

Plan Explanations that Exploit a Cognitive Spatial Model

Raj Korpan

The Graduate Center
City University of New York
rkorpan@gradcenter.cuny.edu

Susan L. Epstein

The Graduate Center and Hunter College
City University of New York
susan.epstein@hunter.cuny.edu

Abstract

Ideally, people who navigate together in a complex indoor space share a mental model that facilitates explanation. This paper reports on a robot control system whose cognitive world model is based on spatial affordances that generalize over its perceptual data. Given a target, the control system formulates multiple plans, each with a model-relevant metric, and selects among them. As a result, it can provide readily understandable natural language about the robot’s intentions and confidence, and generate diverse, contrastive explanations that reference the acquired spatial model. Empirical results in large, complex environments demonstrate the robot’s ability to provide human-friendly explanations in natural language.

1 Introduction

Inspired by recent recommendations for spoken language interaction with robots (Marge et al., 2020), this paper introduces WHY, an approach to communicate a robot’s planning rationales, intentions, and confidence in human-friendly spatial language. Our thesis is that a plan based on spatial representations acquired from travel experience can ground its objectives and support explainable path planning. The principal results of this paper are empirical demonstrations of WHY’s ability to explain and contrast plans in readily-understandable natural language.

Given sensor data and a metric map (e.g., a floor plan), the task of our autonomous robot navigator is to travel to target locations in a large, complex, human-centric, indoor space (henceforward, *world*). The robot’s control system integrates acquired spatial knowledge into a cognitively-based architecture that combines planning with reactivity, heuristics, and situational reasoning. Given a target, the control system creates a *plan*, a sequence of intermediate locations (*waypoints*) to reach it. This plan is expected to balance multiple objectives, combine continuous and discrete spatial representations, and encourage a human’s trust.

Traditional navigation planners use a *cost graph* (also known as a *costmap*) where each node is a point in unobstructed space and each edge connects a pair of nodes with a weight for the cost to move between them. A popular cost graph is based on an *occupancy grid*, uniform square cells superimposed on a two-dimensional metric map. Each edge in the graph represents two adjacent unobstructed cells, labeled with the Euclidean distance between their centers. In a fine-grained grid, however, optimal planners (e.g., A^* (Hart et al., 1968)) hug obstacles so tightly that their plans require tight maneuvers to reach some waypoints and may fail as actuator and sensor errors accumulate near them.

To bias plans toward its particular *objective* (a spatial representation or commonsense rationale), a planner modifies the weights in its own copy of the occupancy-grid graph. The fixed underlying graph structure allows our approach to evaluate a plan within any such modified graph. Voting then selects the plan that best satisfies all the objectives. This approach facilitates contrastive natural-language explanations of the chosen plan with respect to each objective. The control system reports on its beliefs, intentions, and confidence with spatial language. For example, “Although there may be another way that is somewhat shorter, I think my way is a lot better at going through open areas.”

The next sections provide related work and describe the acquired spatial model. Subsequent sections cover the modified graphs, vote-based planning, and how WHY explains plans. The last sections describe empirical results and future work.

2 Related work

A spatial representation of its world is essential to a robot control system that navigates efficiently and explains its behavior clearly. Grounded communication between a robot and a person, however, requires a shared spatial representation. This section first describes work on human cognitive maps

that inspired our control system’s spatial model. It then details approaches that describe and explain the robot’s behavior.

A *cognitive map* is a compact, mental spatial representation of a world, built by a person as she moves through that world (Golledge, 1999). To reduce her cognitive load, a person reasons from a cognitive map that incorporates landmarks, route knowledge, and survey knowledge (Tversky, 1993). Landmarks represent locations in the map, routes represent lines that connect them, and survey knowledge captures spatial relations. Although it has been suggested that cognitive maps use metric distances and angles (Gallistel, 1990), more recent work indicates that cognitive maps have a non-metric, qualitative topological structure (Foo et al., 2005). Other recent work suggests that people use a cognitive graph with labeled metric information that captures connectivity and patterns (Chrastil and Warren, 2014; Warren et al., 2017).

An *affordance* is a characteristic of the world that enables the execution of some action (Gibson, 1977). Affordance-based theories of spatial cognition posit a tight relationship between the specific dynamics of a world and the decisions made by an individual there (Fajen and Phillips, 2013). Here, a *spatial affordance* is an abstract representation of the world that facilitates navigation. This paper introduces path planning in cost graphs based on acquired spatial affordances. People generalize structured representations across domains on similar tasks (Pouncy et al., 2021) much the way the spatial model described here generalizes affordances for use in different worlds.

