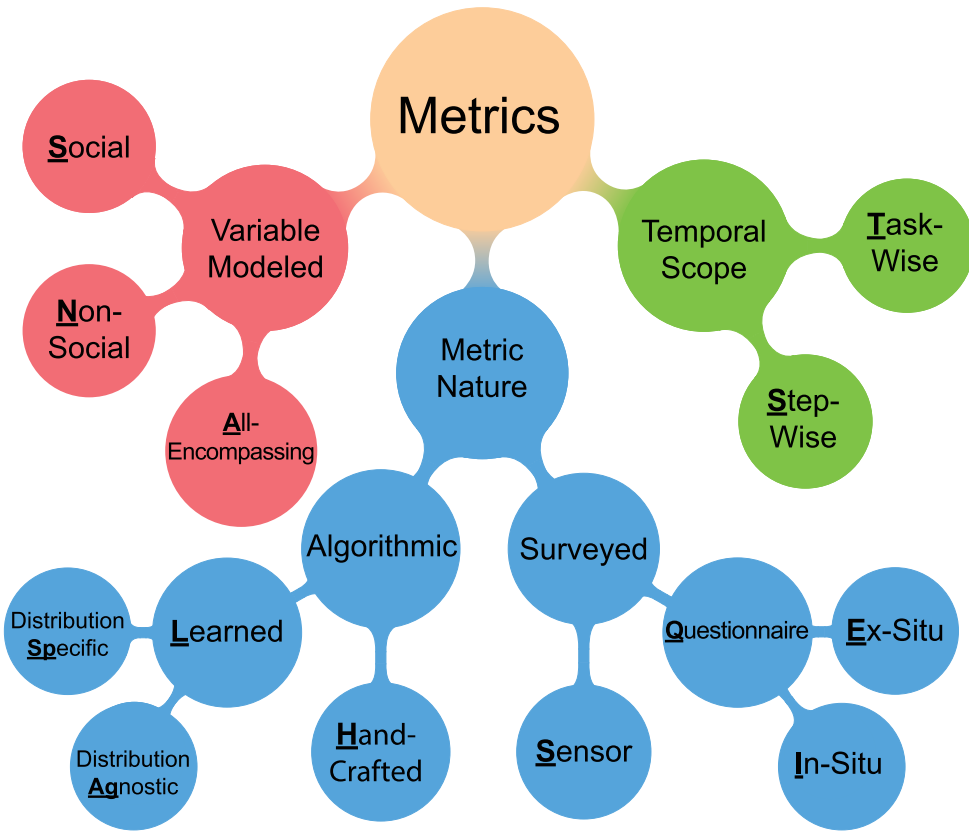


Fig. 6. The proposed taxonomy suggests classifying metrics according to three aspects: the type of variable (or variables) they model, their nature, and their temporal scope. To quickly identify metric types, we suggest using a three-letter code based on these factors. For example, Success Rate (SR) is a Non-Social, Hand-Crafted, Task-Wise metric; a sensor metric gauging moment-to-moment human facial reactions to robot behavior would provide a Social, Sensor, Step-wise metric; and a questionnaire asking about the overall quality of a robot's navigation would be an All-Encompassing, Questionnaire, Task-Wise metric.

Algorithmic metrics can be subsequently classified into hand-crafted and learned, based on whether they are the result of intuition and experience, or modeled using statistical analysis or machine learning. *Hand-crafted* metrics are objective (i.e., in what they compute, not necessarily their interpretation), scalable, and can be easily computed given certain assumptions, yet oftentimes they cannot fully capture the desired property of social robots. *Learned metrics* can be considered a compromise between survey-based and hand-crafted metrics. These evaluative models can be trained on large-scale offline datasets of human ratings and then used to score robot behavior. They are reproducible and have minimal inference cost, but compiling the necessary datasets can be very time-consuming. Learned metrics can be further split into *distribution-specific* metrics, which rely on assumptions on the properties of the dependent variable to model [123], and *distribution-agnostic* metrics, which aim to model these variables without making any relevant assumptions [4, 94].

6.1.2 Taxonomy Based on the Variable Being Modeled. Regardless of the nature of the metrics, algorithmic or surveyed, learned or hand-crafted, the variables they model can refer to different phenomena. Most common metrics assess either social or non-social aspects, but a metric could

also combine both into an all-encompassing metric (see Figure 6). *Non-social* metrics have generally been developed with PointGoal navigation in mind and focus on aspects such as path and energy efficiency or SR. They generally have the advantages of being objective, reproducible, and are usually fast to compute, but do not provide any social performance information. *Social* metrics focus on one or more social aspects of robot navigation, such as comfort, acceptance, trustworthiness, or predictability. *All-encompassing* metrics aim to model the overall scores humans would provide in robot navigation, considering both social and non-social aspects. Although these metrics would arguably be the most desirable, very few have been proposed [4, 30].

6.1.3 Taxonomy Based on Temporal Scope. It is also useful to consider metrics' temporal scopes, as they determine where a metric can be applied. Here we distinguish task-wise and step-wise metrics (see Figure 6). *Step-wise* metrics provide a score per timestep and are well-suited for path planning. *Task-wise* metrics are the most appropriate for benchmarking social navigation algorithms, as they provide a single score per task. Step-wise metrics can potentially be combined into task-wise metrics, such as by averaging across all steps within a task. However, not all moments within a social navigation task have an equal impact on social performance. To address this, task-wise metrics can also combine step-wise data with temporal data to capture features such as reversals in step-wise metrics over time, more heavily weight task-relevant periods of the task, or measure how long a robot was able to navigate with high-quality step-wise metrics [154]. For reinforcement learning-based social navigation, task-wise metrics (specifically, All-encompassing, Algorithmic, Task-Wise, or AAT, according to the proposed taxonomy) can be preferable over step-wise metrics, depending on their properties. Although using task-wise metrics would produce delayed rewards, a step-wise metric would only be advisable if its cumulative value reflects task performance, which is generally not the case.

6.2 The Challenges of Measuring Social Navigation

The evaluation of robot navigation has evolved as the field has matured, moving from success metrics to quality metrics to social metrics. Early work focused on success metrics that gauged whether the robot did its task, such as SR or kilometers without incident [95]. Later work proposed quality metrics that gauged how well the robot did its task, such as SPL [2]. More recent work such as [167] proposes social metrics that gauge how the robot behavior affects other agents, such as **Personal Space Compliance (PSC)** [164] or questionnaire-based metrics [34, 126, 171].

However, because social metrics involve robots interacting in complex real-world environments with humans whose learning changes their behavior over time, several additional factors must be considered to evaluate these social metrics accurately and reliably. These include (1) the challenges of dynamic environments, (2) the impact of long-term exposure on study participants, (3) how robot behavior may be changed by deployments, and (4) the limitations of metrics themselves.

6.2.1 The Challenges of Dynamic Environments. When measuring the performance of a social robot, an important consideration is the dynamic nature of the *environment* and of the other *pedestrians* around. These elements are often controlled when performing in-lab studies, but evaluation in the wild is much more intricate, especially when looking at longer periods of time, as robots become more and more capable of long-term deployment [13]. Results of performance measures like accuracy might be affected by simple changes such as lighting conditions and weather. Speed may be similarly affected by the percentage of remaining battery. While almost impossible to mitigate such effects, they should be acknowledged and highlighted when reporting relevant results.

6.2.2 Challenges Based on How Robots Change People. Yet, a more challenging aspect to measure is the dynamic nature of pedestrians when interacting, even casually, with a robot [135]. When people interact with a navigating robot for the first time, they adapt their beliefs and expectations to the robot's behavior, which can cause immediate improvement in various metrics such as fluency, time, and efficiency.

This phenomenon can be leveraged to train pedestrians around robots, rather than adapting the robot to the pedestrians. One such use of passive demonstrations was shown to significantly reduce the number of conflicts between a person and a robot passing one another in a hallway, by having the robot demonstrate in advance how it signals its intentions [43].

Beyond the novelty effect of first encounters, people will refine their behavior around a robot as they interact with it over time. A person will behave and react differently to a navigating robot on the tenth interaction than the hundredth [52, 58]. More research is needed on how people adapt to the presence of navigating robots [81], but studies of other social interactions such as asking favors [140] and information delivery [75] indicate that adaptation is likely.

When conducting research on social navigation in academia, it is not uncommon to rely on students, especially those with STEM and robotics backgrounds, as participants in an empirical study; this reliance on students has long been known to the psychological field as a potential source of bias [144, 148]. Measuring acceptance, animacy, and fluency can all be affected by this biased population that has been exposed to robots as part of their studies. Moreover, as robots are being deployed around campus, other students are also being exposed to these robots and thus over longer periods of time might also be biased in their expectations regarding the behavior of robots, based on their past encounters. Industry researchers in the symposium also reported differences in socialness ratings between naive subjects and robot researchers, and even between robot researchers and experienced robot “wranglers” who logged far more hours of direct robot time.

6.2.3 Challenges Based on How People Change Robots. When a robot is deployed for a long period of time, people may become familiar with it and thus more willing to accept risky behavior from it. For example, it may be able to drive at higher velocities, which will affect speed measurements, or it might get closer to others, which may affect efficiency and acceptance. Some symposium participants noted multiple instances of robots colliding with visitors after a good track record of avoiding collisions around the development team; a postmortem revealed that this was likely due to the development team implicitly learning to keep a collision-free distance.

Moreover, people's attitudes toward the robot over longer periods may require additional metrics that better capture how people perceive a long interaction with a robot. For example, in a long-term study of a socially assistive robot, the faults of the system did not affect the overall acceptance of the system by the participants [42]. Similar phenomena are likely to be observed in a social navigation context and thus should be reasoned about when measuring long-term interactions.

6.2.4 The Limitations of Social Metrics Themselves. In addition to these concerns, metrics themselves have challenges, including subjectivity and scaling, relevance and weighting, and the transferability of results between robot morphologies.

- (1) *Human Ratings Are Subjective.* Human ratings are by their very nature subjective, and they depend on many factors such as cultural context, environmental context, goals or priorities within a scenario, or their overall familiarity with robots. It is important to account for the factors that all human participants experience, as well as attempt to characterize unique factors relevant to the scenario that can affect how they perceive the scenario.
- (2) *Subjective Metrics Are Difficult to Scale.* Expanding to a larger participant pool can help to mitigate variations between individuals, but it can be hard to execute complicated scenarios

with a large number of participants. Creating analytical models of certain sub-elements of human reactions, such as how comfortable observers are with the proximity of the robot, can potentially be done with studies of a more targeted scope, and then used in broader models of human responses to robot behavior.

- (3) *Real-World Evaluation Is Difficult to Scale.* The closer a study can get to emulating a real-world scenario such as a busy street, a crowded airport, or a packed restaurant, the more it can capture the effectiveness of a robot in this domain. However, creating these scenarios in a laboratory environment is difficult. Eliciting natural behavior can be challenging, and many social environments have a large volume of people entering and exiting which can be hard to represent. Therefore, efforts are being made to record natural human behavior for use in simulations to address this issue, along with blending multiple metrics to account for the many aspects of a real-world deployment.
- (4) *Choosing Which Variables Are Relevant.* Measuring all possible signals humans generate in response to a robot is difficult. However, selecting any subset can neglect other useful signals. For example, using only 2D poses disregards other very important inputs such as facial expressions, gestures, or gaze [59]. Putting thought into which signals are most relevant to a scenario and able to be robustly collect them is important.
- (5) *Weighting Multiple Metrics.* The variety of useful metrics and their context dependence suggest applying an ensemble of metrics, weighted to account for the parameters of a specific scenario. The optimal method for doing this, however, remains open. It is worth considering if this weighting may vary not only across different environments but also over the course of a single path as the audience or priorities of the robot change.
- (6) *Non-Homogeneous Hardware.* Robots have varying sensors and actuators. While some only have access to their wheels' motors and a LiDAR, others can inform pedestrians of their presence and intentions, or share information using sound and visual cues. It is difficult to consider these additional aspects analytically, so standardized metrics do not take them into account. Unfairly, this limitation can make robots able to share such information appear less socially capable than they are.

6.3 Assessment of Existing Social Navigation Metrics

For all the reasons outlined above, quantifying the quality of different social navigation strategies is difficult. Beyond the inherent subjectivity of human ratings, social scenarios include many stakeholders with varying priorities and thus different assessments of the importance of metrics. For example, a passerby may primarily be focused on metrics of discomfort as the robot passes them, while the recipient of a package may be focused on metrics of task success. A warehouse manager may focus more on expediency, while a restaurant's customers may find excess speed or urgency unnerving. Social preferences also vary across cultures and groups. Therefore, there is no 'best' metric, only metrics appropriate for a given application or use case.

In all of these cases, the gold standard for evaluation is to collect subjective metrics reported by humans directly experiencing the interaction. However, subjective metrics can be difficult to scale to larger numbers of participants. A secondary issue is that the higher the density of feedback requested, the more disruption to the social scenario being measured. Both of these issues increase the demand for analytical or learned subjective metrics, and we discuss the considerations for this in Section 6.2.4.

New metrics are often created to address issues that come up in new scenarios, and as social navigation is being deployed in increasingly many new environments, more metrics are being created to address these scenarios. It is also unsurprising that new metrics will be of particular

value in the environments that demand their creation. This means that the number of metrics available to assess performance can be daunting, and their value is very context-dependent.

After reviewing the related literature, we did not find a convincing method to quantitatively compare different metrics to determine whether one is strictly better than the other. We suggest that when using any social navigation metric it is essential to note both the metric's original context and the current one it is being applied to. As mentioned, survey-based metrics are generally preferred for benchmarking, though their results are difficult to reproduce if they are not run correctly and they are resource-intensive. All-encompassing learned metrics would be the next best option for benchmarking, but unfortunately, none of the existing ones (see [4, 94, 123]) satisfy the requirements of all applications and scenarios. Metrics focusing on specific phenomena are of great importance when debugging and diagnosing an algorithm's flaws, but are sometimes difficult to use to compare disparate algorithms.

6.4 Recommendations for Metric Usage and Development

While many in the symposium argued that surveyed metrics are the gold standard, others pointed out that they are challenging to get right, expensive to collect and sometimes inappropriate (e.g., for evaluating ablation studies where safety cannot be guaranteed). Learned metrics have been proposed as a solution, but are not ready for adoption (e.g., no task-wise learned metrics are yet available). Therefore, to measure social robot navigation, we recommend a balanced approach, involving a common subset of hand-crafted metrics, recommendations for the iterative validation of surveys, and suggestions for future metric development.

6.4.1 Recommendations for Hand-Crafted Metrics. To ensure a systematic and objective comparison of social navigation algorithms, we suggest using a subset of existing hand-crafted navigation metrics. The suggestion includes success-related metrics accounting for success itself, collisions, and failures, as well as metrics related to trajectory properties and social aspects. These recommended metrics can be found in Table 1, along with descriptions of the phenomena accounted for, their required parameters, units, ranges, and references where a full mathematical definition can be found.

A relevant characteristic of many of these metrics is that their values, and more importantly what would be considered good ones, heavily depend on the task, the context where the experiments take place, and the parameters of the metric (see Table 1). It is therefore good practice to explicitly state the parameters used and context when reporting results.

It is also worth noting that the metrics in Table 1 are frequently reported as averages for a number of experiments rather than for a single trajectory (e.g., Success (S) is often found as the Success Rate (SR)). When reporting experimental results for multiple trajectories, providing distributional information in addition to averages allows us to show valuable information, including outliers. This is key when consistency is important, as it is the case of safety. Distributional information can be provided, for instance, as histograms.

6.4.2 Recommendations for Survey Development. Gathering human perception with surveys has a long history in HRI (see for example the discussion in [149]), but there is not yet a unified approach to questionnaire development in social robot navigation. Following the social scenario development approach of [34, 126, 171], we recommend an iterative approach in which versions of questionnaires are proposed and then empirically tested to determine their validity [106].

While survey validity is a complex topic worthy of its own book [88], several concerns for the design of questionnaires include assessing test-retest reliability (whether a survey gives stable results over time), construct validity (whether a survey measures what it purports to measure), and

Table 1. Suggested Hand-Crafted Metrics for the Evaluation of Social Navigation Systems

	Metric	Short	Description	Class	Parameters	Unit	Range	Cited
Success metrics	Success	S	Binary variable describing whether the robot reaches the goal. When averaged, it is referred as Success Rate (SR).	NHT	-	Boolean	{0, 1}	[2]
	Collision	C	Number of collisions in the trajectory. When averaged it is referred to as Collision Rate (CR).	NHT	Collisions to terminate episode	Collision	$0, \infty$	[72]
	Wall collisions	WC	Number of collisions against walls.	NHT	-	Collision	$0, \infty$	R@G
	Agent collisions	AC	Number of collisions against humans or robots.	NHT	-	Collision	$0, \infty$	R@G
	Human collisions	HC	Number of collisions against humans. Also called H-collisions [18].	NHT	-	Collision	$0, \infty$	R@G [18]
	Timeout before reaching goal	TO	Binary variable accounting for failures caused by a timeout.	NHT	Time threshold	Timeout	{0, 1}	R@G
	Failure to progress	FP	Number of failures caused by not decreasing the distance to the goal for a given period of time.	NHT	Distance and time thresholds	Failure	$0, \infty$	R@G
	Stalled time	ST	Time where the magnitude of the speed of the robot falls within a given threshold.	NHT	Distance and time thresholds	s	$0, \infty$	[161]
	Time to reach goal	T	Time between task assignment and completion.	NHT	-	s	$0, \infty$	[46, 72]
	Path length	PL	Length of the trajectory.	NHT	-	m	$0, \infty$	[46, 72]
Success weighted by path length	SPL	Success weighted using normalized inverse path length, i.e., weighted using path length divided by the max of the min distance and path length [2].	NHT	-	Success	$0, 1$	[2]	
Quality and social metrics	Velocity-based features	$V_{min}, V_{avg}, V_{max}$	Minimum, average, and maximum linear velocity on a trajectory.	SHT	-	m/s	$-\infty, \infty$	[72]
	Linear acceleration based features	$A_{min}, A_{avg}, A_{max}$	Minimum, average, and maximum linear acceleration on a trajectory.	SHT	-	m/s ²	$-\infty, \infty$	[72]
	Movement jerk	$J_{min}, J_{avg}, J_{max}$	Minimum, average and maximum linear jerk (i.e., the second-order derivative of the linear speed).	SHT	-	m/s ³	$-\infty, \infty$	[72]
	Clearing distance	CD_{min}, CD_{max}	Minimum and average distance to obstacles in a trajectory.	SHT	-	m	$0, \infty$	[72]
	Space compliance	SC	Ratio of the trajectory with the minimum distance to a human under a given threshold. If the threshold is $0.5m$, it is referred to as Personal Space Compliance (PSC) [83].	SHT	Distance threshold	m	$0, 1$	[164]
	Minimum distance to human	DH_{min}	Minimum distance to a human in a given trajectory.	SHT	-	m	$0, \infty$	
	Minimum time to collision	TTC	Minimum time to collision with a human agent at any point in time in the trajectory should all robots and humans move in a linear trajectory.	SHT	-	m	$0, \infty$	[12]
	Aggregated time	AT	Time taken for a subset of cooperative agents to meet their goals.	SHT	Cooperative agents' set	t	$0, \infty$	[167]

The first tranche in the table is traditional navigation metrics, included to ensure that social navigation systems do not regress on traditional navigation performance; the second tranche concerns aspects of the quality and socialness of navigation. Citations refer to either papers or challenges defining the term, or R@G for metrics from an unpublished Robotics at Google [53] robot deployment.

sources of bias (distorting factors that make the results hard to interpret). Assessing these factors involves reviewing both individual questions and the design of the survey as a whole.

For surveys as a whole, the longer a survey is, the less reliable the answers are [5, 49], and the more frequently surveys are given, the less likely people are to participate [128] a phenomenon known as *survey fatigue* or more generally *response burden* on participants. Reducing response burden is important not just to improve the quality of results but to respect the time of participants; nevertheless, issuing surveys multiple times can help measure test-retest reliability, issuing surveys to multiple populations can help measure bias, and including redundant questions can help measure construct reliability and question utility.

For individual survey questions, it is important to ask them using techniques that have been validated. For example, Likert scales [87] are widely used and provide a range of options like “Strongly Agree, Agree, Disagree, or Strongly Agree.” While it is tempting to use consistent wording

between questions, to reduce cognitive load on participants it is arguably better to formulate Likert scale responses so they form direct responses to each question, along with an option to indicate the question is not applicable. For example, to assess Principle P1, Safety, a question might ask “How safe was the robot’s motion?” and give the responses “Unsafe, Somewhat Unsafe, Somewhat Safe, Safe, or Not Applicable.” To evaluate overall navigation quality, some researchers have explored Likert scales similar to performance-based employee rating systems (e.g., “Outstanding, Very satisfactory, Satisfactory, Unsatisfactory, Poor”²) but no consensus yet exists here.³

Statistical analysis of experiments is discussed in depth in standard textbooks such as [32, 106], but we highlight some key concerns for social navigation. Terms such as “significance” often refer to *statistical significance*, a specific and contentious term in psychological literature [91] which should not be used unless the proper statistical tests are conducted. To do so, experimental conditions tested should be properly balanced counts (especially if questions are presented in multiple orders to reduce first-response bias, which creates sub-conditions within the experiment). Properly balanced experiment conditions enable the analysis of variance with tools like ANOVAs [32, 106, 152] and Cronbach’s alpha [106, 153]. Cronbach’s alpha in particular can help determine whether a given question is a reliable factor (see for example the discussion in Appendix D.4 of [171]) or should be dropped in future surveys in favor of more reliable questions.