A control system can learn and use a cognitive map of its world for robot navigation. For example, the Spatial Semantic Hierarchy (SSH) modeled a cognitive map with hierarchical metric and topological representations (Kuipers, 2000). Although SSH’s cognitive map bears some similarity to the one used here, it did not explain plans. Other approaches used semantics to create a meaningfully-labeled metric map (Kostavelis and Gasteratos, 2015). While these maps provide a qualitative context in which to ground a controller’s language, they do not necessarily align with human cognitive maps. Moreover, control systems often use semantic maps for communication but another representation for reasoning and decision-making. Instead, this paper shows how a single, affordance-based representation supports all of those processes.

Indoors, an autonomous robot may interact with people as it navigates to its target. A human collaborator is more likely to accept, trust, and understand a robot that can explain its behavior (Rosenfeld and Richardson, 2019). Rather than describe an event or summarize its causes, an explanation compares counterfactual cases, includes causes selectively, and recognizes people as social beings with beliefs and intentions (Miller, 2019). A *contrastive* explanation compares the reason for a decision to another plausible rationale (Hoffmann and Magazzini, 2019). Human subjects generally prefer such explanations that focus on the difference between the robot’s planned route and their own (e.g., “my route is shorter, but overlaps more and produces less reward”) (Perelman et al., 2020).

Detailed technical logs of a robot’s experience were originally available only to trained researchers (Landsiedel et al., 2017; Scalise et al., 2017). Recent work, however, has generated natural language descriptions of a robot’s travelled path from them. These focus on abstraction, specificity, and locality (Rosenthal et al., 2016; Perera et al., 2016) or on sentence correctness, completeness, and conciseness (Barrett et al., 2017). All, however, required a labeled dataset or a semantic map. Other recent work partitions a plan into actions and uses language templates to generate descriptions of each action in the context of a collaborating robot team (Singh et al., 2021). WHY focuses on explanations for the reasons behind the robot’s decisions rather than descriptions of the robot’s behavior.

To produce explanations, others have selected potentially suboptimal plans (Fox et al., 2017; Chakraborti et al., 2019) or readily understandable behaviors (Huang et al., 2019), or relied on classical planning (Magnaguagno et al., 2017; Grea et al., 2018; Krarup et al., 2019) or on logic (Seegebarth et al., 2012; Nguyen et al., 2020). None of that work, however, explains in natural language. The approach closest to the one presented here provides contrastive explanations for multi-objective path planning in natural language as a Markov decision process (Sukkerd et al., 2020), but considers fewer objectives, requires a hand-labeled map, and has been evaluated only in much smaller worlds.

3 Spatial affordances

The context of this work is *SemaFORR*, a cognitively-based control system for autonomous indoor navigation (Epstein et al., 2015; Epstein and

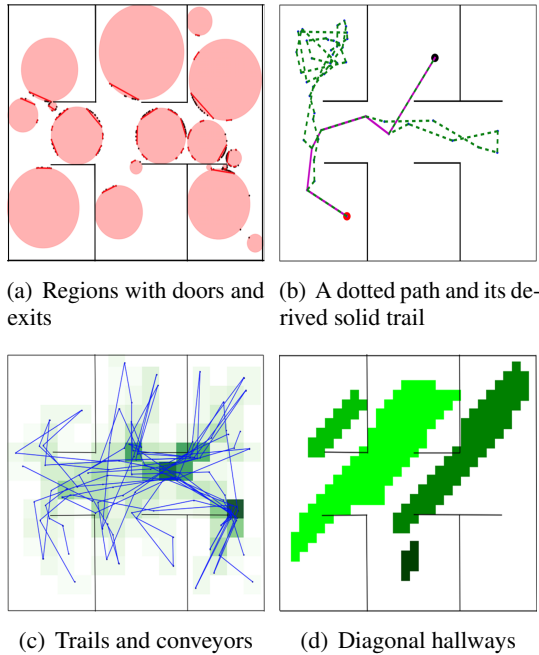


Figure 1: Affordances in a simple artificial world

Korpan, 2019). At decision point $d = \langle x, y, \theta, V \rangle$, SemaFORR records the robot’s location (x, y) , its orientation θ , and its view V , the data from its on-board range finder. After each target, SemaFORR identifies spatial affordances for its acquired model of *freespace*, the unobstructed areas in a world. The model can be used alone or with a metric map.

At decision point d , SemaFORR learns a *region*, a circle in freespace with center at (x, y) and radius equal to the minimum distance reported by V . Accumulated contradictory or overlapping regions are resolved after each target. An *exit* represents access to freespace, a point where the robot’s path once crossed the region’s perimeter. A *door* is an arc on a region’s perimeter, a continuous generalization of finitely many, relatively close exits between its endpoints. Figure 1(a) shows acquired regions with exits and doors (drawn for clarity as secants to their respective arcs). Although regions approximate what appear to be rooms in the figure, they record only freespace, not walls.

A *trail* is a refined version of the robot’s path toward its target. The algorithm that creates trails heuristically smooths the robot’s paths and eliminates digressions. The remaining (usually far fewer) decision points are *trail markers*. As in Figure 1(b), the sequence of line segments defined by consecutive trail markers is typically more direct than the original path, but rarely optimal. A *conveyor* is a freespace cell in a $2 \times 2m$ grid super-

Table 1: SemaFORR’s planners and their objectives

Planner	Objective
FASTP	Minimize distance traveled
SAFEPP	Avoid obstacles
EXPLOREP	Avoid paths
NOVELP	Avoid spatial model
CONVEYP	Exploit conveyors
HALLWAYP	Exploit hallways
REGIONP	Exploit regions, doors, exits
TRAILP	Exploit trail markers

imposed on the world’s footprint. Conveyors tally how often trails pass through them. Higher-count cells represent locations that frequently support travel. They appear darker in Figure 1(c).

A hallway represents well-travelled routes in some *angular direction* (vertical, horizontal, major diagonal, or minor diagonal). A hallway generalizes line segments between consecutive decision points to find relatively straight, narrow, continuous freespace with both length and width. Figure 1(d) shows some acquired minor-diagonal hallways.

4 Modified cost graphs

Planning for navigation requires a graphical representation of the world’s freespace. To produce an optimal plan, A* searches a cost graph G based on an occupancy grid with edge weights for Euclidean distance. SemaFORR constructs a set of graphs; each begins with G but modifies its edge weights to align with a particular objective. This biases search toward that objective but still considers plan length. In practice, an occupancy grid should be sufficiently fine to represent obstacles accurately.

Table 1 lists SemaFORR’s planners and their objectives. Given a target, each planner formulates its own plan to reach it, one biased toward its own objective. Two planners focus on common-sense: FASTP searches the original G , but SAFEPP increases G ’s edge weights based on an edge’s proximity to obstacles. Two others focus on exploration to acquire more knowledge about their world. EXPLOREP creates a grid that tallies how frequently the robot’s path history passes through each cell, and uses those values to increase edge weights where it has already traveled. Because the acquired spatial model summarizes experience more compactly than a path, NOVELP explores areas not covered by the model. It increases a weight if the edge overlaps an acquired affordance.

Four planners exploit a particular kind of spatial affordance with changes to edge weights. (Values based on preliminary testing bias plans to pursue but not overemphasize affordances.) REGIONP’s cost graph modifies each edge’s weight w based on the location of its endpoints. If both lie in the same region, w goes unchanged; if neither lies in a region w becomes $10w$. Otherwise, for the one endpoint v not in a region, w becomes $1.5w$ if v is within $0.5m$ of a door and an exit, $1.75w$ if v is within $0.5m$ of a door or an exit, and otherwise $2w$. This biases plans to pass through regions because it increases edge costs outside them.

HALLWAYP and TRAILP modify their weights similarly, with respective conditions “lie in one hallway” and “lie within $0.5m$ of a trail marker.” If both endpoints of an edge meet the condition, w goes unchanged; if neither does, w becomes $10w$. Otherwise, when just one endpoint meets the condition, w becomes $1.5w$. To bias plans toward high-count conveyors, CONVEYP considers the counters c_1 and c_2 for the cells where the endpoints of an edge with weight w lies. If both are non-zero, w becomes $w + 2/(c_1 + c_2)$; otherwise, w becomes $10w$.

Because SemaFORR’s spatial model focuses on freespace, these modified cost graphs allow a robot control system to encourage travel there but also incorporate the metric cost graph where the model lacks knowledge. The region-based cost graph, for example, imposes relatively lower costs only for doors and exits that the robot has successfully exploited earlier, and thus prioritizes them. Because weights only increase, Euclidean distance remains an *admissible* heuristic for A^* , that is, it never overestimates the actual cost to the target’s location.