6.4.3 Recommendations for Future Metric Development. Because conducting human surveys is expensive, symposium participants expressed interest in finding hand-crafted or learned algorithmic proxies. For example, to gauge safety, some benchmarks measure ‘time-to-collision’ [12]. To gauge comfort, some researchers [157, 158] have proposed some metrics to measure and limit the unnecessary motion and direction changes by the robot in the presence of humans; others have proposed ‘visibility indices’ which gauge the distance and angle at which robots first impinge on a human’s field of view [146, 147]. Legibility is also highly connected to the field of view, as observers need to be able to see a robot to make inferences about its movements and goals [154].

Future metric development should continue to explore learned or hand-crafted algorithmic proxies for surveyed metrics that can be efficiently computed, enabling the development of more efficient, repeatable, and scalable benchmarks. Validating these metrics might require collecting and annotating a large-scale dataset with both algorithmic and surveyed metrics, which could be used to compute the correlations between algorithmic proxies and their surveyed counterparts. This dataset could also be used to learn metrics that capture the surveyed results, as done in [4, 94]. Another approach to learning social metrics could be AutoRL [30], which learns dense reward functions useful for learning based on a sparse true objective; conceivably, data from surveys could be used as the true objective to train a learned social reward.

6.5 Metric Guidelines

In general, social navigation systems should not just be good social systems, but robust navigation systems, with a high SR, low collision rate, and a good SPL to ensure efficient experiments and the safety of human participants. Many of these features can be determined in simulation before deploying policies on potentially dangerous robots, but how social these policies are can only be determined with reference to human reactions to robot behavior - either through direct human surveys or learned metrics derived from human data.

Our recommendations for social metrics expand on these insights and summarize our broader recommendations from Section 6.4: use a broad set of navigation metrics to ensure robustness, attempt to use human survey metrics where feasible to evaluate socialness, validate those metrics

²<https://helpjuice.com/blog/employee-evaluation-form>

³<https://www.performyard.com/articles/performance-review-ratings-scales-examples>

with standard tools, guard against sources of bias, but use metrics appropriately in each stage of development.

- (1) *M1—Ensure Robustness Using Standard Metrics*: To ensure social navigation algorithms are good navigation systems, evaluations should report as many of the standard metrics of Table 1 as feasible.
- (2) *M2—Validate Policies with Algorithmic Metrics in Simulation*: Prior to deployment, algorithmic metrics such as those in Table 1 can enable fast evaluation to filter out bad policies prior to deployment.
- (3) *M3—Parameterize Metrics Appropriately in Context*: Social metrics with parameters, such as failure to progress or space compliance, should be appropriately parameterized given the current context, and parameters should be reported for those metrics that require them (see Table 1).
- (4) *M4—Use Learned Metrics to Help Iterate on Behavior*: Where learned metrics based on human data are available, they can provide insights to improve robot behavior, or acceptance tests prior to deploying policies on a robot.
- (5) *M5—Use Validated Surveys to Evaluate Social Performance*: Human surveys using validated instruments should be used to test the social navigation scenarios once the system is sufficiently robust and reliable.
- (6) *M6—Set Up Experiments Consistently to Avoid Bias*: Environmental complexity, subject selection, robot familiarity, survey fatigue, and differing experimental setups can all distort metrics. Use well-designed scenarios (see Section 7) to make metrics easier to compare.
- (7) *M7—Analyze Experiments Iteratively*: Social contexts are complex and getting metrics and surveys right are difficult; therefore, researchers should analyze experiments and iteratively improve them.
- (8) *M8—Report Results in Depth*: Point-wise estimates of single metrics can provide a distorted view of the performance of a system. Experimenters should report a battery of traditional, learned, and surveyed metrics, including both step-wise and task-wise metrics, as well as histograms or other distributional information.

7 Social Navigation Scenarios

Social navigation scenarios are specifications of categories of HRIs that facilitate the collection of data on human–robot behaviors and the communication of that data between researchers in a common language.

Fundamentally, social robot navigation involves robots interacting with humans. The situations in which we study these interactions range from controlled scenarios in the laboratory with small numbers of humans and robots to large-scale in-the-wild studies with dozens of robots in many uncontrolled pedestrian encounters. Following the symposium, participants engaged in substantial discussion regarding the relative importance of studies along this spectrum.

- Proponents of large-scale in-the-wild studies argue they have good ecological validity, uncover long-tail behaviors, enable more reliable assessments of human perceptions of robot behavior, and enable data collection for unsupervised and reinforcement learning. These studies can have good statistical reliability and can generate large datasets; however, they are expensive, time-consuming, require heavyweight software architectures, and are suitable for policies that are already reliable.
- Proponents of controlled in-the-lab scenarios argue they can also have good ecological validity, prevent regressions on known issues, enable scientific analysis of algorithms and behavior, and enable data collection for supervised and imitation learning. These studies are cheaper

to run, generate data quickly, require less complex software, and are more appropriate for iterating on policies in an earlier stage of development; however, it is harder to uncover long-tail behaviors or to generate large datasets.

Social navigation scenarios are a research tool to help bridge the gap between in-the-wild studies and controlled laboratory experiments by defining a clearly specified set of scenarios that can be identified in data collected in the wild, set up as experiments in the lab, and analyzed consistently based on the common definition. In the following, we will use the common `FRONTAL APPROACH` scenario (Table 3), which involves a robot and a human traveling in opposite directions in an environment large enough for them to pass each other, to illustrate how the same scenario can be used in field studies, robot deployments, laboratory experiments, and even imitation learning:

- *Field Studies*: A `FRONTAL APPROACH` definition could be used to identify HRIs in data collected from an in-the-wild field study, perhaps using the Behavior Graph method for analysis to distinguish them from other interactions such as intersections or overtaking. This suggests social navigation scenarios should be construed broadly so that long-tail behavior can be analyzed. For example, if during a `FRONTAL APPROACH` a pedestrian trips and is helped up by the robot, the pedestrian and robot may not exit the environment normally, but this is nevertheless an event that happens in `FRONTAL APPROACH` scenarios and should be captured in the data.
- *Robot Deployments*: A `FRONTAL APPROACH` definition could be used to set up a deployment to elicit desired behaviors—for example, a robot could be deployed traveling back and forth on a well-trafficked corridor. Thus, social navigation scenarios should be well-specified enough to eliminate counterexamples (for example, a corridor must be wide enough for both robot and human to pass to be considered `FRONTAL APPROACH`).
- *Laboratory Experiments*: A `FRONTAL APPROACH` definition can be used to set up a laboratory experiment (or regression test) to evaluate the performance of a given policy compared to alternatives—for example [126]. For the statistical analysis of this experiment to be successful, both metrics and criteria for a successful test need to be defined. For example, in a laboratory experiment, both the robot and the human need to attempt to cross the scenario environment, whereas in a robot deployment or field study, humans may stop to take a phone call, or a robot navigation stack may crash.
- *Dataset Generation*: When creating datasets for social navigation, scenario definitions can be used to curate existing data for inclusion into the dataset or to guide the setup of robot deployments or laboratory experiments designed to build that data. This scenario categorization can then be used to capture information about the dataset. For example, a pedestrian dataset could collect episodes each with a single scenario like `FRONTAL APPROACH`. In contrast, a crowd dataset, with larger numbers of pedestrians interacting in a larger area, might have episodes with several scenarios happening at once, like `FRONTAL APPROACH`, `INTERSECTION`, and `BLIND CORNER`.
- *Imitation Learning*: A `FRONTAL APPROACH` definition can also be used to collect episodes for imitation learning. For this to be successful, additional criteria must be defined—for example, which behaviors are considered successes or failures, or quality metrics which enable rating episodes as better or worse—so a high-quality set of episodes can be collected to enable training of a policy. In other words, while creating imitation learning is dataset generation, not all datasets are good for imitation learning. In the imitation learning use case, `FRONTAL APPROACH` episodes in which the robot or human fail to cross the scenario environment may be marked as failures so they can be excluded by the learning algorithm.

Table 2. Scenario Card for FRONTAL APPROACH

Social Navigation Scenario Card	
Scenario Metadata	
Scenario name	FRONTAL APPROACH
Scenario description	A robot and a human approach head-on in a passable space.
Scientific purpose	Low-density pedestrian scenario applicable indoors and outdoors.
Scenario Definition	
Geometric layout	A space wide enough for the robot and human to pass each other.
Intended robot task	The robot navigates from one side of the space to the other
Intended human behavior	The human navigates in the opposite direction of the robot.
Scenario Usage Guide	
Success metrics	SR, (No) Collisions
Quality metrics	P2: COMFORT, P3: LEGIBILITY, P4: POLITENESS
Ideal outcome	Robot goes around humans in a socially acceptable manner.
Failure modes	1. Robot collides with human 2. Robot fails to exit in time limit
Labeling criteria	1. Robot and human face each other 2. Robot and human move toward each other at the start of episode 3. Sufficient clearance exists for robots and humans to pass each other

Note that in theory, all social navigation principles could apply to the “Quality Metrics” section, but the Scenario Card should focus on the ones most relevant to this scenario. Arguably, all scenarios should focus on P1: SAFETY, but a scenario like ENTERING ROOM might focus on P5: SOCIAL NORMS (allow others to exit first), P6: AGENT UNDERSTANDING (determine if occupants are leaving), and P7: PROACTIVITY (moving to allow others to exit).

In the next section, we outline a methodology for identifying and specifying social navigation scenarios that support this breadth of usage, present a “Social Navigation Scenario Card” which enables scenarios to be clearly defined and disseminated, list common scenarios in the literature, and conclude with guidelines for scenario development and usage.

7.1 Scenario Design Methodology

Interactions occur between humans and robots wherever robots are deployed. Many of their interactions are unique, but others are common enough or important enough to warrant special treatment—whether we are looking for them in data collected from field studies, trying to recreate them in robot deployments and laboratory experiments, or trying to make them happen at scale for dataset generation or imitation learning. Having a clear definition of what behavior we want to identify, recreate or scale can ensure that we have good data, and can communicate it to other researchers.

To facilitate this, we propose the use of scenarios defining HRIs and propose the following methodology for defining scenarios relevant to social navigation. This consists of a three-step process:

- (1) *Define the Scenario*: Scenario definitions should be clearly specified enough to be identified in data or set up as an experiment. Thus, the scenario designer should consider:

- (a) *Intended Research Context*: The research topic the scenario is designed to explore, for example, low-density indoor pedestrian navigation or high-density outdoor crowd navigation. Many scenarios are general enough to apply to most research contexts.
 - (b) *Intended Robot Task*: The high-level objective of the robot, for example, navigation between two points, visual navigation, or guiding a person to a goal location.
 - (c) *Intended Human Behavior*: The high-level objectives of nearby people, for example, navigating between two points, delivering a package, or following the robot to a goal location.
 - (d) *Success Metrics*: the criteria that define the successful completion of the robot's task. While scenarios may play out in the wild in a variety of ways, the robot's task should be well-specified enough that it is unambiguous whether it succeeded.
- (2) *Evaluate the Definition*: The way scenarios are designed affects the aspects of robot behavior that they evaluate and what behaviors they elicit in humans, sometimes in unexpected ways. Therefore, we propose that designers should evaluate scenarios after their initial design, assessing their ability to measure the desired robot behaviors. Well-designed scenarios should have the properties of commonality, flexibility, and fitness to purpose.
- (a) *Commonality*: Well-designed scenarios should be designed to evaluate the designer's intended criteria while maintaining identifiable characteristics that allow it to be grouped and compared with similar scenarios in use in the community. Common categories of scenarios are listed as sections of Table 3, and include "approach" or "hallway" scenarios involving robots approaching people or objects from specific directions, "intersection" scenarios where robots and humans cross paths, and "interpersonal" scenarios such as robots leaving or joining conversational groups. Scenario designers should compare their scenarios with these common scenarios to avoid introducing redundant scenarios when existing scenarios are available.
 - (b) *Flexibility*: Well-designed scenarios should be broadly specified enough to capture the full range of behavior that occurs in the wild. It is important to avoid "solutionizing" in which scenarios prescribe intended robot or human behavior so narrowly that naturally occurring variants are included. Instead, scenarios should have broad, flexible definitions that enable them to capture behaviors that happen, along with clear success metrics to evaluate whether that behavior came out as intended.
 - (c) *Fitness to Purpose*: Well-designed scenarios should allow the scenario designer to evaluate the adaptations of robot behaviors with which they are concerned. For example, researchers have explored how proactive robot behaviors can improve social interactions during navigation [73]. To evaluate robots that exhibit proactive cooperation, the scenario must be flexible enough to allow proactive cooperativeness interactions to occur. Early drafts of scenarios should be piloted to confirm that desired behaviors can be detected and elicited and that success metrics measure what is intended.
- (3) *Communicate the Definition*: Once a scenario has been evaluated, it should be communicated clearly and consistently. A scenario definition should be specific enough to replicate, so other researchers can identify occurrences of the scenario in their data, recreate it in the laboratory, and determine whether instances of a scenario correspond to the intended outcome for the human or robot.

Social navigation scenario development can be seen as a step toward the more formal scenarios engineering approach being adopted in intelligent vehicle research [84–86]. To facilitate communicating scenarios, we propose a social navigation scenario card, presented next.

Table 3. Example Social Navigation Scenarios

Scenario Name	Scenario Description	Phys. Env.	Geom. Layout	Scientific Purpose	Robot Role	Robot Task	Human Behavior	Ideal Outcome	Related Scenarios	Cited In
Hallway Scenarios										
FRONTAL APPROACH	Pedestrian and robot approach head-on.	Generic	Passable space	Pedestrian interaction	Any	Navigate A to B	Navigate B to A	Robot / humans pass	PED. OBSTRUCT	[50, 126, 167]
PEDESTRIAN OVERTAKING	Pedestrian overtakes moving robot.	Generic	Passable space	Pedestrian interaction	Any	Navigate A to B	Navigate A to B	Human passes robot	DOWN PATH	[26]
ROBOT OVERTAKING	Robot overtakes moving pedestrian.	Generic	Passable space	Pedestrian interaction	Any	Navigate A to B	Navigate A to B	Robot passes human		[50, 167]
INTERSECTION NO GESTURE	Robot and human cross at intersect.	Indoor	Intersection	Pedestrian interaction	Any	Navigate A to B	Cross navigate	Both pass no collision		[27, 50, 161, 167]
INTERSECTION GEST. WAIT	Robot told to wait at intersection.	Indoor	Intersection	Pedestrian interaction	Servant	Navigate A to B	Cross navigate	Human goes then robot	GESTURE PROCEED	[126]
BLIND CORNER	Robot and human meet at blind corner.	Indoor	Corner	Pedestrian interaction	Any	Navigate A to B	Navigate B to A	No collision / obstruction		[126, 171]
Doorway Scenarios										
NARROW DOORWAY	Robot and human at narrow doorway.	Indoor	Room and door	Pedestrian interaction	Any	Navigate A to B	Navigate B to A	No collision / obstruction	NARROW ARCH	[126]
ENTERING ROOM	Robot enters room occupied by human	Indoor	Room and door	Pedestrian interaction	Any	Navigate out to in	Navigate in to out	Robot lets human exit	ENTERING ELEVATOR	R@G
EXITING ROOM	Robot exits room while person enters.	Indoor	Room and door	Pedestrian interaction	Any	Navigate in to out	Navigate out to in	Robot exits first	EXITING ELEVATOR	R@G
Interpersonal Scenarios										
JOINING A GROUP	Robot joins group of robots or people.	Generic	Open space	Group Interaction	Any	Navigate to group	Continue convers.	Robot joins group	LEAVING A GROUP	[50, 161]
FOLLOWING	A robot follows a person.	Generic	Walking space	Joint navigation	Servant	Follow human	Lead robot	Robot follows person	ACCOMPANY PEER	[50]
LEADING	A robot leads a person.	Generic	Walking Space	Joint Navigation	Leader	Lead human	Follow robot	Robot leads person	TOUR GUIDE	[50]
Crowd Scenarios										
CROWD NAVIGATION	A robot navigates through a crowd.	Generic	Passable space	Crowd navigation	Any	Navigate thru	Mill about	No collision / obstruction	ROBOT CROWDING	Various
PARALLEL TRAFFIC	Crowd moves parallel to the robot.	Generic	Passable space	Crowd navigation	Any	Navigate A to B	Mill from A to B	No collision / obstruction	CIRCULAR CROSSING	[167]
PERPENDIC. TRAFFIC	Crowd moves perpendicular to robot.	Generic	Intersection	Crowd navigation	Any	Cross navigate	Mill from A to B	No collision / obstruction	PLAZA CROSSING	[167]
Specialized Scenarios										
OBJECT HANDOVER	A robot hands an object to a human.	Generic	Passable space	Interactive navigation	Servant	Deliver object	Receive object	Human takes object	ROBOT COURIER	[161]
CRASH CART	Robot delivering a medical product.	Indoor	Passable space	Interactive navigation	Leader	Deliver object	Receive object	Delivery of medicine	FOOD DELIVERY	This article

For illustrations of the geometric layout, see Figure 7. Closely related scenarios are listed in the second-to-last column. Citations refer to either papers or challenges defining the scenario, or R@G for scenarios from an unpublished Robotics at Google [53] deployment, developed according to the protocol in [126].

7.2 Social Navigation Scenario Cards

Ideally, a social navigation scenario consists of a well-defined social interaction including robots performing tasks, people performing behaviors, and relevant features of their environment. This definition should be specific enough that an encounter can be labeled as an instantiation of the scenario, but loose enough that it captures a wide variety of behaviors. For ease of reusability,

scenarios should ideally be realistic in that they represent real-world scenarios, scalable in that they can be set up at low cost, and repeatable in that the same scenario could be conducted many times under similar conditions. However, scenarios may encompass a wide variety of situations, from a simple `FRONTAL APPROACH` of a robot and human passing each other up to the complexity of a robot navigating a crowd exiting a stadium, and scenario cards should remain flexible enough to capture these use cases.

Following work on “model cards” in the machine learning community [104], we propose a “Scenario Card” approach to defining scenarios which labels the scenario with a set of features that unambiguously define it. scenario cards have the following three major elements: (a) scenario metadata that defines the name, description, and scientific purpose of the scenario; (b) scenario definition which clearly describes the environment, intended human behavior, and intended robot task; and (c) a scenario usage guide, which provides additional information for specialized usages such as evaluation metrics, success, and failure criteria.

7.2.1 Scenario Metadata. The scenario metadata identifies a scenario in an unambiguous way for other researchers, including the type of the scenario (doorway, hallway, etc.), its name, its description, and its scientific purpose (crowd navigation, low-density pedestrian, interactive, etc.). For example, a head-on pedestrian approach scenario might be labeled `FRONTAL APPROACH`, which we will use as a running example.

- *Scenario Type:* Scenarios can be grouped into broad classes such as head-on approaches vs. intersections, doorways, and elevators, crowd vs. group, interactive and accompanying, and so on. Identifying the group a scenario belongs to can help researchers decide whether to include it for coverage or exclude it as redundant.
- *Name:* The scenario should be given a unique name that does not conflict with existing scenarios used within the community.
- *Description:* The scenario should have a brief description that communicates what is intended to happen in it.
- *Research Context:* Scenarios often are targeted at specific scientific purposes along various dimensions of research interest - for example, indoor low-density pedestrian scenarios or outdoor high-density crowd scenarios. Key elements that are often distinguished include:
 - *Location: Indoor, Outdoor, or General.* Indoor and outdoor navigation have different constraints and are often studied separately; however, some scenarios, like `FRONTAL APPROACH`, can occur in many contexts.
 - *Density: Pedestrian or Crowd.* Low-density pedestrian studies (where robots encounter only a few individuals at a time) are often studied separately from high-density crowd scenarios (in which people exhibit qualitatively different behavior).
 - *High-Level Task: Navigation, Delivery, or Interaction.* Many scenarios focus on pure navigation tasks, but others involve object delivery, interacting with humans, leaving and joining groups, and so on.

7.2.2 Scenario Definition. The scenario definition defines roughly what is meant by the scenario, in a precise but broad way that allows scenarios to be identified but not so restrictive as to prevent recording important behaviors. For example, `FRONTAL APPROACH` scenario definition should enable us to recognize that a robot and human are approaching head-on, but at the same time capture an interaction where the human changes their direction or stops to answer their phone.