5 Voting among planners

To choose paths, people use many different objectives that reflect their motivation (Golledge, 1999). A cognitively-based robot navigator should also incorporate and balance a variety of path-selection heuristics. SemaFORR’s planners can be used together because they originate from the same cost graph. This section explains Algorithm 1, pseudocode for how voting balances the planners’ objectives to select a plan.

SemaFORR constructs multiple plans that optimize a single objective and then uses voting to select the plan that maximally satisfies the most objectives. First, each planner j constructs an op-

Algorithm 1: Voting-based planning

Input: *planners* J , *spatial model* M , *basic cost graph* G

for each planner $j \in J$ **do**

Set j ’s cost graph G_j to a copy of G
 Modify G_j ’s weights based on j and M
 With A^* , find optimal plan P_j in G_j

for each planner $j \in J$ **do**

for each planner $i \in J$ **do**
 $C_{ij} \leftarrow$ cost of plan P_i in G_j
 Normalize plan scores C_{ij} in $[0,10]$

for each plan P_i **do**

$Score_i \leftarrow \sum_{j=1}^J C_{ij}$

$best \leftarrow \operatorname{argmin}_i Score_i$

return P_{best}

timal plan P_j for its objective as a sequence of waypoints in its modified cost graph G_j . This guarantees that each submitted plan is optimal for at least one objective.

Next, each planner’s objective is used to evaluate every plan. All the cost graphs have the same nodes and edges, so to evaluate planner i ’s plan P_i from the perspective of planner j , SemaFORR simply sums the edge weights in G_j for the sequence of edges specified by P_i . The resultant scores C_{ij} are then normalized in $[0, 10]$ for each j . SemaFORR seeks to minimize its objectives. Thus a C_{ij} value near 0 indicates that plan P_i closely conforms to objective j , while a score near 10 indicates that plan P_i conflicts with objective j . Voting selects the plan with the lowest total score across all objectives and breaks ties at random.

6 Contrastive explanations

SemaFORR uses WHY to explain its long-range perspective in natural language. WHY exploits differences among planners’ objectives to produce clear, concise, contrastive explanations for a plan quickly. WHY assumes that the robot’s human companion seeks a shortest-length plan, and compares that to SemaFORR’s plan. Although we assume here that a goal-directed human navigator would seek to minimize travel distance, another objective, including those in Table 1, could label the foundational cost graph G instead.

Throughout this section, \mathcal{N} represents a function that translates its argument (a planner or a metric value) into natural language. Given a real-valued

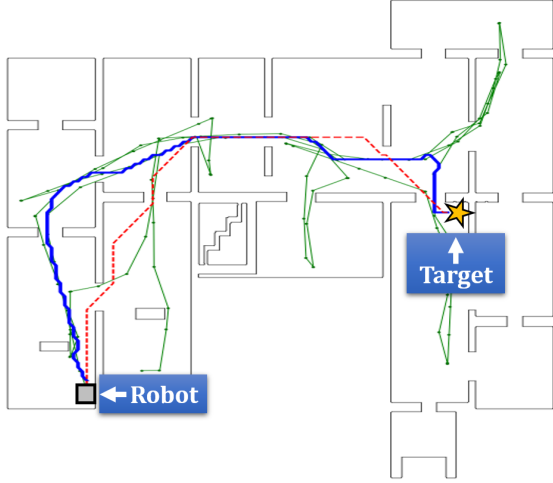


Figure 2: WHY compares FASTP’s (red) plan to TRAILP’s (blue) one biased by SemaFORR’s (green) trails. It explains, “Although there may be another way that is a lot shorter, I think my way is a lot better at following ways we’ve gone before.”

metric m for some aspect (e.g., confidence or enthusiasm) of the decision process, $\mathcal{M}(m)$ bins m ’s value into an ordered partition of m ’s range and $\mathcal{N}(\mathcal{M}(m))$ translates that bin to a natural language phrase. For example, m could measure the desire to select one plan over the others, and the value partition could distinguish a strong preference for that plan from a weak one. Thus, if $m \in (0, +\infty)$ were partitioned as $\{(0, 5), [5, +\infty)\}$, $\mathcal{N}(m < 5)$ could be “a little” and $\mathcal{N}(m \geq 5)$ “a lot.” This allows WHY to hedge in its responses, much the way people explain their reasoning when they are uncertain (Markkanen and Schröder, 1997).

6.1 Why does your plan go this way?

Human and robot plans to reach the same target may differ because they lack a common objective. WHY’s response to this question presumes that a human plans from one perspective, objective β_H , while the robot plans from another perspective, objective β_R . Explanations for a plan assume a human has an alternative objective. Henceforward, β_H is “take the shortest path.”

WHY models the human questioner with β_H to produce plan P_H , a prediction of the human’s implicit plan. Algorithm 2 is pseudocode for WHY’s plan-explanation procedure. WHY takes as input the robot’s plan P_R and objective β_R , and the alternative plan P_H and objective β_H it attributes to the human questioner. $\beta_H(P)$ measures plan length and $\beta_R(P)$ measures plan cost in P_R ’s graph. In the running example shown in Figure 2, WHY ex-

Algorithm 2: Explanation procedure

Input: *planning objectives* β_R and β_H ,
plans P_R and P_H
Output: *explanation*
 $\mathcal{D}_R = \beta_R(P_R) - \beta_R(P_H)$
 $\mathcal{D}_H = \beta_H(P_R) - \beta_H(P_H)$
switch *mode*($\mathcal{D}_R, \mathcal{D}_H$) **do**
 case $\mathcal{D}_R = \mathcal{D}_H = 0$ **do**
 | *explanation* \leftarrow sentence based on
 | template for equivalent plans
 case $\mathcal{D}_R < 0$ and $\mathcal{D}_H > 0$ **do**
 | *explanation* \leftarrow sentence for β_R, β_H
 case $\mathcal{D}_R < 0$ and $\mathcal{D}_H = 0$ **do**
 | *explanation* \leftarrow sentence for β_R
return *explanation*

plains SemaFORR’s preference for its plan P_R from TRAILP where β_R is TRAILP’s objective (“exploit trail markers”). WHY translates β_H and β_R with Table 2 as “short” and “follows ways we’ve gone before,” respectively.

If voting selected the plan constructed by FASTP (i.e., the shortest-length plan), then Why responds with “I decided to go this way because I agree that we should take the shortest route.” Otherwise, to compare P_R with P_H , WHY calculates their difference from two perspectives: \mathcal{D}_H from the human’s perspective (e.g., length), and \mathcal{D}_R from the robot’s perspective (e.g., proximity to trails). WHY places these differences in user-specified bins that represent a human perspective on the objectives. Table 3 provides language for these differences.

The relative size of the differences determines an applicable template. If both \mathcal{D}_H and \mathcal{D}_R , as defined in Algorithm 2, are 0, then the plans equally address the two objectives, and WHY explains “I decided to go this way because I think it’s just as $\mathcal{N}(\mathcal{D}_H)$ and equally $\mathcal{N}(\mathcal{D}_R)$.” Otherwise, the plans differ with respect to one or both objectives. If \mathcal{D}_R is negative (e.g., P_R is more aligned with trails), then WHY instantiates this template:

- 1: Although there may be another way that is $\mathcal{N}(\mathcal{M}(\mathcal{D}_H)) \mathcal{N}^*(\beta_H)$,
- 2: I think my way is $\mathcal{N}(\mathcal{M}(\mathcal{D}_R)) \mathcal{N}^*(\beta_R)$.

where $\mathcal{N}^*(\beta)$ is a comparator for β (e.g., “shorter” or “better at following ways we’ve gone before”). For example, “Although there may be another way that is somewhat shorter, I think my way is a lot better at following ways we’ve gone before.” WHY omits line 1 in the template if $\mathcal{D}_H = 0$. Other cases,

Table 2: Language for the planners’ objectives. $\mathcal{N}^*(\beta)$ and $\mathcal{N}'(\beta)$ values for FASTP and EXPLOREP are as shown. For the others, $\mathcal{N}^*(\beta) \approx \mathcal{N}'(\beta)$, where $\mathcal{N}^*(\beta)$ begins with “better at” and $\mathcal{N}'(\beta)$ begins with “worse at.”

Planner	$\mathcal{N}(\beta)$	$\mathcal{N}^*(\beta)$	$\mathcal{N}'(\beta)$
FASTP	short	shorter	longer
EXPLOREP	goes a new way	newer	familiar
SAFEP	stays far from obstacles	staying far from obstacles	
NOVELP	learns something new	learning something new	
CONVEYP	goes through well-traveled areas	going through well-traveled areas	
HALLWAYP	follows hallways	following hallways	
REGIONP	goes through open areas	going through open areas	
TRAILP	follows ways we’ve gone before	following ways we’ve gone before	

Table 3: Language for value intervals for the difference \mathcal{D} . For affordance-based planners $a=150$ and $b=25$, for SAFEP $a=0.35$ and $b=0.15$, for EXPLOREP $a=100$ and $b=15$, and for NOVELP $a=350$ and $b=100$.