- *Geometric Layout:* Scenarios often occur in specific physical environments, such as corridors, doorways, blind corners, or near elevators. The important features of the environment should be noted; features that can vary should also be noted so the scenario is not overspecified.

- *Intended Robot Task*: The number of robots and their desired behaviors should be recorded. A robot simply navigating around a pedestrian has different behaviors than one which is specifically attempting to navigate to a given target. Typical robot tasks are a robot heading to a pre-defined position, a robot guiding a person to a destination, or a robot delivering an item.
- *Intended Human Behavior*: The expected human behavior should be specified. In the scenario definition, behaviors should be specified clearly enough to recognize the behavior in data or to enable a human to attempt to perform it, but not too specific that diverse behaviors could not be collected.

7.2.3 Scenario Usage Guide. The scenario usage guide specifies how the scenario is used in practice and contains additional information that goes beyond the definition, such as idealized outcomes or instructions for human confederates for experimental setups. This is the place where a **FRONTAL APPROACH** scenario would express that the ideal outcome is that the robot and human pass each other without incident and exit on the opposite sides of the scenario area.

- *Labeling Criteria*: A clear set of criteria should be provided so that scenarios can be labeled in logs data or rejected in the event of a structured run. For example, for an intersection scenario, one could demand that the robot passes within two meters of the human and that their intended paths at least potentially cross.
- *Success and Quality Measures*: To evaluate how well the robot performed in the scenario, we may also want to specify “Success Measures” and “Quality Measures” specific to a scenario such as the ability of the robot to ensure legibility of its behavior, to limit and control disturbance, to facilitate human action and situation understanding, and so on.
- *Ideal Outcome and Failure Modes*: To enable researchers to evaluate robot performances in episodes for imitation learning or data analysis, the ideal outcome should be outlined, for example, that a robot should not collide with a human at a blind corner. Also, to help debug scenarios and guard the safety of human participants, failure modes such as colliding with walls, or stopping dead after a near-collision, should be outlined. We include failure modes in ideal outcomes in the scenario usage guide and not the definition because researchers interested in data collection do not want to artificially exclude arbitrary outcomes that can occur in the wild; however, this is critical information for imitation learning researchers trying to craft behavior.
- *Human Behavior Playbook*: If a scenario is designed to be created in a repeatable way as part of an experiment, a specific script or rubric should be provided so that the participants can perform their roles appropriately. For example, intended human behavior might be travelers in a crowded railway station, or workers going alone or with colleagues in an office context. These could include variations in the behaviors: for instance, some travelers might be in a hurry while others have more time. Also, there could be several categories of users in a given context that might act and react differently.
- *Contextual Information*: Principle P7 notes that a robot’s behavior should depend on context: for instance, a robot should behave differently if the place is very calm and needs silence or if it is a busy place, so ideally robots should recognize in which contextual situation a scenario is happening. Success metrics, ideal outcomes, failure modes, human behavior, and more can be altered by the context, so it can be useful to outline any important contextual variants of the scenario and how they affect intended robot or human behavior.

7.3 Example Social Navigation Scenarios

To effectively evaluate social navigation policies, they should be exercised in a set of scenarios that address the common use cases that come up in their intended context. For example, policies for

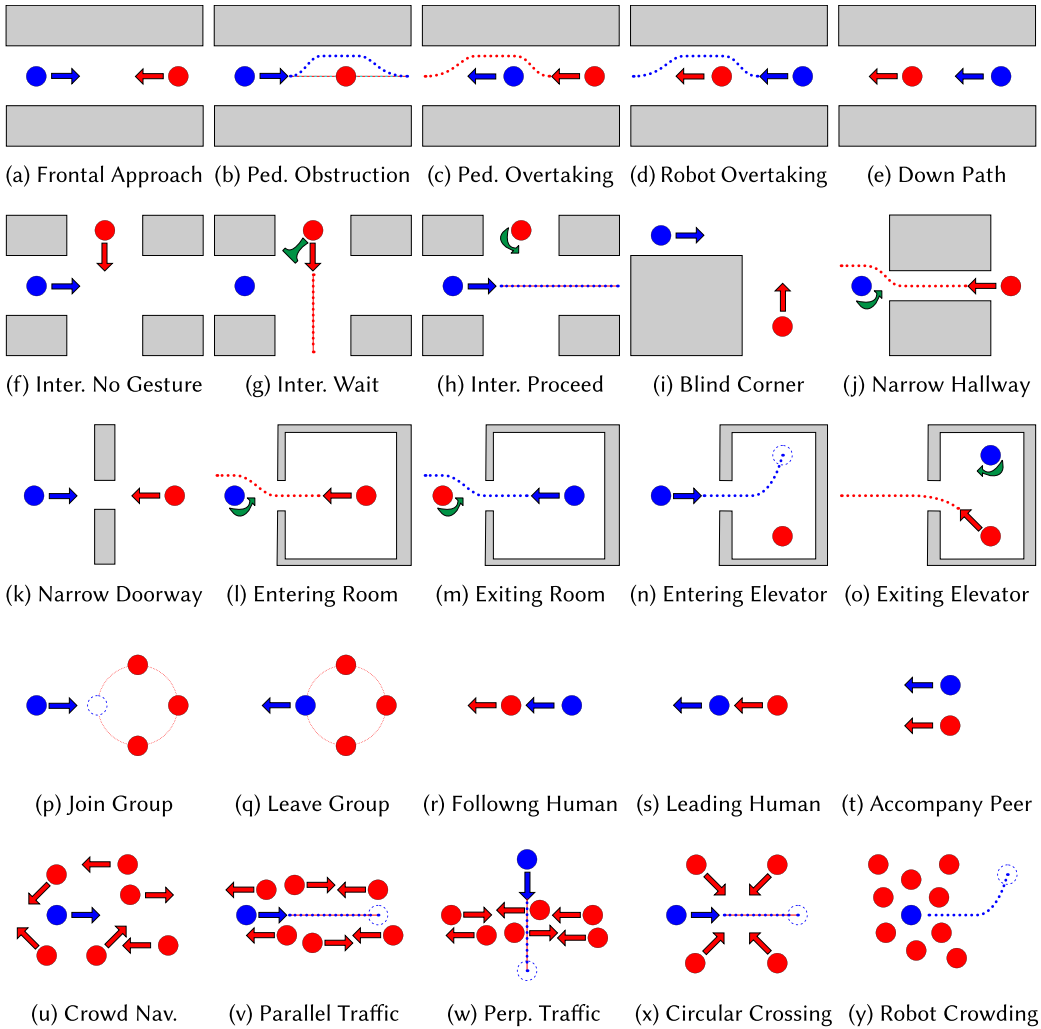


Fig. 7. Geometric layout and intended human and robot behavior for example social navigation scenarios from Table 3. The blue circles, arrows, dotted lines, and dotted circles represent the robot, its direction of motion, its intended path, and its intended destination, respectively; red figures represent the corresponding items for humans. Gray backgrounds represent obstructions, while green figures represent signals or gestures emitted by an agent such as a gesture to stop or go ahead. Note that sample paths and gestures are provided as examples to make the graphics clear; the actual scenario card definition should be flexible enough to capture a range of behaviors.

interacting with pedestrians in low-traffic areas should be tested in common hallway and doorway scenarios, policies designed to navigate through crowds should handle common scenarios like traveling parallel to or perpendicular to the flow of traffic, and policies for interaction should handle scenarios like leaving and joining groups.

Ideally, policies should be evaluated using standard benchmarks as discussed in Section 8; however, for a new research purpose a suitable benchmark may not yet exist. Nevertheless, researchers should try to find scenarios that are already in use in the field and apply them as comprehensively

as possible so that the evaluation of policies is meaningful and can be reasonably compared to other work in the literature.

To facilitate this process, we summarize common social navigation scenarios in Table 3. Dozens of social navigation scenarios have been proposed in the literature, and we cannot list them comprehensively; however, we provide references above to relevant prior work using these scenarios where available. Also note that scenarios can be grouped into a broad variety of scientific purposes, including pedestrian navigation, crowd navigation, and interaction scenarios, which can help guide researchers in their selections.

- (1) *Pedestrian Navigation*: Low-density pedestrian navigation scenarios study how pedestrians interact with robots a few at a time and include common hallway and door interaction scenarios. Common pedestrian scenarios often include **FRONTAL APPROACH** where a human pedestrian approaches a moving robot head-on, **ROBOT OVERTAKING** where a robot overtakes a slower-moving human, **INTERSECTION** where a robot passes a human at a right angle, **BLIND CORNER** where a robot and a human pass each other at an angle with poor visibility, **NARROW DOORWAY** where which a robot and a human attempt to exit a doorway in opposite directions, and so on.
- (2) *Crowd Navigation*: High-density crowd navigation scenarios study how robots can navigate dense human crowds. While there exist commonalities shared by navigating pedestrians [107], qualitatively unique behaviors can emerge during crowd navigation, for example when walking agents form social groups [108]. Common scenarios for crowd navigation include **PARALLEL TRAFFIC** where a robot is going with or against the flow of moving pedestrians, **PERPENDICULAR TRAFFIC** where the robot must cross a flow of pedestrians, **CIRCULAR CROSSING** and **RANDOM CROSSINGS** where pedestrians are crossing a plaza or room, and even **ROBOT CROWDING** where a robot is surrounded by stationary pedestrians and must extricate themselves.
- (3) *Interaction*: Interaction scenarios involve a task that places constraints on robot navigation, such as group navigation skills like **JOINING GROUPS** of pedestrians in conversations, **LEAVING GROUPS** of pedestrians, or interactive skills such as **OBJECT HANDOVER** where robots deliver or receive an item, **QUESTION ANSWERING** where robots answer or ask questions, and **CONTINUOUS MONITORING** where a robot observes individuals exercising or performing another activity.

The columns of Table 3 capture many of the features of the social navigation scenario card, though we cannot list all of them for space. Referring back to our running example, Table 2 shows an example of how the social navigation scenario card could be retroactively applied to one of the most common social navigation scenarios, **FRONTAL APPROACH**, which appears in [50, 126, 167] among others.

7.4 Scenario Guidelines

Scenario guidelines can be broken into three groups following the methodology outlined above: guidelines for new scenario development, guidelines for evaluating scenarios for research purposes, and guidelines for communication. For new scenarios, we propose the following guidelines:

Guideline N1: Specify Research Context. New social navigation scenarios should clearly define the research context under which they are expected to apply.

Guideline N2: Define Intended Robot Task. New social navigation scenarios should clearly define the task the robot is expected to accomplish and not just the start and end poses for navigation alone.

Guideline N3: Define Intended Human Behavior. Scenarios should specify what human participants are intended to do in the scenario.

Guideline N4: Define Success Metrics. Scenarios should include metrics to gauge the success or failure of the task.

To evaluate the usefulness of scenarios, we recommend:

Guideline N5: Cover Common Scenarios. To adequately evaluate social navigation algorithms, researchers to try to include good coverage of scenarios which are used commonly in the field, such as those listed in Table 3.

Guideline N6: Ensure Scenario Flexibility. Scenarios should be broadly specified enough to capture the full range of behaviors that can occur.

Guideline N7: Evaluate Fitness for Purpose. Scenarios should identify or elicit the desired behaviors and enable the desirable properties of robot behavior to be evaluated.

Finally, we recommend the use of scenario cards as a standard communication format:

Guideline N8: Use Scenario Cards. When communicating scenarios—either new scenarios, or specializations of scenarios used for specific research purposes—use the scenario card format to clearly communicate scenario content.

8 Social Navigation Benchmarks

Social navigation benchmarks improve upon individual laboratory experiments or well-defined scenarios by collecting a set of scenarios into a benchmark suite with well-specified metrics, enabling the comparison of a variety of different methods against each other. However, existing benchmarks focus on different aspects of the social navigation problem outlined in Section 3, using different permutations of the factors we outlined in Section 5. Hence, the results of these benchmarks may be more or less useful for researchers investigating different aspects of the social navigation problem. In this section, we advocate a set of criteria to make benchmarks useful across the social navigation community and review existing benchmarks in use with regard to these criteria.

First, we analyze benchmarks in use in the social navigation community, grouping them into benchmarking protocols, benchmarking environments, and benchmark challenges. Then, we analyze the strengths and weaknesses of these benchmarks, abstracting out criteria for good social navigation benchmarks, including evaluating social behavior using quantitative metrics and well-validated questionnaires grounded in human data. Finally, we make recommendations on how to improve the state of social navigation benchmarking and discuss how social benchmarking could be integrated with standard navigation benchmarks as regression tests of navigation behavior, which ensure that previously successful behaviors do not degrade as changes are made [113, 165].

8.1 Expanding the Factors for Benchmark Analysis

In addition to the factors listed in Section 5.1, additional aspects must be considered for benchmarks:

Simulation Platform. Benchmarks must specify how to set up an evaluation, but are more useful if that evaluation is already set up on a commonly available simulation platform.

Associated Dataset. Some benchmarks specify one or more datasets of reference behaviors used for comparisons.

Provided Baselines. Some benchmarks specify a set of baseline policies that can be used for comparisons.

Challenge Leaderboard. Benchmark challenges may also provide a leaderboard to enable policies among different teams to be compared publicly.

Downloadability. Ideally, a benchmark should include a downloadable software suite to enable replication of results.

Most Recent Update. Because software platforms evolve, benchmarks should be updated recently to ensure they are usable with current hardware and software.

Robot Hardware Platform. To make benchmarks most useful, they should support a wide variety of robot morphologies or custom robot morphologies so researchers have the best chance of generating comparisons for their target platform.

Human Behavior Authoring Methods. Benchmarks must include agents other than the robot, whether human or other robots. Support for realistic human behavior or replayed datasets can improve a benchmark's fidelity and usefulness.

8.2 Existing Social Navigation Benchmarks

In the social navigation literature, the term “benchmark” is sometimes applied to labeled datasets of reference behavior, which we discuss in Section 9. In this section, we focus specifically on social navigation benchmarks that combine at least three components: (a) a social navigation system (such as a simulator) that can run algorithms and pedestrians (b) in well-defined scenarios (c) with metrics for evaluation; these benchmarks may optionally specify datasets of human or robot behavior for comparisons. Full benchmarks can be broken into three classes: (1) *benchmarking protocols* which enable the construction of experiments along well-specified principles, like the SOCIAL NAVIGATION PROTOCOL [126], (2) *benchmarking environments* which enable comparison of algorithms against baselines in environments, including DYNABARN [109], GYM-COLLISION-AVOIDANCE [39], HU NAVSIM [122], and SocNAV BENCH [12], and (3) *benchmark challenges* which also provide a platform or forum to share results, including CROWDBOT [45], IGIBSON [83, 145], and SEANAVBENCH.⁴ In the following, we describe these benchmarks; see Table 4 for a side-by-side comparison based on the previously described factors and Figure 8 for a visual description of some of the more commonly used benchmarks.

8.2.1 Benchmarking Protocols. The SOCIAL NAVIGATION PROTOCOL [126] is an industry benchmark proposed by Robotics at Google [53] and used in [34, 126, 171] to evaluate the performance of a series of learning-based model predictive control policies for social robot navigation (though the protocol was intended to be applicable to the evaluation of any policy, learning or not). This protocol involves selecting social navigation scenarios of interest, such as Frontal Approach, Blind Corner, Corridor Intersection, and so on. Each scenario's HRI is defined by the start and end of the robot's trajectory and a short description of what is expected to happen for the human. This serves two purposes: enabling the collection of expert human trajectories for training social navigation policies, and evaluating policies on the same scenarios with low variability. Over the course of [34, 126, 171], the protocol was iteratively improved. For example, the questionnaire proposed in [126] was analyzed in [171] to identify reliable factors according to Cronbach's alpha, which were used to update the questionnaire for [34], which enabled more extensive analysis. While the SOCIAL NAVIGATION PROTOCOL can be applied to a wide variety of setups, it does not provide a downloadable, simulated environment, and must be manually set up for each experiment.

8.2.2 Social Navigation Benchmarks. ARENABENCH [72] is a downloadable, simulated social navigation benchmark designed to test how navigation algorithms perform under different tasks. Building on the 2D Flatland⁵ and 3D Gazebo [76] simulators and the Pedsim [51] implementation of

⁴<https://seanavbench.interactive-machines.com/>

⁵<https://flatland-simulator.readthedocs.io/en/latest/>

Table 4. Characteristics of Existing Social Navigation Benchmarks

Benchmark	Arena-Bench	CrowdBot	DynaBarn	gym-coll-avoidance	Hu-NavSim	iGibson	SocNav-Bench	SeaNav-Bench	Social Nav. Protocol
Factors for Analysis									
Benchmark classification	Benchmark	Challenge	Benchmark	Benchmark	Benchmark	Challenge	Benchmark	Challenge	Protocol
Benchmark context and scope	Dynamic obstacle benchmark	Crowd simulation benchmark	Dynamic obstacle benchmark	Collision avoidance benchmark	Human simulation benchmark	Social navigation benchmark	Social navigation benchmark	Social navigation benchmark	Human-robot expt. design
Physical environment	Indoor	Indoor	Synthetic	Synthetic	Indoor	Indoor	Indoor and outdoor	Indoor and outdoor	Principally indoor
Intended human user type	Synthetic pedestrian	Synthetic pedestrian	Varied human motion	Synthetic pedestrian	Varied resp to robot	synthetic pedestrian	Synthetic pedestrian	Synthetic pedestrian	Human coworkers
Supported robot tasks	Navigation	Navigation	Navigation	Navigation	Navigation	Navigation	Navigation	Navigation	Navigation
Social scenarios evaluated	3 worlds, 5/10 peds	Basic crowd scenarios	60 crowd scenarios	Multi-agent scenarios	House, cafe, warehouse	15 house scenes	5 curated environments	TBD	6 social nav. scenarios
Coverage of corner cases	Diversity, random	Not tested	Diversity, random	Not tested	Not tested	Not tested	Not tested	TBD	Not specified
Simulation platform	Flatland, gazebo	Unity	Gazebo	Custom	Gazebo	iGibson	Soc-NavBench	SEAN 2.0	None
Benchmarking dataset	None	CrowdBot	None	None	None	None	UCY and ETH	UCY and ETH	None
Human behavior authoring	Pedsim	UMANS	Multiple algorithms	Baseline policies	Soc. force, behav. tree	ORCA	Replay, planned	Replay, soc force	Scripted
Human simulation fidelity	Walking humans	Walking humans	Moving cylinders	Moving cylinders	Walking humans	Moving humans	Walking humans	Walking humans	Real humans
Supported robot embodiments	Jackal, burger, robotino	Pepper, wheelchair, CuyBot, Qolo	Custom robots, ClearPath Jackal	Cylinders	ROS Gazebo-Compatible	8 real, 2 mujoco	Simulated mobile robot	Fetch, Jackal, Turtlebot, Warthog	Human-scale robots
Communication modalities	None	None	None	None	None	None	None	None	Human gestures
Challenge leaderboard	None	None	None	None	None	2021	None	2022	None
Benchmark last updated	2022	2021	2023	2022	2023	2021	2022	2022	2022
Guidelines for Benchmarks									
B1: Evaluate Social Behavior	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
B2: Quantitative metrics provided	Many	Many	Succ. rate	Succ. rate, Time2Goal	Many	Succ. rate, PSC	Many	Many	No
B3: Baseline policies Provided	SOA nav, RL policies	No	SOA nav, RL policies	SOA social, worst-case	No	SOA RL	SOA social, worst-case	SOA social, worst-case	No
B4: Scalable, repeatable	Simulated, download	Simulated, download	Simulated, download	Simulated, download	Simulated, download	Simulated, download	Simulated, download	Simulated, download	Setup req., phys. eval
B5: eval grounded in Human Data	No	No	Demo. pipeline	No	No	No	No	SEAN-EP extension	Yes
B6: Use validated instruments	No	No	No	No	No	No	No	No	Validation in process

the **Social Forces Model (SFM)** [62], ARENABENCH provides the ability to evaluate both classical and learning-based approaches in the **Robotic Operating System (ROS)** [129] framework. In addition to providing tools for automatically and manually creating scenarios, ARENABENCH supplies both the non-learned baselines MPC [137], DWA [44], TEB [136] and the learned baselines NAVREP [38], Gring [55] as well as ARENABENCH’s own trained ROSNAV approach. Supported robots include the Robotis Turtlebot3, ClearPath Jackal, and Festo Robotino 4. ArenaBench provides a variety of navigation metrics including SR, collision, time to goal, path length, velocity, acceleration, jerk, curvature, angle over length, roughness, and clearing distance. However, ARENABENCH does not at this time support human evaluation of robot behavior.