Planner	Intervals $\mathcal{M}(\mathcal{D})$	$\mathcal{N}(\mathcal{M}(\mathcal{D}))$
	$(0, 1]$	a bit
FASTP	$(1, 10]$	somewhat
	$(10, +\infty)$	a lot
	$(-\infty, -a]$	a lot
All others	$(-a, -b]$	somewhat
	$(-b, +\infty)$	a bit

where $\mathcal{D}_H < 0$ or $\mathcal{D}_R > 0$ cannot occur because each planner is optimal with respect to its own cost graph and objective, as described in Section 5.

6.2 Why do you prefer your plan?

WHY also addresses the question “Why do you prefer your plan?” Unlike the previous response, which contrasted the human’s objective with the robot’s, this response has the robot explain its objective. If voting selects the FASTP plan, which the robot assumes has the same objective as its human companion, WHY would respond “Actually, I agree that we should take the shortest route.” Otherwise, WHY uses the differences \mathcal{D}_H and \mathcal{D}_R from Algorithm 2. If they are both 0, then WHY replies, “I think both plans are equally good.” Otherwise, WHY responds with the template “I prefer my plan because it’s $\mathcal{N}(\mathcal{M}(\mathcal{D}_R)) \mathcal{N}^*(\beta_R)$.” For example, to explain why SemaFORR chose TRAILP’s plan, WHY might say “I prefer my plan because it’s a lot better at following ways we’ve gone before.”

6.3 What’s another way we could go?

Figure 3 shows an example where WHY responds to “What’s another way we could go?” Because WHY has access to two plans from SemaFORR

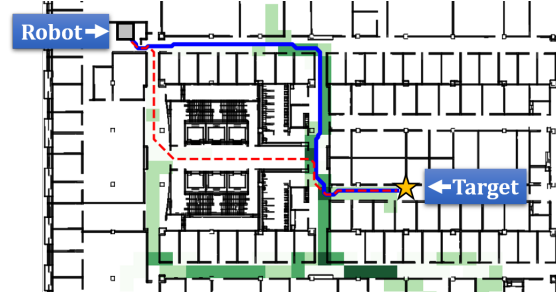


Figure 3: Acquired conveyors in green, with darker higher-count cells. Voting chose CONVEYP’s (blue) plan which is drawn to high-count cells. In response to “What’s another way we could go?” WHY compares the conveyor plan with FASTP’s (red) plan: “We could go that way since it’s a bit shorter but it could also be a bit worse at going through well-traveled areas.”

(P_R and P_H), it can provide P_H , the shortest-path plan, as the alternative plan in response. If voting selects the FASTP plan, which uses the same objective as the robot’s human companion, then WHY responds “Yours is the best way to go.” Otherwise, it instantiates the template: “We could go your way since it’s $\mathcal{N}(\mathcal{M}(\mathcal{D}_H)) \mathcal{N}^*(\beta_H)$ but it could also be $\mathcal{N}(\mathcal{M}(\mathcal{D}_R)) \mathcal{N}'(\beta_R)$.” Here \mathcal{N}' denotes an opposite comparator (e.g., “longer” or “worse at following ways we’ve gone before”). For example, an explanation is “We could go that way since it’s somewhat shorter but it could also be a lot worse at following ways we’ve gone before.”

6.4 How sure are you about your plan?

In response to “How sure are you about your plan?” WHY explains its confidence that P_R meets its objective. Figure 4 shows an example. WHY uses the language for $\mathcal{M}(\mathcal{D}_R)$ and $\mathcal{M}(\mathcal{D}_H)$ from Table 3 to extract a value $\mathcal{C} = \mathcal{N}(\mathcal{M}(\mathcal{D}_R, \mathcal{D}_H))$ from Table 4. WHY then instantiates “I’m $\mathcal{N}(\mathcal{C})$ sure because” followed by line \mathcal{C} :

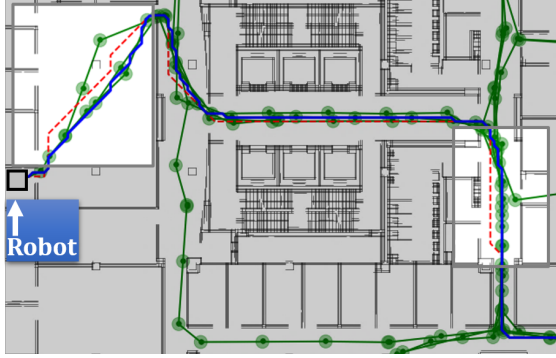


Figure 4: Highlighted sections of FASTP’s (red) plan and TRAILP’s (blue) plan to follow acquired (green circle) trail markers. WHY explains “I’m really sure because my plan is a lot better at following ways we’ve gone before and only a bit longer than your plan.”

Table 4: Language $\mathcal{N}(\mathcal{M}(\mathcal{D}_R, \mathcal{D}_H))$ for confidence compares $\mathcal{M}(\mathcal{D}_R)$ and $\mathcal{M}(\mathcal{D}_H)$ from Table 3. Here, 1 denotes “really,” 2 = “only somewhat,” and 3 = “not.”

$\mathcal{N}(\mathcal{M}(\mathcal{D}_R))$	$\mathcal{N}(\mathcal{M}(\mathcal{D}_H))$		
	“a lot”	“somewhat”	“a bit”
“a lot”	2	1	1
“somewhat”	3	2	1
“a bit”	3	3	2

1: my plan is $\mathcal{N}(\mathcal{M}(\mathcal{D}_R)) \mathcal{N}^*(\beta_R)$ and only $\mathcal{N}(\mathcal{M}(\mathcal{D}_H)) \mathcal{N}'(\beta_H)$ than yours.

2: even though my plan is $\mathcal{N}(\mathcal{M}(\mathcal{D}_R)) \mathcal{N}^*(\beta_R)$, it is also $\mathcal{N}(\mathcal{M}(\mathcal{D}_H)) \mathcal{N}'(\beta_H)$ than yours.

3: my plan is $\mathcal{N}(\mathcal{D}_H) \mathcal{N}'(\beta_H)$ and only $\mathcal{N}(\mathcal{D}_R) \mathcal{N}^*(\beta_R)$ than yours

6.5 How are we getting there?

“How are we getting there?” shows a human companion’s uncertainty about the route planned to reach their shared target. Rather than reference the planner’s objective, WHY treats this as a request for a high-level description of P_R itself, and uses the segments between consecutive waypoints in SemaFORR’s plan P_R to produce natural language that describes it. Figure 5 shows an example.

WHY anticipates travel with P_R as an ordered sequence of locations from the robot’s current location through P_R ’s waypoints and then to the target. First, WHY forms plan segments from consecutive locations in P_R and computes each segment’s length and angular direction χ (based on the angle between its endpoints relative to a fixed horizontal axis). It then bins χ within an interval $\mathcal{M}(\chi)$ and assigns a label $\mathcal{N}(\mathcal{M}(\chi))$ as shown in Table 5.

These labels are allocentric, and therefore less

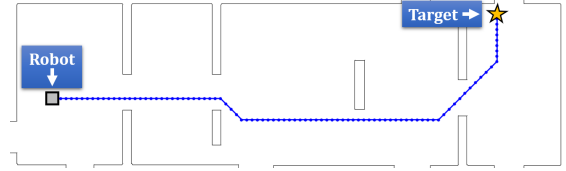


Figure 5: SemaFORR’s FASTP plan with 92 waypoints from the robot to its target. WHY explains in 9 clauses, “We will go straight about 20 meters, turn right a little, go straight about 4 meters, turn left a little, go straight about 20 meters, turn left a little, go straight about 8 meters, turn left a little, and go straight about 4 meters to reach our target.”

Table 5: Labels $\mathcal{N}(\mathcal{M}(\chi))$ for segment angle intervals $\mathcal{M}(\chi)$. Language $\mathcal{N}(\alpha)$ adjusts the change in consecutive angular directions for full 2π rotation: $\alpha = \mathcal{N}(\mathcal{M}(\chi_k)) - \mathcal{N}(\mathcal{M}(\chi_{k-1})) \bmod 8$.

$\mathcal{M}(\chi)$	$\mathcal{N}(\mathcal{M}(\chi))$	α	Phrase $\mathcal{N}(\alpha)$
$[-\frac{7\pi}{8}, -\frac{5\pi}{8})$	2	0	go straight
$[-\frac{6\pi}{8}, -\frac{3\pi}{8})$	3	1	turn left a little
$[-\frac{3\pi}{8}, \frac{\pi}{8})$	4	2	turn left
$[-\frac{\pi}{8}, \frac{\pi}{8})$	5	3	turn hard left
$[\frac{\pi}{8}, \frac{3\pi}{8})$	6	4	turn around
$[\frac{3\pi}{8}, \frac{5\pi}{8})$	7	5	turn hard right
$[\frac{5\pi}{8}, \frac{7\pi}{8})$	8	6	turn right
otherwise	1	7	turn right a little

appropriate indoors. WHY translates them to an egocentric frame of reference, as if the robot and its companion faced the same way along the intended route. The change in consecutive $\mathcal{N}(\mathcal{M}(\chi))$ labels represents the change in direction from one path segment to the next. $\mathcal{N}(\alpha)$ is language for α , the angular change in χ from one segment to the next. For example, if the first segment in P_R were labeled 2 and the second segment labeled 7, then $\alpha = 5$ which Table 5 translates as “turn hard right.”