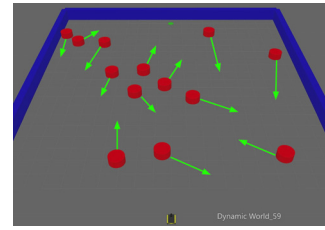
DYNABARN [109] is a downloadable, simulated social navigation benchmark designed to test how algorithms respond to a variety of different pedestrian models. Building on the BARN navigation benchmark [124], DYNABARN provides 60 environments in the Gazebo simulator. DYNABARN



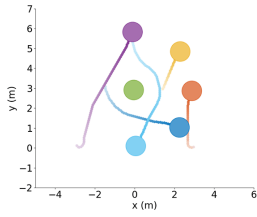
(a) ArenaBench



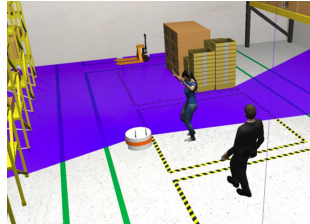
(b) CrowdBot



(c) DynaBarn



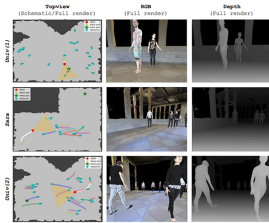
(d) gym-collision-avoidance



(e) HuNavSim



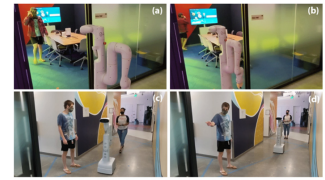
(f) iGibson



(g) SocNavBench



(h) SEANavBench in SEAN 2.0



(i) Social Navigation Protocol

Fig. 8. Commonly used social benchmarks. Benchmarks range from abstract tests of dynamic obstacle avoidance to simulated interactions with moving humans of varying degrees of fidelity to protocols for setting up physical experiments in well-specified scenarios.

evaluates algorithms against social behavior through crowds of cylindrical pedestrians controlled by motion trajectories specified by polynomials of different orders and different numbers of pedestrians. It is customizable to different robot platforms, with a Jackal provided. Only SR (collision-free navigation reaching the goal) is provided as a metric, though the platform is extensible. DYNABARN provides several baselines including DWA [44], TEB [136], a behavior cloned [127] policy, and a TD3 [48] RL policy. While DYNABARN does not support human evaluation of robot behavior, it includes a demonstration pipeline to collect human teleoperation baselines of navigation in dynamic environments.

GYM-COLLISION-AVOIDANCE [39] is a downloadable, simulated benchmark used to evaluate multi-agent collision avoidance. Created to evaluate the GA3C-CADRL algorithm [39, 40] against the baselines ORCA [162], SA-CADRL [26], and DRLMACA [90], this benchmark provides a variety of multi-agent scenarios involving cylinders in simplified synthetic environments and measures SR, collisions, stuck and time-to-goal metrics. However, it focuses on policy-controlled agents interacting with each other and does not support human evaluation of robot behavior.

HU NAVSIM [122] is a downloadable, simulated benchmark focused on improving the development of social navigation systems around a variety of human behaviors. HU NAVSIM combines

Behavior Trees (BT) [33] and the SFM [62] to provide a variety of human behaviors ranging from indifferent, surprised, curious, fearful and aggressive. HuNAVSim is implemented as a framework that can work with various simulators and provides a plugin to work with ROS2 and Gazebo. HuNAVSim provides a variety of metrics comparable to those used in the SEAN simulator [161] and other benchmarks but does not provide baseline policies or a way to evaluate robot behavior with human ratings.

SocNAVbench [12] is a downloadable, simulated benchmark used to evaluate social navigation algorithms against pre-recorded episodes of human pedestrian behavior drawn from the UCY [82] and ETH [119] datasets. SocNAVbench provides visually realistic pedestrians and environments, as well as baselines based on the SFM [62], ORCA [162], and SA-CADRL [26] as well as a pedestrian-unaware policy. SocNAVbench provides a wide variety of metrics in areas such as path quality, motion quality, robot–pedestrian interaction, and episode statistics. However, SocNAVbench’s purpose is to automatically generate scores, so it makes the design decision to focus on automatically generated metrics that approximate human ratings instead.

8.2.3 Social Navigation Challenges. The CROWDBOT [45] Challenge is an effort to develop a benchmarking platform for social robot navigation in dense crowds. CROWDBOT supports four different robot morphologies interacting with simulated crowds of walking humans controlled by a flexible framework called UMANS [163], with several crowd setups provided in the initial benchmark. CROWDBOT is a downloadable, simulated challenge;⁶ initial phases were held in 2020 and 2021 but a full public challenge has not yet been held.

The iGIBSON Challenge at the CVPR 2021 Embodied AI Workshop⁷ is a social navigation benchmark based on the eponymous iGIBSON [83, 145] simulation environment for navigation and manipulation tasks in household scenes. In this benchmark challenge, robots must navigate to targets without collision among pedestrians [120], which are simulated via the ORCA model [163] in fifteen interactive indoor household scenes. Evaluation metrics include Success weighted by Time Length for reaching the goal quickly, and PSC for maintaining a comfortable distance from all pedestrians. This benchmark enabled quantitative comparison of approaches from over a dozen teams, including methods based on techniques like DD-PPO [168], PPO [143], SAC [56], and so on, providing a clear picture of which algorithms were superior for the task. iGIBSON is a downloadable, simulated challenge, but it does not include human ratings, and in 2021 did not include on-robot tests.

The SEANAVBENCH Challenge is a social navigation benchmark created for the SEANavBench workshop⁸ held at ICRA’22. SEANavBench combines SocNavBench [12] within the SEAN 2.0 [161] simulator which enables social navigation algorithms to run on simulated robots via ROS in environments rendered in the Unity game engine. Social navigation algorithms can be evaluated in simulated environments across a variety of environment sizes, crowd densities, and pedestrian behavior, including simulated pedestrians and replay of pedestrian datasets. This enables the analysis of how algorithms can succeed or fail as environmental conditions change and the measurement of performance using a variety of metrics. SEANAVBENCH is a simulated benchmark to which users can upload their code and compare performance against other submissions and baselines. While the public version of the challenge did not use human ratings, SEANAVBENCH uses SEAN-EP [159] to run the SEAN 2.0 simulation environment on the web, which could be used to collect human feedback.

⁶<https://gitlab.inria.fr/CrowdBot/CrowdBotUnity/-/tree/master>

⁷<https://svl.stanford.edu/igibson/challenge2021.html>

⁸<https://seanavbench.interactive-machines.com/>

8.3 Strengths and Limitations of Existing Benchmarks

As we can see from Table 4, social navigation benchmarks support a variety of scopes, from dynamic obstacle avoidance to HRIs to navigation through crowds. All attempt to address features of social behavior and many of them are downloadable, simulated benchmarks that can be efficiently deployed and which provide metrics for evaluation and sometimes baselines for comparison.

Broadly speaking, however, different types of benchmarks have characteristic limitations: (a) scalable benchmarks tend not to ground their evaluations in human data, (b) benchmarks that use human data tend to need manual setup or additional components, (c) protocols for designing experiments focus only on human evaluations, and (d) few benchmarks have meaningful coverage of edge cases of navigation behavior.

We believe these limitations are resolvable, and next present our recommendations for how good benchmarks should be designed and outline steps the community could take to improve existing benchmarks.

8.4 Properties of a Good Social Navigation Benchmark

Existing social navigation benchmarks have many purposes, from testing in large crowds, smaller social scenarios, algorithm improvements, and even tests of benchmark fidelity themselves. However, for the results of one benchmark to be useful to the rest of the community, it is important to have a common language for benchmarking and to have a shared understanding of what it is that a benchmark tests.

To ensure that social navigation benchmarks evaluate approaches for social navigation in a way that communicates their results broadly in the social navigation community, we argue that benchmarks themselves should be evaluated against a set of commonly agreed-upon criteria.

Based on how benchmarks are used in the field and what results they need to communicate, we recommend that benchmarks (1) evaluate social behavior, (2) include quantitative metrics, (3) provide baselines for comparison, (4) be efficient, repeatable, and scalable, (5) ground human evaluations in human data, and (6) use well-validated evaluation instruments. Next, we unpack these criteria and explain how they should guide the development and usage of benchmarks.

- (1) *Guideline B1: Evaluate Social Behavior:* A good social benchmark should evaluate the properties of algorithms in social scenarios which involve humans and robots interacting. Therefore, a social benchmark should have metrics related to social behavior and not just contain pure navigation metrics such as SPL [2] or pure task metrics such as SRs.
- (2) *Guideline B2: Include Quantitative Metrics:* The benchmark should provide a breadth of quantitative metrics, enabling researchers with different goals to use the benchmark to evaluate their algorithms with respect to their task and context and to compare to other approaches in the literature; common metrics can be found in Section 6. Quantitative metrics are ideal to enable comparisons between approaches; including those which are objectively measurable (e.g., PSC [164]) and those assessed with validated instruments (such as Likert scale evaluation with validated questions). Benchmark metrics should measure not just socially relevant concerns but also traditional navigation performance, such as task success, speed of performance, safety, and proximity to humans.
- (3) *Guideline B3: Provide Baselines for Comparison:* At a minimum, it is recommended to have baseline policies that show worst-case performance (e.g., a straight line planner that stops at obstacles) to serve as a lower bound for the benchmark. An upper bound oracle performance (e.g., demonstrations from a human, or an appropriate state-of-the-art algorithm) can also be provided if feasible. Ideally, if a state-of-the-art approach exists, it should be compared, but it is not always feasible to include these in a given benchmark due to availability or cost.

- (4) *Guideline B4: Be Efficient, Repeatable, and Scalable:* To democratize benchmarks and promote productive competition and collaboration among different scientists, efficient, repeatable, and scalable benchmarks are preferable. For example, the cost to run the benchmark should not be prohibitively expensive. While some benchmarks explicitly seek to reveal unique in-the-wild variations, the benchmark should nevertheless be repeatable such that it can be repeated multiple times with comparable results when scaled to a large number of trials. A good rule of thumb is at least 30 samples for real robot trials, but this number can be determined in a more principled statistical way from data if means and variances are available.
- (5) *Guideline B5: Ground Human Evaluations in Human Data:* At this point, many researchers agree that we do not have a good enough model of how humans react to robots to predict how they will react from other observables. Therefore, many researchers propose benchmarks should measure socialness based on human evaluations. An alternative approach is to use a learned model to predict human perception of the socialness of robot behaviors using a dataset of labeled examples; some researchers argue this provides a more validated metric than an *ad hoc* social score; other researchers argue the context that makes these learned metrics can be lost if used in other scenarios. Nonetheless, learned metrics could offer repeatable and scalable approximations of human responses, which could be evaluated via user studies.
- (6) *Guideline B6: Use Well-Validated Evaluation Instruments:* Ideally, human questionnaires should be standardized or empirically validated and should be ecologically valid for the task at hand; validating metrics is an iterative process which involves proposing metrics, conducting studies, statistically analyzing responses, and exposing metrics to peer review in the community. Objective metrics should also be empirically validated to ensure they measure what they purport to measure.

To address the shortcomings of existing benchmarks against these criteria, we recommend the following:

- (1) *Promote More Human Evaluation:* Many benchmarks use proxies of human ratings; while this is reasonable to enable fast evaluations, the community should encourage benchmark developers to collect human ratings and should push for broader adoption of rating pipelines such as SEAN-EP [159] to facilitate this collection.
- (2) *Standardize Social Questionnaires:* While it is useful to have well-defined scenarios as in the SOCIAL NAVIGATION PROTOCOL, the improvements to the questionnaires made by subsequent work in this area should be standardized and made available to inform labeling pipelines.
- (3) *Standardize Quantitative Metrics:* While some existing benchmarks and protocols specify minimum quantitative metrics, SOCNAV BENCH, SEANAV BENCH, and HUNAVSIM are converging on metrics similar to CROWDBOT's metrics; the community should encourage adopting a minimum set of these metrics.
- (4) *Test Corner Cases on Standard Benchmarks:* While social metrics are important, ensuring safe, reliable navigation performance is also important. Navigation benchmarks such as BARN [124] or BENCH-MR [61] should be used to validate traditional navigation behaviors.

Finally, it is worth noting that there are additional multi-agent benchmarks focused on gridworlds such as ASPRILO⁹ for logistics and MAPF¹⁰ for multi-agent pathfinding which we did not discuss as they do not focus on aspects of social behavior; however, as social navigation approaches become integrated into multi-agent or logistically complex domains, features from these benchmarks may also be useful for testing corner cases.

⁹<https://asprilo.github.io/>

¹⁰<https://movingai.com/benchmarks/mapf.html>

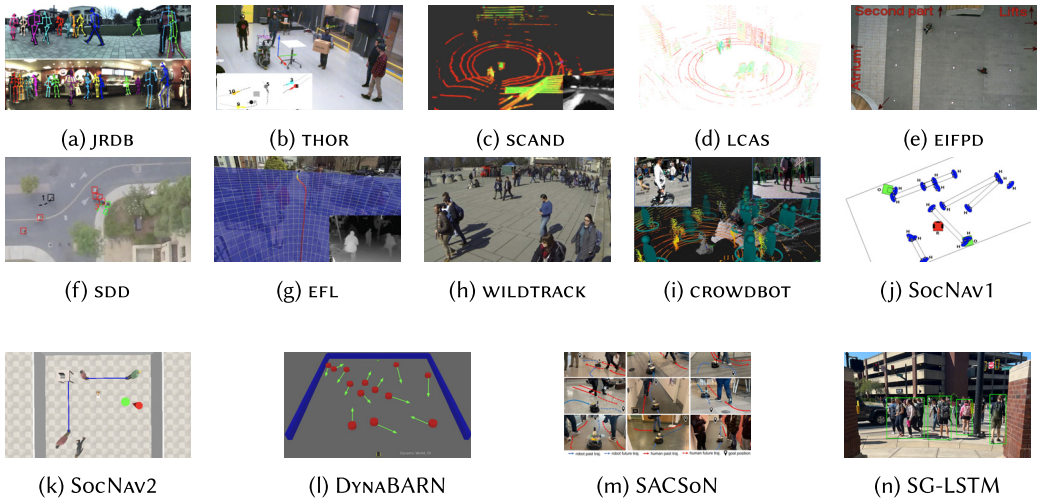


Fig. 9. Illustration of various social navigation datasets. See the text and Tables 5 and 6 for details.

9 Social Navigation Datasets

In this section, we provide a deeper look at datasets with regard to the factors listed in Section 5. First, we review desired dataset characteristics, noting that analyzing datasets requires drilling deeper into factors such as robot hardware, sensors, and behavior authoring methods, as well as additional factors for analysis such as data collected, dataset coverage, sampling distribution, annotations, and privacy and fairness handling. Then we use these factors to analyze several datasets, including JRDB [96], THOR [138], TRAJNET++ [78], ETH/UCY [82, 119], **Edinburgh Informatics Forum Pedestrian Database (EIFPD)** [92], **Stanford Drone Dataset (SDD)** [134], **Egocentric Future Localization (EFL)** [118], WILDTRACK [24], SCAND [70, 71], MUSoHU [111], CROWDBOT [116], DYNABARN [109], SocNAV1 [94], SocNAV2 [4], SACSON [63], and SG-LSTM [10], reviewing them with respect to the criteria (Figure 9).

9.1 Expanding the Factors for Dataset Analysis

In addition to the factors listed in Section 5.1, additional aspects must be considered for datasets:

Robot Hardware Platform. As different robot morphologies might elicit different human responses, it might be of importance to consider a *larger set of robots to collect data with*. Further, it might be useful to utilize props, e.g., engaging face, human-like head, and eye appearance and movement, to elicit stronger engagement with humans.

Sensors. In addition to robot sensors, a good practice is to *record teleoperation commands*, e.g., joystick controls, together with the data.

Robot Behavior Authoring Methods. The core of a social navigation dataset is demonstrations of desired socially aware robot behaviors. How are these demonstrations defined (see Section 7 for a deeper discussion of this topic). Should the robot behave as a human or as a different social agent (see Section 3.1 for a deeper discussion on this topic). In the case of a dataset, some of the options are as follows:

- (1) Pedestrians/humans: If the definition of a social robot is to behave as a human, recordings of moving humans/pedestrians might suffice.

- (2) Teleoperators: If behavior is desired that might be different from human behaviors, then data can be collected via robot teleoperation. Hence, an important principle in creating a dataset is to have *explicit and clear instructions to teleoperators on how to control the robot*. These instructions should cover the following topics:
- Is the teleoperator visible to humans?
 - Where is the teleoperator positioned w.r.t. the robot?
 - Instructions should ideally guarantee that the teleoperator does not affect the HRI.
 - Utilize multiple teleoperators, especially for the same scenarios, to encourage diversity in the data.

Data Collected. When it comes to dataset creation one of the major questions is for what social scenarios does one collect data. Therefore, the guidelines in Section 7 apply here. Note that for datasets in the wild, there is limited ability to control the scenarios. On the one side, one can opt for a completely unconstrained collection in a given environment, e.g., building, city, or area. On the other side, one can target specific events/activities, e.g., busy areas around campus, campus cafeteria, boardwalk crowds.

An important guideline is to *define the scope of the dataset such that the available dataset resources (hours of collection) are sufficient to collect data that thoroughly explores this scope*. The scope should be broad enough to present interesting challenges for the community to study. Therefore, it is desirable to *make the dataset scope as broad as possible*.

At the same time, one needs sufficient data for the dataset to be useful. More concretely, each *scenario within the dataset scope should be well sampled in the dataset*. This can help ensure that methods developed on the dataset can be deployed in the real world within the scope of the dataset, as they are less likely to encounter out-of-distribution scenarios.

Annotations. A question specific to a dataset is the annotations generated after the data has been collected. When it comes to social navigation, there aren't existing taxonomies of human-robot or human-human interactions. Existing computer vision datasets and benchmarks for activity recognition can provide a good starting point, e.g., ActivityNet [17].

Another consideration is the granularity of annotation. When it comes to activities, one can annotate whole navigation episodes with global labels, or segments within these episodes. Similarly, for human tracks, one can annotate tracks only, tracks with bounding boxes, skeletal tracking and gaze, and so on.

Privacy and Fairness. As social navigation datasets contain humans, privacy is an important concern. Decisions must be made on whether to anonymize humans and how to comply with privacy protection regulations.

9.2 Existing Social Navigation Datasets

In this section, we review some of the existing datasets in the context of our social navigation characteristics. These are presented in Tables 5 and 6. We review the following datasets.

JRDB [96] is a multi-modal dataset containing stereo 360 RGB video, 3D lidar scans, audio, and wheel encoder measurements from both indoor and outdoor environments. It provides annotations for human tracking and detection along with a benchmark and metrics to compare different algorithms.

THOR [138] is a public dataset providing motion trajectories of robots and humans in a range of curated scenarios of humans visiting and inspecting areas or carrying objects.

SCAND [70, 71] is a public dataset providing socially compliant navigation demonstrations recorded via teleoperating two different mobile robots in a socially compliant manner by human

Table 5. Characteristics of Existing Social Navigation Datasets, Part 1

Benchmark Dataset	THOR	SCAND	ETH/UCY	Trajnet++	EIFPD	SocNav1	SocNav2	SDD
Dataset context and scope	Three pre-defined human-robot social role-play activities.	Socially compliant navigation in-the-wild	Human trajectories recorded from a bird's eye view vantage point	Human trajectories recorded from a bird's eye view vantage point	Human trajectories recorded from a bird's eye view vantage point	Scenarios with interactions labeled with a social score	Short sequences with interactions labeled with social and holistic scores	Human trajectories recorded from a bird's eye view vantage point using a drone
Environment	Curated indoor environment: two rooms with arranged furniture and motion caption system	Indoor and outdoor campus-scale environment	Outdoor, fixed environment	Indoor and outdoor, fixed environment	Outdoor, fixed environment	Indoor, abstract	Indoor, abstract	Outdoor environment, focusing on diverse social navigation scenarios
Data collected	60 minutes of motion tracking across 600 human trajectories	522 minutes, consisting of 138 trajectories of teleoperation data from four demonstrators	Bird's eye view frames, with annotated human trajectories across time on five scenes	>200K human trajectories across dozens type scenes with high-density crowds	Bird's eye view frames, w/ annotated human trajectories across time	9,280 static scenario descriptions with social scores	53,600 dynamic scenario descriptions with social and all-encompassing scores	Bird's eye view frames w/ annotations of pedestrians, bikes, cars, and so on in 100 scenes w/ annotations of social interactions
Scenarios	Visiting areas, carrying boxes, inspecting targets	Goal-oriented social navigation	Pedestrian navigation	Pedestrian navigation	Pedestrian navigation	Evaluating robot disturbance	Evaluating robot trajectories	Real-world navigation with social interactions
Robot platform	Linde CitiTruck robot (W 1.56 m x L 0.55x x H 1.17 m)	Boston dynamic spot, ClearPath Jackal	N/A	N/A	N/A	Turtlebot-sized robot	Turtlebot-sized robot	3DR solo drone
Robot behavior	Programmed to follow a pre-defined path in a socially unaware manner	Human teleoperation in a socially compliant manner	N/A	N/A	N/A	Static placement	Teleoperation and policy	N/A
Human behavior	Follow a pre-defined path, and solving tasks in presence of other humans	Demonstrators teleoperate robot in open environment with other humans	Open world navigation	Open world navigation	Open world navigation	Static placement	Randomized simulated trajectories	Performing navigation activities such as walking, driving, and biking in a socially compliant manner
Sensors	Stationary Velodyne 3D LiDAR, Qualis Oqus 7+ motion tracking system, Tobii Pro Glasses for gaze tracking	Velodyne 3D LiDAR, Azure RGB, Odometry, Joystick	Stationary RGB camera overlooking pedestrians	Stationary RGB camera overlooking pedestrians	Camera fixed overhead 23 meters from the floor	Overhead view of robots and humans	Robot and human poses, with 53,600 short videos	RGB camera
Tasks and metrics	N/A	N/A	Human trajectory prediction	Human trajectory prediction	Human trajectory prediction	Acceptability of robot disturbance of humans	Acceptability of robot movement around humans	Social activity recognition, planning and trajectory prediction

demonstrators. The objective behind the SCAND dataset is to study the social navigation behavior of robots in the presence of human crowds. Similar to SCAND, MuSoHu [111] includes 3D lidar scans, RGBD camera images, 360° camera images, IMU data, and ambient sound collected from a sensor suite mounted on a helmet worn by humans walking around public spaces (instead of on a teleoperated robot), from which social robot navigation can be learned. This allows social human navigation data to be collected in the wild with a low setup cost, making MuSoHu easily extendable.

Table 6. Characteristics of Existing Social Navigation Datasets, Part II

Benchmark Dataset	EFL	LCAS	WILD-TRACK	JRDB	CrowdBot	DynaBARN	SACSoN	SG-LSTM
Dataset context and scope	Human trajectories recorded from human perspective	Online human detection from 3D lidar scans	Multi-camera detection and tracking of moving humans	Dataset of social interactions in indoor and outdoor envs. for solving perceptual tasks	Outdoor pedestrian tracking around a personal mobility robot	Diverse set of moving agent scenarios	Autonomous policy interacting with humans.	Curated interaction, movements, and formation of pedestrian groups
Environment	Outdoor scenes such as parks, malls, and a university campus	Outdoor environment	Outdoor environment	Indoor and outdoor environment	Crowded outdoor scenes	Indoor, abstract	Indoor environment	Outdoor university campus environment
Data collected	RGBD frames recorded from an egocentric perspective	49 minutes of 3D lidar scans	Multi-camera synchronized frames at 10fps	60,000 annotated frames of humans, recorded from an egocentric robot view	250K frames / 200 minutes from an egocentric POV	Moving polynomial agent scenarios	75 hours of visual navigation with 4,000 HRIs collected from robot POV	Color and depth frames, pedestrian and group bounding boxes
Scenarios	Egocentric real-world navigation	Real-world navigation in crowded environments	Third-person view open-world pedestrian navigation in crowded environments	Navigation in a campus-scale crowded environment	Real-world navigation in a crowded outdoor environment	Multiple moving agents	Real-world navigation in a crowded indoor environment	Paths, green spaces, study spaces, cafes, gatherings, weather events
Robot platform	N/A	Pioneer 3-AT mobile robot	N/A	JackRabbit mobile manipulator	Qolo personal mobility robot	N/A	iRobot Roomba	GO1 Edu robot
Robot behavior	N/A	Human teleoperation	N/A	Teleoperated	Both shared-control and reactive control	N/A	Policy-controlled	Unobtrusive data collection
Human behavior	Socially compliant navigation in the open world	Open-world navigation	Open-world navigation	Open-world navigation	Open-world navigation	Navigation among moving obstacles	Open-world navigation	Many types of pedestrian groups
Sensors	GoPro Hero 3 stereo cameras with 100 mm baseline	Velodyne VLP-16 3D LiDAR	three GoPro Hero 4 and four GoPro Hero 3	RGBD, fisheye and 360 RGB cameras, Velodyne and Sick LiDARs, microphone, wheel encoders	Point clouds, RGBD, people trackers, pose, contact forces	N/A	Spherical RGBD, fisheye RGB, 2D LiDAR, odometry, bumper	RGB-D cameras
Task and metrics	Trajectory prediction	Online human detection and tracking	Trajectory prediction, person tracking, and re-identification	Benchmark and metrics for 2D and 3D Person detection and tracking	Benchmark for crowd navigation	Dataset of scenarios for benchmarking dynamic navigation	Learning dataset of autonomous policy interacting with humans	Group size, walking speeds, proximity, cohesiveness, interactions

Also like SCAND, LCAS [172] is a public dataset containing 3D lidar scans collected using a mobile robot teleoperated in heavily crowded environments. Unlike SCAND, the robot is not necessarily teleoperated in a socially compliant manner. The focus of LCAS is to solve perception-related challenges in social navigation, such as online human detection.

The ETH/UCY [82, 119] dataset consists of human trajectories recorded in public spaces from a bird's eye view vantage point using an RGB camera. The trajectories are extracted by tracking humans from the bird's eye view images. The motivation behind the ETH/UCY dataset is to provide real-world trajectories of humans navigating among other humans in the scene so one can replicate-by-copying such trajectories in a simulator. Trajectories from the ETH/UCY dataset can be used to simulate a diverse set of realistic social scenarios. Pellegrini et al. [119] propose conditioning the predicted future trajectory also on scene knowledge and social interactions among agents.

The TRAJNET++ [78] dataset is composed of several existing datasets such as ETH/UCY [82, 119], CFF crowd dataset [1] with other synthetic data generated with ORCA [162]. Kothari et al. [78] have shared a benchmark and challenge focusing on agent-agent scenarios. They provide a proper sampling of trajectories and a unified extensive evaluation system to test the gathered methods for a fair comparison.

The EIFPD [92] is again similar to ETH/UCY, while providing a much higher number of humans captured in the dataset, the camera is fixed overhead roughly about 23 meters from the floor. Humans are detected by processing this bird's eye view image from the scene and tracking them in the scene.

SDD [134] is similar to ETH/UCY since it also provides a bird's eye view frame, recorded using a drone (unlike ETH/UCY that uses a statically mounted camera). Compared to ETH/UCY, the unique selling point of this dataset is large-scale images and videos of diverse scenarios including bicyclists, skateboarders, cars, buses, and golf carts navigating in the real world.

EFL [118] provides RGBD sequences of frames from the perspective of a human in various indoor and outdoor scenes such as parks, malls, and a campus, with various activities such as walking, shopping, and social interaction. EFL's focus is human trajectory prediction in novel scenes.

WILDTRACK [24] is similar to ETH/UCY. A GoPro camera is mounted in an outdoor environment scene consisting of crowds of people walking around. This dataset focuses on person detection in the presence of severe obstacles such as other humans and static obstacles in the scene.

The CROWDBOT [116] consists of egocentric RGBD and point-cloud data from a Qolo robot [117], [54] captured in autonomous and teleoperated modes in outdoor scenes.

Several datasets present synthetic trajectories for benchmark comparison. SocNAV1 [94] and SocNAV2 [4] are datasets of human-labeled simulated HRIs used for both benchmarking algorithms and as training datasets for learning algorithms. DYNABARN [109] includes 300 synthetic environments with agents with different motion profiles.

The SACSoN [63] dataset is a collection of egocentric RGB, RGBD, LIDAR, odometry, and bumper data from a policy-controlled iRobot Roomba navigating autonomously under policy control in indoor human environments. The dataset was created by a scalable system wrapping the policy control with a help-and-rescue module enabling long-term data collection, resulting in 75 hours of data and 58 kilometers of interaction with over 4,000 individual HRIs. The dataset supported a continual-learning architecture which showed the ability to learn from collected data. Interestingly, the experimenters collected an "interaction-rich" subset of data in which the robot was encouraged to drive closer to humans—then negated this objective and used these data to train a socially compliant policy.

The SG-LSTM [10] dataset is a curated dataset designed to provide insight into pedestrian behavior collected on Purdue University's West Lafayette, Indiana campus using a GO1 Edu robot by Unitree Robotics navigating unobtrusively among pedestrians. SG-LSTM focuses on the interactions, movements, and group formations of pedestrians in a variety of scenarios including campus thoroughfares, gatherings, dining and study areas, green spaces, and inclement weather events.

9.3 Guidelines for Datasets

Guideline D1: Make Datasets as Broad as Possible. This will ensure the dataset is useful to the community and will ensure investment in the data collection is well spent.

Guideline D2: Scope Datasets Based on Resources. Ensure the available dataset resources are sufficient to collect data that thoroughly explore the dataset scope.

Guideline D3: Ensure Each Scenario Is Well-Sampled. This ensures that methods trained on the dataset do not encounter out-of-distribution scenarios and the dataset is representative.

Guideline D4: Use Robots If Robot Behavior Is Desired. While datasets of pedestrians are useful, if robots are expected to behave differently than people, recording actual robot behavior rather than just pedestrians is desirable.

Guideline D5: Use Diverse Robot Platforms: Different robot morphologies may elicit different human responses, so if feasible datasets should use more than one robot morphology.

Guideline D6: Record Behavior Generation Commands. In addition to normal robot sensors, teleoperation commands should be recorded if robots are human-driven, or policy actions should be recorded if the behavior is authored.

Guideline D7: Collect Annotations Systematically. While standards for social navigation annotation are still being developed, formalizing data collection and modeling it on existing benchmarks in other fields can help. Data should be well labeled: methods used for generating human and robot behavior and collecting labels should be specified.

Guideline D8: Consider Privacy Issues Early. The collection of data involving humans involves privacy, policy, legal and moral issues. Considering these issues early can ensure that the dataset does not face legal challenges.

10 Simulation-Based Evaluation

The fundamental requirement for a social navigation simulator is the ability to simulate two agents at one time in a social encounter—without that, it’s just traditional navigation. Beyond this core requirement, social navigation simulators span the gamut from supporting crowds of simplified agents that test dynamic navigation algorithms to simulators that recreate human appearances, footsteps, behavioral diversity, and environmental interactivity. Most benchmarks discussed in Section 8 rely on a simulator to make benchmarking efficient, repeatable, and scalable.

In this section, we expand the social navigation factors particular to simulators, review existing simulators including CROWDBOT, CROWDNAV, DYNABARN, GYM-COLLISION-AVOIDANCE, HUNAVSIM, IGIBSON, INHUS, IMHUS, MENGEROS, PEDSIMROS, SEAN 2.0, SOCIALGYM 2.0, and SOCNAV BENCH, analyze the properties of these simulators and how they may be improved. We then we attempt to find a common ground between simulators and benchmarks for social navigation by proposing a unified API in order to compute metrics along a single code path, including discussions of its high-level requirements, implementation of the high-level API, and implementation efforts in representative simulation environments. We conclude with guidelines for simulator usage and development.

10.1 Expanding the Factors for Simulator Analysis

In addition to the factors listed in Section 5.1, additional aspects must be considered for simulators, including:

Abstraction Level. Some social simulations model large-scale crowds and do not attempt to model humans or robots in detail. For our purposes here and in Table 7, we discuss only simulations that are at least capable of modeling individual HRIs.

Simulation Focus. Similar to the notion of context, social simulations can be targeted at large-scale crowd simulation, social navigation interaction between humans and robots, or more narrowly focused on dynamic obstacle avoidance.

Simulation Platform. Some social simulations are standalone codebases; others are built atop of existing simulators such as Gazebo or MORSE.

Agent Representation. Some simulations represent only one kind of interacting agent (generally, presumed to be all humans or all robots); others represent robots and pedestrians separately.

Table 7. Characteristics of Existing Social Navigation Simulations

Sim Name	Crowd-Bot	Crowd-Nav	Dyna-BARN	gym-collision-avoidance	Hu-Nav-Sim	iGibson	InHuS	IMHuS	Menge-ROS	Ped-Sim-ROS	SEAN 2.0	Social-Gym 2.0	Soc-Nav-Bench
Sim focus	Crowd sim	Crowd sim	Dyn. obstacles	Collision avoid.	Human sim	Social nav	Human sim	Human sim	Crowd sim	Crowd sim	Social nav	Social nav	Social nav
Sim platform	Gazebo	Crowd-Nav	Gazebo	gym-collision-avoidance	Gazebo	iGibson	MORSE and Stage	Gazebo	Menge-ROS	Gazebo	SEAN 2.0	Social-Gym 2.0	Soc-Nav-Bench
Agent repr.	Ped. and robot	Ped. and robot	Ped. and robot	Robots	Ped. and robot	Ped. and robot	Ped. and robot	Ped. and robot	Ped. and robot	Ped. and robot	Ped. and robot	Ped. and robot	Ped. and robot
Scene repr.	Model 3D	Geom. 2D	Geom. 2D	Geom. 2D	Model 3D	Scanned mesh	Geom. 2D	Geom. 2D	Geom. 2D	Model 3D	Model 3D	Geom. 2D	Scanned mesh
Scene fidelity	Realist. 3D	2D	Abstr. 3D	2D	Abstr. 3D	Abstr. 3D	Realist. 3D	Abstr. 3D	2D	Abstr. 3D	Realist. 3D	2D	Realist. 3D
Physical fidelity	Physics model	Kinematic	Kinematic	Kinematic	Kinematic	Force and mass	Kinematic	Kinematic	Kinematic	Physics model	Physics model	Kinodynamic	Kinematic
Robot fidelity	Robot dyn.	Disc	Robot shape	Disc	Robot dyn.	Robot shape	Robot shape	Robot shape	Disc	Robot shape	Robot dyn.	Kinodynamic	Robot shape
Ped. sim fidelity	UMANS	Move w/o gait	Polynomial	Policy-based	Gait and attitude	Move w/o gait	Gait and attitude	Gait and attitude	Move w/o gait	Gait and activity	Move w/o gait	Move w/o gait	Gait and pose
Ped. viz fidelity	Detail'd	Disc	Cylind.	Disc	Detail'd	Detail'd	Detail'd	Detail'd	Disc	Sensors	Detail'd	Polyg.	Detail'd
Ped. reaction	UMANS	ORCA	Policy-based	Policy-based	SFM and attitude	ORCA	React and attitude	ORCA, attitude	SFM and ORCA	SFM	SFM	SFM	Replay only
Sim interop	Unity interf.	Gym	ROS bag	Gym	ROS interf.	Gym	ROS interf.	ROS interf.	ROS interf.	ROS interf.	ROS and unity	Gym, ROS	ROS interf.

See Section 10.1 for details.

Scene Representation. Environmental assets for simulators include 2D geometry, modeled 3D geometry, and scanned meshes of real scenes; these scenes can represent abstract, indoor, or outdoor environments.

Scene Visual Fidelity. Some simulations are purely 2D; others use abstracted 3D representations; others attempt to render realistic 3D scenes with rich shaders.

Physics Simulation Fidelity. Some simulations only model the kinematics of moving agents in static environments; others model forces and object mass or kinodynamic constraints; others incorporate full physics models.

Robot Simulation Fidelity. Some simulations model robots as points or cylinders; others support detailed robot morphologies or even full robot simulation.

Pedestrian Simulation Fidelity. Some simulations model human movement as point movement controlled by a crowd algorithm; others model humans as 3D objects, and some model the human walking gait. Some add variability based on the human's personality or attitude.

Pedestrian Visual Fidelity. Pedestrians can be represented by 2D points, discs or polygons, 3D cylinders, basic human meshes that don't change shape, animated meshes with basic walking movements, or photorealistic agents. As photorealistic is subjective, we lump all human meshes into "detailed" for the purpose of Table 7.

Pedestrian Reactivity. Pedestrians can move on pre-recorded trajectories without reacting to other agents, or may react using a model such as the SFM [62] or ORCA [162]. Pedestrian behavior may be also modulated with individual attitudes, behavioral styles, or social activities specified by higher-level modules or using techniques such as BT [33].

Simulation Interoperability. Some simulators are standalone; others support the OpenAI Gym API [16] or have interfaces to integrate with environments such as ROS [129].

10.2 Existing Social Navigation Simulators

A variety of social navigation simulators have been used in the literature, from simple simulators designed to test individual algorithms to complex standalone simulators used in multiple contexts, shown in Figure 10 and described in Table 7. These simulators include:

The CROWDBOT [45] simulator supports four different robot morphologies interacting with simulated crowds of walking humans controlled by a flexible framework called UMANS [163] built on the Gazebo simulator [76].

CROWDNAV [25] is a 2D simulator for multi-agent scenarios using ORCA to orchestrate pedestrian discs around policy-controlled discs in simplified environments.

DYNABARN [109] is the simulator used in the DYNABARN benchmark. DYNABARN models crowds of pedestrians controlled by polynomial motion trajectories moving through simulated environments. Humans are represented by cylinders but robots are represented with full morphologies.

GYM-COLLISION-AVOIDANCE¹¹ is a 2D simulator for multi-agent scenarios using policy-controlled cylinders in simplified environments. Humans and robots are not distinguished.

The HUNAVSIM [122] benchmark contains a simulator using SFM and BT to provide a variety of human behaviors ranging from indifferent, surprised, curious, fearful, and aggressive. It can work with various simulators and represents both human gait and robot morphologies.

The IGIBSON [83, 145] simulation environment supports navigation and manipulation tasks in household scenes. Pedestrians are represented with moving mannequins controlled via ORCA [89, 163] but robots are represented with full morphologies, and objects in the environment can be moved.

INHUS [41] is a simulator for testing social navigation algorithms against a variety of human behaviors called attitudes. It provides a general interface to ROS simulators and is currently integrated with the MORSE and Stage simulators.

The InHuS system is extended to simulate multiple human agents with modulated behaviors. This new system, called IMHuS [60], uses ORCA for motion planning of agents and is built atop of Gazebo. The behaviors are modeled and controlled using a supervisor module.

MENGEROS [3] is a 2D simulation designed to support very large crowds. Robots are discs, but several pedestrian reactivity models are supported including SFM and ORCA. A ROS interface allows this to be used with a variety of systems.

PEDSIMROS¹² is a ROS package for pedestrian simulation based on SFM augmented with group behaviors and social activities. PedSimROS simulates behaviors in 2D, but can integrate with 3D simulators like Gazebo to incorporate physics models. Robot and pedestrian models are realistic enough for point-cloud sensors but pedestrians are visually simplified.

The SEAN 2.0 [159, 161] simulator enables social navigation algorithms to run on simulated robots via ROS in environments rendered in the Unity game engine; pedestrians are represented with full gaits and environments can be detailed.

SOCIALGYM 2.0 [65, 150] is a simulation supporting diverse robot types and human behaviors in a 2D simulation that respects kinodynamic constraints, built atop the PettingZoo [155] multi-agent reinforcement learning environment.

¹¹<https://github.com/mit-acl/gym-collision-avoidance>

¹²https://github.com/srl-freiburg/pedsim_ros

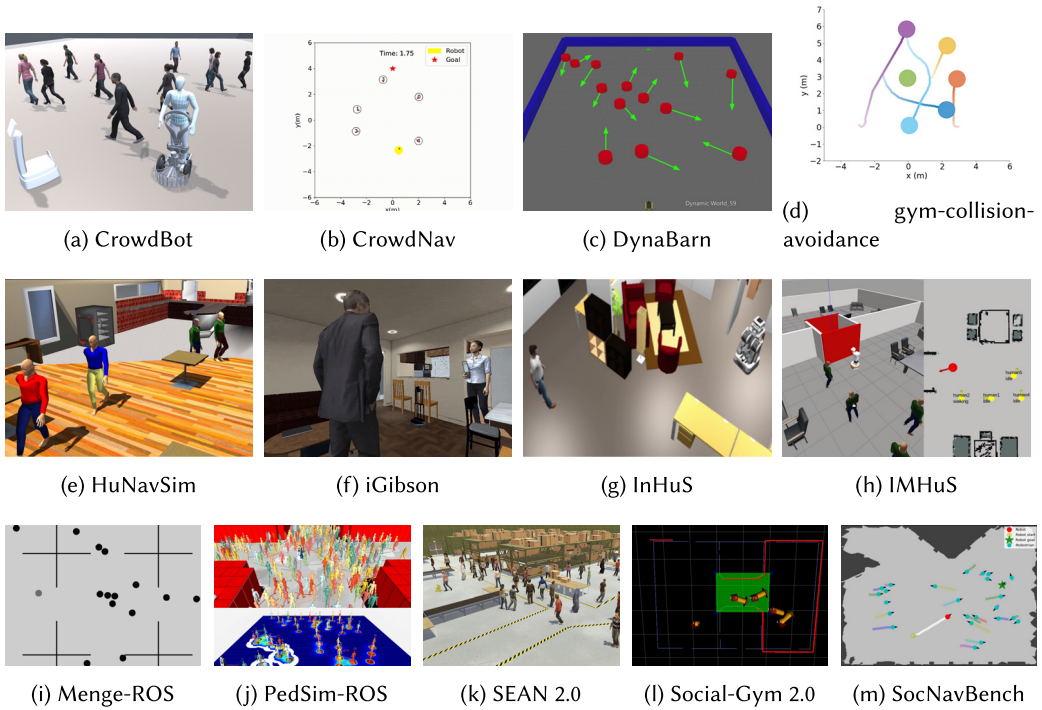


Fig. 10. Visual description of select simulators. See Section 10.2 for details.

The SOCNABENCH [12] benchmark contains a simulator to replay pre-recorded episodes of human pedestrian behavior drawn from the UCY [82] and ETH [119] datasets. SOCNABENCH provides visually realistic pedestrians and environments as well as robot morphologies.

10.3 Analysis of Simulation Platforms

10.3.1 Simulation Focus. Each simulation platform has been designed with a focus on a particular problem area. Example areas of focus include crowd simulation, how a robot should deal with dynamic obstacles or specific tasks such as social navigation or collision avoidance. Algorithms developed in different simulators may have a unique focus area as well, which implies we should be mindful when comparing algorithms across different simulators. For example, results from an algorithm trained in a simulator that uses a cylindrical representation of pedestrians may not be directly comparable to an algorithm that incorporates pedestrian gait.

We acknowledge the need for specialized simulators focusing on different problem areas. At the same time, we believe the community could benefit from a common social navigation simulator or a common API for multiple simulators. This common interface would provide access to a shared set of features that span focus areas. A common simulator or common API would enable training and evaluation across different approaches and promote the reuse of features from simulators that are focused on different areas.

10.3.2 Common Platforms. Many of the simulators listed in Table 7 have shared properties. For example, several simulators use Unity, ROS, or Gazebo as an underlying technology; simulators that use the same type of scene representations could share these representations; and methods of pedestrian reactivity could also be shared across simulators.

10.3.3 Pedestrian Reactivity. How pedestrians react to other nearby agents is an important factor to consider because the actions of simulated pedestrians directly influence both the training and evaluation of social navigation policies. Ideally, each simulated pedestrian would act in a manner identical to how a real-world pedestrian would act. Real-world pedestrian behavior can be observed by recording real-world pedestrian trajectories and playing them back in a simulator. However, there are two downsides to this approach. First, some fidelity of human motion is lost in the recording and playback process, for example, it is typical to capture motion only along the ground plane and not incorporate the body pose [82, 119]. Second, because the position of each pedestrian is determined by a pre-recorded trajectory, pedestrians cannot react to changes in the simulation. As soon as some element of the simulation deviates from the original data, such as the robot changing course, pedestrian motion is no longer realistic.

Motion models for pedestrians such as SFM [62] and ORCA [162] enable them to move in reaction to changes in the environment. While no model of human motion is perfect, the modeling of reactive agents in simulation allows researchers to explore how changes made to the environment by different robot policies affect task performance.

Pedestrian motion and reactivity play a critical role in the study of social navigation [25, 98]. The two imperfect solutions we have discussed indicate an opportunity for collaboration with the community to develop better alternatives.

10.3.4 Multi-Agent Policies. In the previous section, we explored various approaches to model pedestrian motion, including the use of SFM or ORCA, as well as playback of recorded trajectory data. However, in many real-world scenarios, the policies of agents are unknown and must be learned simultaneously. The field of robotics literature extensively covers navigation among dynamic obstacles, and there has been significant progress in multi-agent reinforcement learning [26], which has enabled the development of socially aware behavior in robots operating in constrained environments.

10.3.5 Environments. The scenario plays a crucial role in social navigation. Social navigation is not commonly observed in open environments; rather, it predominantly occurs in geometrically constrained or highly dense scenarios. Indoor spaces such as corridors, hallways, and dense areas like malls or airports are typical examples of such environments. These locations share similarities in terms of their physical characteristics. Thus, the simulators discussed thus far incorporate models that capture various aspects of such environments.

10.3.6 Metrics. Simulation can be a cost-effective alternative to the real world when training and evaluating robot control policies, which can in turn promotes scalability and reproducibility. The ability to compute metrics in a fair and comparable way, across robot control algorithms and simulators, is crucial to understanding the state of the field and making progress. Running trials and computing metrics under the same initial conditions in the real world is challenging. Simulation, however, allows the calculation of analytical metrics using ground-truth data, which is provided by the simulator, under common initial conditions when evaluating different algorithms. Moreover, learned metrics can be easily computed in a similar fashion and subjective metrics, which are based on human feedback, can be collected as well [94, 109, 160].

10.4 Toward a Unified API for Social Navigation Simulation

As discussed in Section 8, many benchmarks have been created using a variety of simulators to evaluate different aspects of social navigation. However, these benchmarks lack a unified standard for collecting metrics, making comparisons between benchmarks difficult and fragmenting the

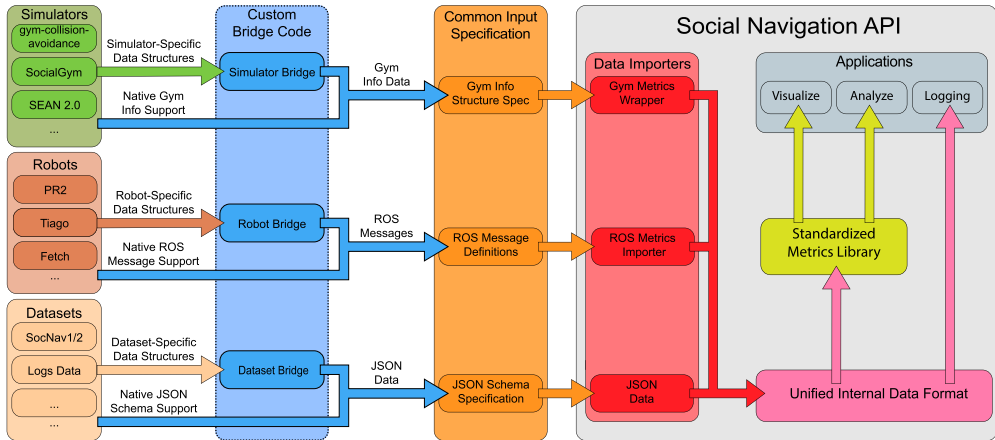


Fig. 11. Proposed social navigation simulation metrics API. A wide diversity of simulators and robot platforms exist, many of them supporting one or more platform APIs such as OpenAPI Gym or ROS. We propose to define a unified API that specifies the inputs needed to generate our recommended social navigation metrics, specified as either Gym observations or ROS messages. A unified metrics API with implementations for Gym and ROS will output a single output representation, enabling post-processing tools to generate visualizations, analytics, and logging with uniform code. To take advantage of these tools, simulator, and robot developers only need to contribute bridge code to output the required Gym or ROS data; dataset developers only need to output the single output representation.

community. While different benchmarks and simulators often have divergent emphases, nevertheless, we argue many common factors could be captured by a single high-level API, which would reduce fragmentation by easing comparisons.

Therefore, we propose a high-level API to calculate social navigation metrics that could be shared between simulators, enabling easier comparison of data collected from benchmarks built from those simulators. For broad adoption, we argue this simulation metrics API should be easy to use with a variety of simulators, real robots, and datasets, either natively or with easy-to-develop bridge code. To facilitate this interoperability, the API will specify both the data that it needs to compute metrics as well as implementations in the ROS and the OpenAI Gym API. These implementations will compute the metrics in Table 1 in a consistent way using common library code. While simulators often have very different code structures and philosophies, the proposed high-level API aims to help unify disparate efforts by defining a common set of data required to compute typical social navigation metrics, common library code, and a common data output format. This will make data from all API-compatible simulators and datasets available to use in shared analysis and visualization tools, which can be implemented in the future.

Figure 11 illustrates the flow of data through the proposed API. Next, we outline the API's design and our preliminary work on implementing it for common simulators.

10.4.1 Design of the High-Level Simulator Metric API. To enable the calculation of the metrics recommended in Table 1 from a variety of real robots, simulators, and datasets, the simulation metric API must specify its expected input data, including the robots under test, human pedestrians and other agents, and static and dynamic obstacles in the environment. The metrics API will enable the development of common downstream tools, but to make it broadly useful to the community it should also clearly its output format, as well as provide mechanisms for extensibility to support novel use cases as they develop.

- *Input Specification*: To compute the desired metrics, the API requires specific data from robots, simulators, or datasets. Specifying these data requires both the format needed for specific implementations, such as ROS message or OpenAI Gym info structures, as well as the content needed for metrics, including human pedestrians and other agents, the trajectories of robots under test, and static and dynamic obstacles in the environment.
 - *Pedestrian Data*: Simulators represent agents such as pedestrians or other robots in different ways. Typical data points include trajectories, teleoperation commands, current goals, and collective data about multiple agents, such as crowd flow. To compute many desired metrics, the proposed API requires at a minimum a pose for each agent over each timestep.
 - *Robot Data*: For the robot (or other agent, such as simulated pedestrian) under test the API needs not just pose but what the robots observed, what actions they performed, and what trajectories resulted—regardless of whether robots are guided by recorded trajectories, teleoperation by humans, or control policies.
 - *Obstacle Data*: The API needs information about the geometry of the physical environment to calculate certain metrics, such as collisions or the safety of an agent’s behavior; obstacle data includes static (wall geometry) and dynamic (doors, chairs) components.
- *Metric Computation*: The API will compute a variety of metrics listed in Table 1, including step-wise and task/episode level metrics as discussed in Section 6. Ideally, this metric computation should be done by standardized libraries so metrics are computed according to common definitions. This library for computing metrics should be extensible by the community, as different metrics are important to different researchers.
- *Output Specification*: The API should have a well-defined output specification so downstream tools can parse data from any system with a compatible format, facilitating the integration of datasets like those in Section 9, even if they cannot readily be replayed in simulators.
- *Downstream Tools*: This common output data format output will allow downstream tools to generate analyses and visualizations in a consistent way, as well as enable other data-driven applications to use data from API-compatible robots, simulators, and datasets. This could enable researchers to not only evaluate their systems in a common way but also analyze, visualize, and train data-driven systems on a variety of data from different sources with minimal feature engineering effort.

10.4.2 Implementation of Social Navigation API. To enable the broad usage of this API, we are developing an open-source implementation at <https://github.com/SocialNav/SocialNavAPI>. This reference implementation will include:

- (1) A JSON schema specification for the data input format, along with implementations that generate these data for ROS and for GymCollisionAvoidance simulator.
- (2) Reference implementations of the metrics in C++ and Python, packaged as libraries so different groups can reuse the same implementation to get comparable results.
- (3) A JSON schema for the output format, with examples generating output data for ROS and OpenAI Gym.

To use the proposed API, researchers must implement bridge code that translates data from their robots, simulators, or datasets into a format the API can consume. To make implementing bridge code easier, we will provide implementations for GymCollisionAvoidance and ROS which can be adapted for other systems, as shown in Figure 11. We are also working with the developers of SEAN 2.0, SocialGym, and DynaBarn to develop bridge code for these systems as well.

10.5 Guidelines for Simulators

Each simulator has its own purpose and scope, but, based on our analysis, we feel that a number of guidelines can be made for social navigation simulators which are intended to have broad use. First among these are guidelines which make it easier for simulators to interoperate:

Guideline S1: Use Standardized APIs. When possible, simulators should use standard APIs that enable approaches to be tested across different simulators.

Guideline S2: Support Standard Metrics. Simulators should provide quantitative metrics on a variety of dimensions of interest to enable different researchers to compare results—ideally, leveraging standard APIs so that metrics are computed in consistent ways, as suggested in Section 10.4.

Guideline S3: Support Extensibility. Regardless of the features a simulator supports, it is impossible to satisfy every use case. Novel research may require specific features that cannot be anticipated. Therefore, simulators should be designed with extensibility in mind, specifically enabling expert users to incorporate new functionality within the existing framework.

Next, we suggest guidelines to make simulators participate in the lifecycle of social navigation research:

Guideline S4: Support Dataset Generation. Simulators should make it easy to create datasets by systematically recording data from large-scale simulated runs.

Guideline S5: Support Benchmark Creation. Simulators should provide an API to create tasks and scenarios and to combine them with metrics and baselines to create a social navigation benchmark.

Guideline S6: Support Human Labeling. Simulators should make it easy to collect human labels of the acceptability or socialness of simulated episodes.

In addition, to support the increasing sophistication of social navigation scenarios and policies, we suggest guidelines for supporting increased visual and behavioral fidelity:

Guideline S7: Support Common Robot Morphologies. Simulators should provide instantiations of common robot morphologies to enable easy comparisons.

Guideline S8: Support Detailed Pedestrians. Where possible, simulators should support detailed pedestrian simulations to enable visual policies to react to walking pedestrian gaits. Ideally, this would extend to full visual realism of backgrounds as well, as well as replay of realistic pedestrians.

Guideline S9: Provide Options for Behavior Authoring. Simulators should provide ways to support behavior authoring, including playback of pedestrian recording, standard simulated models such as ORCA, and controls by custom policies. Supporting behavioral diversity in the generated policies is also important to capture the range of pedestrian behavior.

Finally, it is important to validate the simulation setup against its intended usage. Simulators should be periodically validated and refined to improve the realism and scope of the social navigation behaviors that they support.

11 Conclusions

Social robot navigation is critical to the success of mobile robots in human environments, but challenging because it combines all the problems of traditional robot navigation with the twin challenges of understanding how a robot can and should operate in concert with moving humans and understanding how humans react to this participation. In this article, we have outlined principles for social robot navigation and discussed guidelines for how these principles can be properly evaluated in scenarios, benchmarks, datasets, and simulators.

We defined a socially navigating robot as a robot that acts and interacts with humans or other robots, achieving its navigation goals while modifying its behavior to enable the other agents to better achieve theirs, and identified the key aspects needed to achieve this as safety, comfort,

legibility, politeness, social competency, understanding other agents, proactivity, and responding appropriately to context.

Building on this foundation, we reviewed the methodology of social navigation research and defined a taxonomy of factors used to describe social navigation metrics, scenarios, benchmarks, datasets, and simulators. Based on a review of existing work, we proposed a list of criteria for good benchmarking, including evaluating social behavior, including quantitative metrics, providing baselines for comparison, being efficient, repeatable, and scalable, using human evaluations on human data, and using well-validated evaluation instruments.

Figure 1 summarizes these guidelines to help researchers analyze their own research efforts and make good choices for benchmarking social robot navigation. We hope this framework for understanding social robot navigation will promote clearer benchmarking and faster progress in this field, and to promote this, we also proposed a common API for social navigation metrics to improve the ease of comparison.

Acknowledgment

An excerpt of this article, focusing on the HRI aspects of the social navigation problem, was published [47] in the AAAI 2023 Spring Symposium on HRI in Academia and Industry: Bridging the Gap.¹³

Author Contributions

Alexandre Alahi was a presenter at the symposium and contributed to the Definition working group. *Rachid Alami* was a presenter at the symposium, and contributed to the Metrics, Evaluation and Scenarios working groups. *Aniket Bera* contributed to the Definition working group. *Abhijat Biswas* was a presenter at the symposium and contributed to the Benchmarks working group. *Joydeep Biswas* contributed to the Datasets and Simulation working groups. *Rohan Chandra* contributed to the Simulations working group and the Definition section. *Hao-Tien Lewis Chiang* contributed to the Metrics and Benchmarks working groups. *Claudia Pérez-D'Arpino* co-led the direction and writing of the paper as co-DRI (directly responsible individual), and was an organizer, presenter, and moderator of the symposium, interviewed symposium participants and co-authors, collected their position papers and organized them into a draft, edited the manuscript, led the Evaluation working group, and assisted with the Simulator subgroup. *Michael Everett* was a presenter at the symposium, contributed a position paper, and contributed to the Simulators working group. *Anthony Francis* co-led the writing and direction of the paper as co-DRI (directly responsible individual), was an organizer and moderator of the symposium, interviewed symposium participants and co-authors, collected their position papers and organized them into a draft, helped draft many sections of the manuscript, edited the manuscript, led the Benchmarks working group and its subgroups, assisted with the Evaluation, Metrics and Simulation working groups, helped write and organize the Principles and Guidelines, helped organize, create, edit, and format the tables and figures, and created Figures 2–4, and 7. *Sehoon Ha* was a presenter at the symposium, contributed a position paper, and contributed to the Benchmarks working group. *Justin Hart* was a presenter at the symposium, contributed a position paper, and contributed to the Definition and Datasets working groups. *Jonathan P. How* was a presenter at the symposium, contributed a position paper, and helped edit the document. *Haresh Karnan* contributed to the Datasets working group. *Tsang-Wei Edward Lee* contributed to the Metrics and Evaluation working groups. *Chengshu Li* was an organizer, presenter and moderator for the symposium, helped define the process for the paper, helped collate the participant's comments, helped draft many of the sections of the document, led the Metrics working

¹³<https://sites.google.com/view/aaai-hri-bridge/home>

group, and assisted the Simulator working group. *Luis J. Manso* was a presenter at the symposium, contributed a position paper, and contributed to the Metrics working group. *Roberto Martín-Martín* was an organizer and presenter for the symposium, helped guide the direction for the paper, helped collate the participant position papers, helped draft many of the sections of the document, led the Definition working group, and assisted with the Datasets working group. *Reuth Mirsky* contributed a position paper and contributed to the Definition and Metrics working groups. *Sören Pirk* was a presenter at the symposium and contributed to the Metrics, Benchmarks and Evaluation working groups. *Phani Teja Singamaneni* was a presenter at the symposium and contributed to the Figures, Metrics, Simulation and Evaluation working groups. *Peter Stone* was a presenter at the symposium, contributed a position paper, and contributed to the Definition and Datasets working groups. *Ada V. Taylor* was a presenter at the symposium, contributed to the Metrics working group, Related Work, and editing across sections. *Alexander Toshev* was an initiator, organizer and presenter for the symposium, proposed and helped set the direction for the paper, helped collate the participant position papers, helped draft many of the sections of the document, led the Datasets working group, and assisted with the Metrics working group. *Peter Trautman* was a presenter at the symposium, advised on the paper, contributed a position paper, and contributed to the Datasets and Simulators working groups. *Nathan Tsoi* was a presenter at the symposium, contributed a position paper, and contributed to the Metrics, Simulation, and Benchmarks working groups. *Marynel Vázquez* was a presenter at the symposium, contributed a position paper, and contributed to the Metrics working group. *Fei Xia* was an organizer and moderator for the symposium, helped define the process for the paper, helped collate the participant position papers, helped draft many of the sections of the document, and led the Simulator working group. *Xuesu Xiao* was a presenter at the symposium, contributed a position paper, and contributed to the Benchmarks and Datasets working groups. *Peng Xu* contributed to the Benchmarks working group. *Naoki Yokoyama* was a presenter at the symposium, contributed a position paper, and contributed to the Benchmarks working group. In addition, the authors thank *Henny Admoni*, who presented at the symposium and contributed to the definitions of context for social navigation, but who was unable to participate in the drafting of the paper.

References

- [1] Alexandre Alahi, Vignesh Ramanathan, and Li Fei-Fei. 2014. Socially-aware large-scale crowd forecasting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2203–2210.
- [2] Peter Anderson, Angel Chang, Devendra Singh Chaplot, Alexey Dosovitskiy, Saurabh Gupta, Vladlen Koltun, Jana Kosecka, Jitendra Malik, Roozbeh Mottaghi, Manolis Savva, et al. 2018. On evaluation of embodied navigation agents. arXiv:1807.06757. Retrieved from <https://arxiv.org/abs/1807.06757>
- [3] Anoop Aroor, Susan L. Espstein, and Raj Korpan. 2017. Mengeros: A crowd simulation tool for autonomous robot navigation. In *2017 AAAI Fall Symposium Series*.
- [4] Pilar Bachiller, Daniel Rodriguez-Criado, Ronit R. Jorvekar, Pablo Bustos, Diego R. Faria, and Luis J. Manso. 2022. A graph neural network to model disruption in human-aware robot navigation. *Multimedia Tools and Applications* 81, 3 (2022), 3277–3295.
- [5] Kristen Backor, Saar Golde, and Norman Nie. 2007. Estimating survey fatigue in time use study. In *Proceedings of the International Association for Time Use Research Conference*. Citeseer.
- [6] Rishabh Baghel, Aditya Kapoor, Pilar Bachiller, Ronit R. Jorvekar, Daniel Rodriguez-Criado, and Luis J. Manso. 2021. A toolkit to generate social navigation datasets. In *Advances in Physical Agents II: Proceedings of the 21st International Workshop of Physical Agents (WAF '20)*. Springer, 180–193.
- [7] Aniket Bera, Sujeong Kim, Tanmay Randhavane, Srihari Pratapa, and Dinesh Manocha. 2016. GLMP-realtime pedestrian path prediction using global and local movement patterns. In *Proceedings of the 2016 IEEE International Conference on Robotics and Automation (ICRA '16)*. IEEE, 5528–5535.
- [8] Aniket Bera, Tanmay Randhavane, Emily Kubin, Austin Wang, Kurt Gray, and Dinesh Manocha. 2018. The socially invisible robot navigation in the social world using robot entitativity. In *Proceedings of the 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '18)*. IEEE, 4468–4475.

- [9] Aniket Bera, Tanmay Randhavane, Rohan Prinja, and Dinesh Manocha. 2017. Sociosense: Robot navigation amongst pedestrians with social and psychological constraints. In *Proceedings of the 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '17)*. IEEE, 7018–7025.
- [10] Rashmi Bhaskara, Maurice Chiu, and Aniket Bera. 2023. SG-LSTM: Social group LSTM for robot navigation through dense crowds. In *Proceedings of the 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '23)*. IEEE, 3835–3840.
- [11] Abhijat Biswas, Nathan Tsoi, Allan Wang, Peter Yu, Xiao He, Liyao Fu, Gustavo Silvera, Marynel Vazquez, and Aaron Steinfeld. 2022. SEANavBench @ ICRA 2023 workshop website. Retrieved January 1, 2023 from <https://seannavbench2022.netlify.app/benchmark/overview>
- [12] Abhijat Biswas, Allan Wang, Gustavo Silvera, Aaron Steinfeld, and Henny Admoni. 2022. Socnavbench: A grounded simulation testing framework for evaluating social navigation. *ACM Transactions on Human-Robot Interaction* 11, 3 (2022), 1–24.
- [13] Joydeep Biswas and Manuela Veloso. 2016. The 1,000-km challenge: Insights and quantitative and qualitative results. *IEEE Intelligent Systems* 31, 3 (2016), 86–96.
- [14] Paula Boddington. 2017. EPSRC principles of robotics: Commentary on safety, robots as products, and responsibility. *Connection Science* 29, 2 (2017), 170–176.
- [15] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 632–642. DOI : <https://doi.org/10.18653/v1/D15-1075>
- [16] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. 2016. Openai gym. arXiv:1606.01540. Retrieved from <https://arxiv.org/abs/1606.01540>
- [17] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. 2015. Activitynet: A large-scale video benchmark for Human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 961–970.
- [18] Enrico Cancelli, Tommaso Campari, Luciano Serafini, Angel X. Chang, and Lamberto Ballan. 2022. Exploiting socially-aware tasks for embodied social navigation. arXiv:2212.00767. Retrieved from <https://arxiv.org/abs/2212.00767>
- [19] Rohan Chandra, Aniket Bera, and Dinesh Manocha. 2021. Using graph-theoretic machine learning to predict human driver behavior. *IEEE Transactions on Intelligent Transportation Systems* 23, 3 (2021), 2572–2585.
- [20] Rohan Chandra, Rahul Maligi, Arya Anantula, and Joydeep Biswas. 2023. SOCIALMAPF: Optimal and efficient multi-agent path finding with strategic agents for social navigation. *IEEE Robotics and Automation Letters* 8, 6 (2023), 3214–3221.
- [21] Rohan Chandra and Dinesh Manocha. 2022. Gameplan: Game-theoretic multi-agent planning with human drivers at intersections, roundabouts, and merging. *IEEE Robotics and Automation Letters* 7, 2 (2022), 2676–2683.
- [22] Rohan Chandra, Mingyu Wang, Mac Schwager, and Dinesh Manocha. 2022. Game-theoretic planning for autonomous driving among risk-aware human drivers. In *Proceedings of the 2022 International Conference on Robotics and Automation (ICRA '22)*. IEEE, 2876–2883.
- [23] Konstantinos Charalampous, Ioannis Kostavelis, and Antonios Gasteratos. 2017. Recent trends in social aware robot navigation: A survey. *Robotics and Autonomous Systems* 93 (2017), 85–104.
- [24] Tatjana Chavdarova, Pierre Baqué, Stéphane Bouquet, Andrii Maksai, Cijo Jose, Timur Bagautdinov, Louis Lettry, Pascal Fua, Luc Van Gool, and François Fleuret. 2018. Wildtrack: A multi-camera HD dataset for dense unscripted pedestrian detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 5030–5039.
- [25] Changan Chen, Yuejiang Liu, Sven Kreiss, and Alexandre Alahi. 2019. Crowd-robot interaction: Crowd-aware robot navigation with attention-based deep reinforcement learning. In *Proceedings of the 2019 International Conference on Robotics and Automation (ICRA '19)*. IEEE, 6015–6022.
- [26] Yu Fan Chen, Michael Everett, Miao Liu, and Jonathan P. How. 2017. Socially aware motion planning with deep reinforcement learning. In *Proceedings of the 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '17)*. IEEE, 1343–1350.
- [27] Yu Fan Chen, Miao Liu, Michael Everett, and Jonathan P. How. 2017. Decentralized Non-communicating multiagent collision avoidance with deep reinforcement learning. In *Proceedings of the 2017 IEEE International Conference on Robotics and Automation (ICRA '17)*. IEEE, 285–292.
- [28] Jiyu Cheng, Hu Cheng, Max Q.-H. Meng, and Hong Zhang. 2018. Autonomous navigation by mobile robots in human environments: A survey. In *Proceedings of the 2018 IEEE International Conference on Robotics and Biomimetics (ROBIO '18)*. IEEE, 1981–1986.
- [29] Ernest Cheung, Aniket Bera, Emily Kubin, Kurt Gray, and Dinesh Manocha. 2018. Identifying driver behaviors using trajectory features for vehicle navigation. In *Proceedings of the 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '18)*. IEEE, 3445–3452.

- [30] Hao-Tien Lewis Chiang, Aleksandra Faust, Marek Fiser, and Anthony Francis. 2019. Learning navigation behaviors end-to-end with AutoRL. *IEEE Robotics and Automation Letters* 4, 2 (2019), 2007–2014.
- [31] S. F. Chik, C. F. Yeong, E. L. M. Su, T. Y. Lim, Y. Subramaniam, and P. J. H. Chin. 2016. A review of social-aware navigation frameworks for service robot in dynamic human environments. *Journal of Telecommunication, Electronic and Computer Engineering* 8, 11 (2016), 41–50.
- [32] Paul R. Cohen. 1995. *Empirical Methods for Artificial Intelligence*, Vol. 139, MIT Press Cambridge, MA.
- [33] Michele Colledanchise and Petter Ögren. 2018. *Behavior Trees in Robotics and AI: An Introduction*. CRC Press.
- [34] Catie Cuan, Edward Lee, Emre Fisher, Anthony Francis, Leila Takayama, Tingnan Zhang, Alexander Toshev, and Sören Pirk. 2022. Gesture2Path: Imitation learning for gesture-aware navigation. arXiv:2209.09375. Retrieved from <https://arxiv.org/abs/2209.09375>
- [35] Matt Deitke, Dhruv Batra, Yonatan Bisk, Tommaso Campari, Angel X. Chang, Devendra Singh Chaplot, Changan Chen, Claudia Pérez D’Arpino, Kiana Ehsani, Ali Farhadi, et al. 2022. Retrospectives on the embodied ai workshop. arXiv:2210.06849. Retrieved from <https://arxiv.org/abs/2210.06849>
- [36] Matt Deitke, Winson Han, Alvaro Herrasti, Aniruddha Kembhavi, Eric Kolve, Roozbeh Mottaghi, Jordi Salvador, Dustin Schwenk, Eli VanderBilt, Matthew Wallingford, et al. 2020. Robothor: An open simulation-to-real embodied AI platform. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 3164–3174.
- [37] Anca D. Dragan, Kenton C. T. Lee, and Siddhartha S. Srinivasa. 2013. Legibility and predictability of robot motion. In *Proceedings of the 2013 8th ACM/IEEE International Conference on Human-Robot Interaction*. IEEE, 301–308.
- [38] Daniel Dugas, Juan Nieto, Roland Siegwart, and Jen Jen Chung. 2021. Navrep: Unsupervised representations for reinforcement learning of robot navigation in dynamic human environments. In *Proceedings of the 2021 IEEE International Conference on Robotics and Automation (ICRA ’21)*. IEEE, 7829–7835.
- [39] Michael Everett, Yu Fan Chen, and Jonathan P. How. 2018. Motion planning among dynamic, decision-making agents with deep reinforcement learning. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS ’18)*. Retrieved from <https://arxiv.org/pdf/1805.01956.pdf>
- [40] Michael Everett, Yu Fan Chen, and Jonathan P. How. 2021. Collision avoidance in pedestrian-rich environments with deep reinforcement learning. *IEEE Access* 9 (2021), 10357–10377.
- [41] Anthony Favier, Phani Teja Singamaneni, and Rachid Alami. 2022. An intelligent human avatar to debug and challenge human-aware robot navigation systems. In *Proceedings of the 2022 17th ACM/IEEE International Conference on Human-Robot Interaction (HRI ’22)*. IEEE, 760–764.
- [42] Ronit Feingold-Polak, Oren Barzel, and Shelly Levy-Tzedek. 2021. A robot goes to rehab: A novel gamified system for long-term stroke rehabilitation using a socially assistive robot—Methodology and usability testing. *Journal of NeuroEngineering and Rehabilitation* 18, 1 (2021), 1–18.
- [43] Rolando Fernandez, Nathan John, Sean Kirmani, Justin Hart, Jivko Sinapov, and Peter Stone. 2018. Passive demonstrations of light-based robot signals for improved human interpretability. In *Proceedings of the 2018 27th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN ’18)*. IEEE, 234–239.
- [44] Dieter Fox, Wolfram Burgard, and Sebastian Thrun. 1997. The dynamic window approach to collision avoidance. *IEEE Robotics & Automation Magazine* 4, 1 (1997), 23–33.
- [45] INRIA France. 2020. *Safe Robot Navigation in Dense Crowds*. Technical Report. INRIA France.
- [46] Anthony Francis, Aleksandra Faust, Hao-Tien Lewis Chiang, Jasmine Hsu, J. Chase Kew, Marek Fiser, and Tsang-Wei Edward Lee. 2020. Long-range indoor navigation with PRM-Rl. *IEEE Transactions on Robotics* 36, 4 (2020), 1115–1134.
- [47] Anthony Francis, Claudia Pérez-D’Arpino, Chengshu Li, Fei Xia, Alexandre Alahi, Aniket Bera, Abhijat Biswas, Joydeep Biswas, Hao-Tien Lewis Chiang, Michael Everett, et al. 2023. Benchmarking social robot navigation across academia and industry. In *Proceedings of the AAAI 2023 Spring Symposium on HRI in Academia and Industry: Bridging the Gap*. AAAI.
- [48] Scott Fujimoto, Herke Hoof, and David Meger. 2018. Addressing function approximation error in actor-critic methods. In *Proceedings of the International Conference on Machine Learning*. PMLR, 1587–1596.
- [49] Mirta Galesic and Michael Bosnjak. 2009. Effects of questionnaire length on participation and indicators of response quality in a web survey. *Public Opinion Quarterly* 73, 2 (2009), 349–360.
- [50] Yuxiang Gao and Chien-Ming Huang. 2021. Evaluation of socially-aware robot navigation. *Frontiers in Robotics and AI* 8 (2021), 420.
- [51] Christian Gloor. 2016. Pedsim: Pedestrian crowd simulation. Retrieved from <https://web.archive.org/web/20200215011020/http://pedsim.silmaril.org/>
- [52] Rachel Gockley, Allison Bruce, Jodi Forlizzi, Marek Michalowski, Anne Mundell, Stephanie Rosenthal, Brennan Sellner, Reid Simmons, Kevin Snipes, Alan C. Schultz, et al. 2005. Designing robots for long-term social interaction. In *Proceedings of the 2005 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 1338–1343.
- [53] Google. 2023. Robotics @ Google website. Retrieved January 29, 2023 from <https://research.google/research-areas/robotics/>

- [54] Diego Felipe Paez Granados, Hideki Kadone, and Kenji Suzuki. 2018. Unpowered lower-body exoskeleton with torso lifting mechanism for supporting sit-to-stand transitions. In *Proceedings of the 2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '18)*, 2755–2761. DOI: <https://doi.org/10.1109/IROS.2018.8594199>
- [55] Ronja Guldenring, Michael Görner, Norman Hendrich, Niels Jul Jacobsen, and Jianwei Zhang. 2020. Learning local planners for human-aware navigation in indoor environments. In *Proceedings of the 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '20)*. IEEE, 6053–6060.
- [56] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the 35th International Conference on Machine Learning*. PMLR, 1861–1870.
- [57] Edward Twitchell Hall. 1966. *The Hidden Dimension*, Vol. 609, Garden City, NY: Doubleday.
- [58] Justin Hart, Elliot Hauser, Samuel Baker, Joydeep Biswas, Junfeng Jiao, and Luis Sentis. 2022. Longitudinal social impacts of HRI over long-term deployments. In *Proceedings of the 2022 17th ACM/IEEE International Conference on Human-Robot Interaction (HRI '22)*. IEEE, 1258–1260.
- [59] Justin Hart, Reuth Mirsky, Xuesu Xiao, Stone Tejada, Bonny Mahajan, Jamin Goo, Kathryn Baldauf, Sydney Owen, and Peter Stone. 2020. Using human-inspired signals to disambiguate navigational intentions. In *Proceedings of the 12th International Conference on Social Robotics (ICSR '20)*. Springer, 320–331.
- [60] Olivier Hauterville, Camino Fernández, Phani Teja Singamaneni, Anthony Favier, Vicente Matellán, and Rachid Alami. 2022. Interactive social agents simulation tool for designing choreographies for human-robot-interaction research. In *Proceedings of the 5th Iberian Robotics Conference on Advances in Robotics (ROBOT '22)*, Vol. 2, Springer, 514–527.
- [61] Eric Heiden, Luigi Palmieri, Leonard Bruns, Kai O. Arras, Gaurav S. Sukhatme, and Sven Koenig. 2021. Bench-MR: A motion planning benchmark for wheeled mobile robots. *IEEE Robotics and Automation Letters* 6, 3 (2021), 4536–4543. DOI: <https://doi.org/10.1109/LRA.2021.3068913>
- [62] Dirk Helbing and Peter Molnar. 1995. Social force model for pedestrian dynamics. *Physical Review E* 51, 5 (1995), 4282.
- [63] Noriaki Hirose, Dhruv Shah, Ajay Sridhar, and Sergey Levine. 2023. SACSoN: Scalable autonomous data collection for social navigation. Retrieved from <https://arxiv.org/abs/2306.01874>
- [64] Guy Hoffman and Xuan Zhao. 2020. A primer for conducting experiments in human–robot interaction. *ACM Transactions on Human–Robot Interaction* 10, 1 (2020), 1–31.
- [65] Jarrett Holtz and Joydeep Biswas. 2022. SOCIALGYM: A framework for benchmarking social robot navigation. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '22)*. IEEE, 11246–11252. DOI: <https://doi.org/10.1109/IROS47612.2022.9982021>
- [66] Ohad Inbar and Joachim Meyer. 2019. Politeness counts: Perceptions of peacekeeping robots. *IEEE Transactions on Human-Machine Systems* 49, 3 (2019), 232–240.
- [67] Hiroko Kamide, Yasushi Mae, Koji Kawabe, Satoshi Shigemi, Masato Hirose, and Tatsuo Arai. 2012. New measurement of psychological safety for humanoid. In *Proceedings of the 2012 7th ACM/IEEE International Conference on Human-Robot Interaction (HRI '12)*. IEEE, 49–56.
- [68] Takayuki Kanda and Hiroshi Ishiguro. 2017. *Human-Robot Interaction in Social Robotics*. CRC Press.
- [69] Takayuki Kanda, Hiroshi Ishiguro, and Toru Ishida. 2001. Psychological analysis on human-robot interaction. In *Proceedings of the 2001 IEEE International Conference on Robotics and Automation*, Vol. 4, IEEE, 4166–4173.
- [70] Haresh Karnan, Anirudh Nair, Xuesu Xiao, Garrett Warnell, Soeren Pirk, Alexander Toshev, Justin Hart, Joydeep Biswas, and Peter Stone. 2022. Socially compliant navigation dataset (SCAND). Texas Data Repository. DOI: <https://doi.org/10.18738/T8/0PRYRH>
- [71] Haresh Karnan, Anirudh Nair, Xuesu Xiao, Garrett Warnell, Sören Pirk, Alexander Toshev, Justin Hart, Joydeep Biswas, and Peter Stone. 2022. Socially Compliant Navigation Dataset (SCAND): A large-scale dataset of demonstrations for social navigation. *IEEE Robotics and Automation Letters* 7, 4 (2022), 11807–11814.
- [72] Linh Kästner, Teham Bhuiyan, Tuan Anh Le, Elias Treis, Johannes Cox, Boris Meinardus, Jacek Kmiecik, Reyk Carstens, Duc Pichel, Bassel Fatloun, et al. 2022. Arena-Bench: A benchmarking suite for obstacle avoidance approaches in highly dynamic environments. *IEEE Robotics and Automation Letters* 7, 4 (2022), 9477–9484.
- [73] Harmish Khambhaita and Rachid Alami. 2020. Viewing robot navigation in human environment as a cooperative activity. In *Robotics Research: The 18th International Symposium ISRR (2020)*. Springer, 285–300.
- [74] John F. Kihlstrom. 2021. Ecological validity and “ecological validity”. *Perspectives on Psychological Science* 16, 2 (2021), 466–471.
- [75] Aelee Kim, Jooyun Han, Younbo Jung, and Kwanmin Lee. 2013. The effects of familiarity and robot gesture on user acceptance of information. In *Proceedings of the 2013 8th ACM/IEEE International Conference on Human-Robot Interaction (HRI '13)*. IEEE, 159–160.

- [76] Nathan Koenig and Andrew Howard. 2004. Design and use paradigms for gazebo, an open-source multi-robot simulator. In *Proceedings of the 2004 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '04)*, Vol. 3, IEEE, 2149–2154.
- [77] Hatice Kose-Bagci, Kerstin Dautenhahn, and Chrystopher L. Nehaniv. 2008. Emergent dynamics of turn-taking interaction in drumming games with a humanoid robot. In *Proceedings of the 17th IEEE International Symposium on Robot and Human Interactive Communication*. IEEE, 346–353.
- [78] Parth Kothari, Sven Kreiss, and Alexandre Alahi. 2021. Human trajectory forecasting: A deep learning perspective. *IEEE Transactions on Intelligent Transportation Systems* 23, 7 (2021), 7386–7400.
- [79] Thibault Kruse, Amit Kumar Pandey, Rachid Alami, and Alexandra Kirsch. 2013. Human-aware robot navigation: A survey. *Robotics and Autonomous Systems* 61, 12 (2013), 1726–1743.
- [80] Przemyslaw A. Lasota, Terrence Fong, and Julie A. Shah. 2017. A survey of methods for safe human-robot interaction. *Foundations and Trends in Robotics* 5, 4 (2017), 261–349.
- [81] Iolanda Leite, Carlos Martinho, and Ana Paiva. 2013. Social robots for long-term interaction: A survey. *International Journal of Social Robotics* 5 (2013), 291–308.
- [82] Alon Lerner, Yiorgos Chrysanthou, and Dani Lischinski. 2007. Crowds by example. In *Computer Graphics Forum*, Vol. 26, Wiley Online Library, 655–664.
- [83] Chengshu Li, Fei Xia, Roberto Martín-Martín, Michael Lingelbach, Sanjana Srivastava, Bokui Shen, Kent Vainio, Cem Gokmen, Gokul Dharan, Tanish Jain, et al. 2021. Igbison 2.0: Object-centric simulation for robot learning of everyday household tasks. arXiv:2108.03272. Retrieved from <https://arxiv.org/abs/2108.03272>
- [84] Xuan Li, Yonglin Tian, Peijun Ye, Haibin Duan, and Fei-Yue Wang. 2022. A novel scenarios engineering methodology for foundation models in metaverse. *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 53, 4 (2022), 2148–2159.
- [85] Xuan Li and Fei-Yue Wang. 2023. Scenarios engineering: Enabling trustworthy and effective AI for autonomous vehicles. *IEEE Transactions on Intelligent Vehicles* 8, 5 (2023), 3205–3210. DOI : <https://doi.org/10.1109/TIV.2023.3269421>
- [86] Xuan Li, Peijun Ye, Juanjuan Li, Zhongmin Liu, Longbing Cao, and Fei-Yue Wang. 2022. From features engineering to scenarios engineering for trustworthy AI: I&I, C&C, and V&V. *IEEE Intelligent Systems* 37, 4 (2022), 18–26.
- [87] Rensis Likert. 1932. A technique for the measurement of attitudes. *Archives of Psychology* 22 (1932), 5–52.
- [88] Mark S. Litwin and Arlene Fink. 1995. *How to Measure Survey Reliability and Validity*, Vol. 7, Sage.
- [89] Yige Liu, Siyun Li, Chengshu Li, Claudia Perez-D’Arpino, and Silvio Savarese. 2021. Interactive pedestrian simulation in iGibson. In *RSS Workshop on Social Robot Navigation*.
- [90] Pinxin Long, Tingxiang Fan, Xinyi Liao, Wenxi Liu, Hao Zhang, and Jia Pan. 2018. Towards optimally decentralized multi-robot collision avoidance via deep reinforcement learning. In *Proceedings of the 2018 IEEE International Conference on Robotics and Automation (ICRA '18)*. IEEE, 6252–6259.
- [91] David T. Lykken. 1968. Statistical significance in psychological research. *Psychological Bulletin* 70, 3p1 (1968), 151.
- [92] Barbara Majecka. 2009. *Statistical Models of Pedestrian Behaviour in the Forum*. Master’s thesis, School of Informatics, University of Edinburgh.
- [93] Swathi Mannem, William Macke, Peter Stone, and Reuth Mirsky. 2023. Exploring the cost of interruptions in human-robot teaming. In *Proceedings of the 2023 IEEE-RAS 22nd International Conference on Humanoid Robots (Humanoids)*. IEEE, 1–8.
- [94] Luis J. Manso, Pedro Nuñez, Luis V. Calderita, Diego R. Faria, and Pilar Bachiller. 2020. Socnav1: A dataset to benchmark and Learn social navigation conventions. *Data* 5, 1 (2020), 7.
- [95] Eitan Marder-Eppstein, Eric Berger, Tully Foote, Brian Gerkey, and Kurt Konolige. 2010. The office Marathon: Robust navigation in an indoor office environment. In *Proceedings of the 2010 IEEE International Conference on Robotics and Automation*. IEEE, 300–307.
- [96] Roberto Martín-Martín, Hamid Rezatofighi, Abhijeet Sheno, Mihir Patel, J. Gwak, Nathan Dass, Alan Federman, Patrick Goebel, and Silvio Savarese. 2019. Jrdb: A dataset and benchmark for visual perception for navigation in human environments. arXiv:1910.11792. Retrieved from <https://arxiv.org/abs/1910.11792>
- [97] Christoforos Mavrogiannis, Francesca Baldini, Allan Wang, Dapeng Zhao, Pete Trautman, Aaron Steinfeld, and Jean Oh. 2021. Core challenges of social robot navigation: A survey. arXiv:2103.05668. Retrieved from <https://arxiv.org/abs/2103.05668>
- [98] Christoforos Mavrogiannis, Alena M. Hutchinson, John Macdonald, Patrícia Alves-Oliveira, and Ross A. Knepper. 2019. Effects of distinct robot navigation strategies on human behavior in a crowded environment. In *Proceedings of the 2019 14th ACM/IEEE International Conference on Human-Robot Interaction (HRI '19)*. IEEE, 421–430.
- [99] Ross Mead and Maja J. Matarić. 2017. Autonomous human–robot proxemics: Socially aware navigation based on interaction potential. *Autonomous Robots* 41, 5 (2017), 1189–1201.
- [100] Richard Meyes, Melanie Lu, Constantin Waubert de Puiseau, and Tobias Meisen. 2019. Ablation studies in artificial neural networks. arXiv:1901.08644. Retrieved from <https://arxiv.org/abs/1901.08644>

- [101] Marvin Minsky. 2007. *The Emotion Machine: Commonsense Thinking, Artificial Intelligence, and the Future of the Human Mind*. Simon and Schuster.
- [102] Reuth Mirsky and Einav Shpiro. 2024. Recognition and identification of intentional blocking in social navigation. In *Proceedings of the 2024 International Symposium on Technological Advances in Human-Robot Interaction*, 101–110.
- [103] Reuth Mirsky, Xuesu Xiao, Justin Hart, and Peter Stone. 2024. Conflict avoidance in social navigation—A survey. *ACM Transactions on Human-Robot Interaction* 13, 1 (2024), 1–36.
- [104] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. 2019. Model cards for model reporting. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 220–229.
- [105] Ronja Möller, Antonino Furnari, Sebastiano Battiato, Aki Härmä, and Giovanni Maria Farinella. 2021. A survey on human-aware robot navigation. *Robotics and Autonomous Systems* 145 (2021), 103837.
- [106] Beth Morling. 2014. *Research Methods in Psychology: Evaluating a World of Information*. WW Norton & Company.
- [107] Mehdi Moussaïd, Dirk Helbing, and Guy Theraulaz. 2011. How simple rules determine pedestrian behavior and crowd disasters. *Proceedings of the National Academy of Sciences* 108, 17 (2011), 6884–6888.
- [108] Mehdi Moussaïd, Niriaska Perozo, Simon Garnier, Dirk Helbing, and Guy Theraulaz. 2010. The walking behaviour of pedestrian social groups and its impact on crowd dynamics. *PLoS One* 5, 4 (2010), e10047.
- [109] Anirudh Nair, Fulin Jiang, Kang Hou, Zifan Xu, Shuoze Li, Xuesu Xiao, and Peter Stone. 2022. DynaBARN: Benchmarking metric ground navigation in dynamic environments. In *Proceedings of the 2022 IEEE International Symposium on Safety, Security, and Rescue Robotics (SSRR '22)*. IEEE, 347–352.
- [110] Venkatraman Narayanan, Bala Murali Manoghar, Vishnu Sashank Dorbala, Dinesh Manocha, and Aniket Bera. 2020. Proxemo: Gait-based emotion learning and multi-view proxemic fusion for socially-aware robot navigation. In *Proceedings of the 2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '20)*. IEEE, 8200–8207.
- [111] Duc M. Nguyen, Mohammad Nazeri, Amirreza Payandeh, Aniket Datar, and Xuesu Xiao. 2023. Toward human-Like social robot navigation: A large-scale, multi-modal, Social Human Navigation Dataset. arXiv:2303.14880. Retrieved from <https://arxiv.org/abs/2303.14880>
- [112] Aastha Nigam and Laurel D. Riek. 2015. Social context perception for mobile robots. In *Proceedings of the 2015 IEEE/RSJ International Conference on intelligent robots and systems (IROS '15)*. IEEE, 3621–3627.
- [113] Billy Okal and Kai O. Arras. 2016. Learning socially normative robot navigation behaviors with Bayesian inverse reinforcement learning. In *Proceedings of the 2016 IEEE International Conference on Robotics and Automation (ICRA '16)*. IEEE, 2889–2895.
- [114] Martin T. Orne. 1962. On the social psychology of the psychological experiment: With particular reference to demand characteristics and their implications. *American Psychologist* 17, 11 (1962), 776.
- [115] Maike Paetzel, Giulia Perugia, and Ginevra Castellano. 2020. The persistence of first impressions: The effect of repeated interactions on the perception of a social robot. In *Proceedings of the 2020 ACM/IEEE International Conference on Human-Robot Interaction*, 73–82.
- [116] D. Paez-Granados, Y. He, D. Gonon, L. Huber, and A. Billard. 2021. 3d point cloud and RGBD of pedestrians in robot crowd navigation: Detection and tracking. *IEEE DataPort* 12 (2021).
- [117] Diego F. Paez-Granados, Hideki Kadone, Modar Hassan, Yang Chen, and Kenji Suzuki. 2022. Personal mobility with synchronous trunk-knee passive exoskeleton: Optimizing human-robot energy transfer. *IEEE/ASME Transactions on Mechatronics* 27, 5 (2022), 3613–3623. DOI: <https://doi.org/10.1109/TMECH.2021.3135453>
- [118] Hyun Soo Park, Jyh-Jing Hwang, Yedong Niu, and Jianbo Shi. 2016. Egocentric future localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4697–4705.
- [119] Stefano Pellegrini, Andreas Ess, Konrad Schindler, and Luc Van Gool. 2009. You'll never walk alone: Modeling social behavior for multi-target tracking. In *Proceedings of the 2009 IEEE 12th International Conference on Computer Vision*. IEEE, 261–268.
- [120] Claudia Pérez-D'Arpino, Can Liu, Patrick Goebel, Roberto Martín-Martín, and Silvio Savarese. 2021. Robot navigation in constrained pedestrian environments using reinforcement learning. In *Proceedings of the 2021 IEEE International Conference on Robotics and Automation*, 1140–1146.
- [121] Claudia Pérez-D'Arpino, Rebecca P. Khurshid, and Julie A. Shah. 2024. Experimental assessment of human-robot teaming for multi-step remote manipulation with expert operators. *ACM Transactions on Human-Robot Interaction* 13, 3 (2024). DOI: <https://doi.org/10.1145/3618258>
- [122] Noé Pérez-Higueras, Roberto Otero, Fernando Caballero, and Luis Merino. 2023. HuNavSim: A ROS 2 Human navigation simulator for benchmarking human-aware robot navigation. arXiv:2305.01303. Retrieved from <https://arxiv.org/abs/2305.01303>
- [123] Noé Pérez-Higueras, Rafael Ramón-Vigo, Fernando Caballero, and Luis Merino. 2014. Robot local navigation with learned social cost functions. In *Proceedings of the 2014 11th International Conference on Informatics in Control, Automation and Robotics (ICINCO '14)*, Vol. 2, IEEE, 618–625.

- [124] Daniel Perille, Abigail Truong, Xuesu Xiao, and Peter Stone. 2020. Benchmarking metric ground navigation. In *Proceedings of the 2020 IEEE International Symposium on Safety, Security, and Rescue Robotics (SSRR '20)*. IEEE, 116–121.
- [125] Björn Petrak, Gundula Sopper, Katharina Weitz, and Elisabeth André. 2021. Do you mind if I pass through? Studying the appropriate robot behavior when traversing two conversing people in a hallway setting. In *Proceedings of the 2021 30th IEEE International Conference on Robot & Human Interactive Communication (RO-MAN '21)*. IEEE, 369–375.
- [126] Sören Pirk, Edward Lee, Xuesu Xiao, Leila Takayama, Anthony Francis, and Alexander Toshev. 2022. A protocol for validating social navigation policies. arXiv:2204.05443. Retrieved from <https://arxiv.org/abs/2204.05443>
- [127] Dean A Pomerleau. 1989. Alvin: An autonomous land vehicle in a neural network. In *Advances in Neural Information Processing Systems*, 305–313.
- [128] Stephen R. Porter, Michael E. Whitcomb, and William H. Weitzer. 2004. Multiple surveys of students and survey fatigue. *New Directions for Institutional Research* 2004, 121 (2004), 63–73.
- [129] Morgan Quigley, Ken Conley, Brian Gerkey, Josh Faust, Tully Foote, Jeremy Leibs, Rob Wheeler, and Andrew Y. Ng. 2009. ROS: An open-source robot operating system. In *Proceedings of the ICRA Workshop on Open Source Software*, Vol. 3, 5.
- [130] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. arXiv:1606.05250. Retrieved from <https://arxiv.org/abs/1606.05250>
- [131] Tanmay Randhavane, Aniket Bera, Emily Kubin, Austin Wang, Kurt Gray, and Dinesh Manocha. 2019. Pedestrian dominance modeling for socially-aware robot navigation. In *Proceedings of the 2019 International Conference on Robotics and Automation (ICRA '19)*. IEEE, 5621–5628.
- [132] Carmine Recchiuto and Antonio Sgorbissa. 2022. Diversity-aware social robots meet people: Beyond context-aware embodied AI. arXiv:2207.05372. Retrieved from <https://arxiv.org/abs/2207.05372>
- [133] Jorge Rios-Martinez, Anne Spalanzani, and Christian Laugier. 2015. From proxemics theory to socially-aware navigation: A survey. *International Journal of Social Robotics* 7, 2 (2015), 137–153.
- [134] Alexandre Robicquet, Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. 2016. Learning social etiquette: Human trajectory understanding in crowded scenes. In *Proceedings of the European Conference on Computer Vision*. Springer, 549–565.
- [135] Astrid Rosenthal-von der Pütten, David Sirkin, Anna Abrams, and Laura Platte. 2020. The forgotten in HRI: Incidental encounters with robots in public spaces. In *Companion of the 2020 ACM/IEEE International Conference on Human-Robot Interaction*, 656–657.
- [136] Christoph Rösmann, Wendelin Feiten, Thomas Wösch, Frank Hoffmann, and Torsten Bertram. 2012. Trajectory modification considering dynamic constraints of autonomous robots. In *Proceedings of the 7th German Conference on Robotics (ROBOTIK '12)*. VDE, 1–6.
- [137] Christoph Rosmann, Artemi Makarow, Frank Hoffmann, and Torsten Bertram. 2017. Time-optimal nonlinear model predictive control with minimal control interventions. In *Proceedings of the 2017 IEEE Conference on Control Technology and Applications (CTA '17)*. IEEE, 19–24.
- [138] Andrey Rudenko, Tomasz P. Kucner, Chittaranjan S. Swaminathan, Ravi T. Chadalavada, Kai O. Arras, and Achim J. Lilienthal. 2020. THÖr: Human-robot navigation data collection and accurate motion trajectories dataset. *IEEE Robotics and Automation Letters* 5, 2 (2020), 676–682.
- [139] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. 2015. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision* 115 (2015), 211–252.
- [140] Shane Saunderson and Goldie Nejat. 2021. Robots asking for favors: The effects of directness and familiarity on persuasive HRI. *IEEE Robotics and Automation Letters* 6, 2 (2021), 1793–1800.
- [141] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, et al. 2019. Habitat: A platform for embodied AI research. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 9339–9347.
- [142] Mark A. Schmuckler. 2001. What is ecological validity? A dimensional analysis. *Infancy* 2, 4 (2001), 419–436.
- [143] John Schulman, Filip Wolski, Prfulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal policy optimization algorithms. arXiv:1707.06347. Retrieved from <https://arxiv.org/abs/1707.06347>
- [144] Duane P. Schultz. 1969. The human subject in psychological research. *Psychological Bulletin* 72, 3 (1969), 214.
- [145] Bokui Shen, Fei Xia, Chengshu Li, Roberto Martín-Martín, Linxi Fan, Guanzhi Wang, Claudia Pérez-D'Arpino, Shyamal Buch, Sanjana Srivastava, Lyne Tchapmi, et al. 2021. IGibson 1.0: A simulation environment for interactive tasks in large realistic scenes. In *Proceedings of the 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '21)*. IEEE, 7520–7527.
- [146] Phani Teja Singamaneni. 2022. *Combining Proactive Planning and Situation Analysis for Human-Aware Robot Navigation*. Theses. Université Paul Sabatier - Toulouse III. Retrieved from <https://theses.hal.science/tel-04006782>

- [147] Phani-Teja Singamaneni, Anthony Favier, and Rachid Alami. 2023. Towards benchmarking human-aware robot navigation: A new perspective and metrics. In *Proceedings of the 2023 32nd IEEE International Conference on Robot and Human Interactive Communication (RO-MAN '23)*. IEEE, 914–921.
- [148] Reginald G. Smart. 1966. Subject selection bias in psychological research. *Canadian Psychologist/Psychologie canadienne* 7, 2 (1966), 115.
- [149] Nicolas Spatola, Barbara Kühnlenz, and Gordon Cheng. 2021. Perception and evaluation in human–robot interaction: The human–robot interaction evaluation scale (HRIES)—A multicomponent approach of anthropomorphism. *International Journal of Social Robotics* 13, 7 (2021), 1517–1539.
- [150] Zayne Sprague, Rohan Chandra, Jarrett Holtz, and Joydeep Biswas. 2023. SOCIALGYM 2.0: Simulator for multi-agent social robot navigation in shared human spaces. arXiv:2303.05584. Retrieved from <https://arxiv.org/abs/2303.05584>
- [151] Nilesh Suriyarachchi, Rohan Chandra, John S. Baras, and Dinesh Manocha. 2022. GAMEOPT: Optimal real-time multi-agent planning and control for dynamic intersections. In *Proceedings of the 2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC '22)*. IEEE, 2599–2606.
- [152] Barbara G. Tabachnick and Linda S. Fidell. 2007. *Experimental Designs Using ANOVA*. Vol. 724. Thomson/Brooks/Cole Belmont, CA.
- [153] Mohsen Tavakol and Reg Dennick. 2011. Making sense of Cronbach’s Alpha. *International Journal of Medical Education* 2 (2011), 53.
- [154] Ada V. Taylor, Ellie Mamantov, and Henny Admoni. 2022. Observer-aware legibility for social navigation. In *Proceedings of the 2022 31st IEEE International Conference on Robot and Human Interactive Communication (RO-MAN '22)*. IEEE, 1115–1122.
- [155] J. Terry, Benjamin Black, Nathaniel Grammel, Mario Jayakumar, Ananth Hari, Ryan Sullivan, Luis S. Santos, Clemens Dieffendahl, Caroline Horsch, Rodrigo Perez-Vicente, et al. 2021. Pettingzoo: Gym for multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems* 34 (2021), 15032–15043.
- [156] Peter Trautman and Andreas Krause. 2010. Unfreezing the robot: Navigation in dense, interacting crowds. In *Proceedings of the 2010 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 797–803.
- [157] Xuan-Tung Truong and Trung Dung Ngo. 2017. Toward socially aware robot navigation in dynamic and crowded environments: A proactive social motion model. *IEEE Transactions on Automation Science and Engineering* 14, 4 (2017), 1743–1760.
- [158] Xuan-Tung Truong and Trung-Dung Ngo. 2017. “To approach humans?”: A unified framework for approaching pose prediction and socially aware robot navigation. *IEEE Transactions on Cognitive and Developmental Systems* 10, 3 (2017), 557–572.
- [159] Nathan Tsoi, Mohamed Hussein, Jeacy Espinoza, Xavier Ruiz, and Marynel Vázquez. 2020. Sean: Social environment for autonomous navigation. In *Proceedings of the 8th International Conference on Human-Agent Interaction*, 281–283.
- [160] Nathan Tsoi, Mohamed Hussein, Olivia Fugikawa, J. D. Zhao, and Marynel Vázquez. 2021. An approach to deploy interactive robotic simulators on the web for HRI experiments: Results in social robot navigation. In *Proceedings of the 2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '21)*. IEEE, 7528–7535.
- [161] Nathan Tsoi, Alec Xiang, Peter Yu, Samuel S. Sohn, Greg Schwartz, Subashri Ramesh, Mohamed Hussein, Anjali W. Gupta, Mubbasir Kapadia, and Marynel Vázquez. 2022. Sean 2.0: Formalizing and generating social situations for robot navigation. *IEEE Robotics and Automation Letters* 7, 4 (2022), 11047–11054.
- [162] Jur Van Den Berg, Stephen J. Guy, Ming Lin, and Dinesh Manocha. 2011. Reciprocal n-body collision avoidance. In *Robotics Research: The 14th International Symposium ISRR (2011)*, Springer, 3–19.
- [163] Jur van den Berg, Stephen J. Guy, Jamie Snape, Ming C. Lin, and Dinesh Manocha. 2011. Rvo2 library: Reciprocal collision avoidance for real-time multi-agent simulation. Retrieved from <https://gamma.cs.unc.edu/RVO2>
- [164] Stanford Vision and Learning Lab. 2022. iGibsonChallenge2021. Retrieved from <https://github.com/StanfordVL/iGibsonChallenge2021>
- [165] Nancy J. Wahl. 1999. An overview of regression testing. *ACM SIGSOFT Software Engineering Notes* 24, 1 (1999), 69–73.
- [166] Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in Neural Information Processing Systems* 32 (2019).
- [167] Junxian Wang, Wesley P. Chan, Pamela Carreno-Medrano, Akansel Cosgun, and Elizabeth Croft. 2022. Metrics for evaluating social conformity of crowd navigation algorithms. In *Proceedings of the 2022 IEEE International Conference on Advanced Robotics and Its Social Impacts (ARSO '22)*. IEEE, 1–6.
- [168] Erik Wijmans, Abhishek Kadian, Ari Morcos, Stefan Lee, Irfan Essa, Devi Parikh, Manolis Savva, and Dhruv Batra. 2019. DD-PPO: Learning near-perfect pointgoal navigators from 2.5 billion frames. arXiv:1911.00357. Retrieved from <https://arxiv.org/abs/1911.00357>

- [169] Fei Xia, Chengshu Eric Li, William Shen, Priya Kasimbeg, Roberto Martin-Martin, Alexander Toshev, and Silvio Savarese. 2020. Interactive Gibson: Benchmark for interactive navigation in cluttered environments. *IEEE Robotics and Automation Letters* 5, 2 (2020), 713–720.
- [170] Xuesu Xiao, Bo Liu, Garrett Warnell, and Peter Stone. 2022. Motion planning and control for mobile robot navigation using machine learning: A survey. *Autonomous Robots* 46, 5 (2022), 569–597.
- [171] Xuesu Xiao, Tingnan Zhang, Krzysztof Marcin Choromanski, Tsang-Wei Edward Lee, Anthony Francis, Jake Varley, Stephen Tu, Sumeet Singh, Peng Xu, Fei Xia, et al. 2022. Learning model predictive controllers with real-time attention for real-world navigation. In *Proceedings of the Conference on Robot Learning*. PMLR.
- [172] Zhi Yan, Tom Duckett, and Nicola Bellotto. 2017. Online learning for human classification in 3D LiDAR-based tracking. In *Proceedings of the 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS '17)*. IEEE, 864–871.

Received 27 September 2023; revised 27 September 2023; accepted 20 June 2024