Plan P_R now has a sequence of phrases for the points where two consecutive segments meet. WHY inserts a “go straight” after each “turn” phrase. WHY then summarizes consecutive “go straight” phrases into a single one (since they indicate no change in direction) with a length \mathcal{L} , the sum of the lengths of the segments that induced it. These \mathcal{L} s are binned into intervals and reported in natural language (e.g., 5.7m lies in (4, 6] with language “about 6 meters”).

WHY combines the list of phrases and lengths appropriately to form a succinct explanation with the template “We will [$\mathcal{N}(\alpha)$ {about $\mathcal{N}(\mathcal{M}(\mathcal{L}))$ },] to reach our target.” It repeats the material in square

Table 6: How often planners won the vote

Planner	M5	H10	G5	Total
FASTP	25.0%	42.9%	32.4%	33.4%
SAFE P	37.0%	25.7%	27.5%	30.1%
EXPLOREP	9.0%	6.9%	4.9%	6.9%
NOVELP	0.0%	0.0%	0.0%	0.0%
CONVEYP	14.0%	7.4%	16.5%	12.6%
HALLWAYP	6.0%	9.1%	6.6%	7.2%
REGIONP	5.5%	6.3%	0.5%	4.1%
TRAILP	3.5%	1.7%	11.5%	5.6%

brackets for each $\mathcal{N}(\alpha)$, and includes the material in curly brackets only when $\mathcal{N}(\alpha)$ is “go straight.”

In summary, WHY produces natural explanations for a robot’s plan as it travels through a complex world. These explanations are essential for human-friendly autonomous indoor navigation and require an assumption about its human collaborator’s objective. Our approach explains the robot’s plan, responds to questions about alternatives, and expresses a human-friendly level of confidence.

7 Empirical Evaluation

SemaFORR with WHY is evaluated on three challenging real worlds: M5, H10, and G5. M5 is the fifth floor of New York’s Museum of Modern Art. It is $54 \times 62m$ and has $1585m^2$ freespace. H10 is the $89 \times 58m$ tenth floor of an academic building with $2627m^2$ of freespace and 75 rooms. G5 is the $110 \times 70m$ fifth floor of a renovated Manhattan building. G5 has about $4021m^2$ of freespace, 180 rooms, and many intersecting hallways. It is known for its ability to perplex human navigators, despite color-coded walls and art introduced as landmarks. All testing was in simulation with ROS, the state-of-the-art robot operating system (Quigley et al., 2009). MengeROS manages the simulation and deliberately introduces error into both the sensor data and action execution (Aroor et al., 2017).

To evaluate WHY we randomly sampled 5 sequences of 40 targets in each world’s freespace. Table 6 reports how often voting selected each planner’s submission. Two-thirds of the selected plans were based on a modified cost graph, about half of them biased by SemaFORR’s spatial model. Because SemaFORR revises its model incrementally, as the robot addresses more targets, it begins to value EXPLOREP’s plans less than model-based ones. For example, by the second 20 targets in each sequence of 40, plans based on the spatial model

Table 7: Analysis of explanation results with number of unique phrasings and average readability scores

Unique phrasings	M5	H10	G5	All
Why this way?	38	30	39	49
How sure are you?	24	19	26	30
Another way?	24	19	26	30
Why yours?	17	15	16	18
How to get there?	199	175	182	556
Average readability	M5	H10	G5	All
Why this way?	4.7	5.3	5.3	5.1
How sure are you?	6.6	6.6	6.7	6.7
Another way?	3.8	2.7	3.5	3.3
Why yours?	6.8	7.0	7.2	7.0
How to get there?	7.7	7.8	7.8	7.8

were chosen 8.2% more often, and EXPLOREP’s plans 5.4% less often. No plan from NOVELP was ever selected because its plans typically performed poorly in the four affordance-based graphs. Voting, however, included NOVELP to preserve a potential trade-off between exploration and exploitation.

We evaluated WHY for its efficiency (average computation time) and diversity (number of unique explanations produced in response to each question). We also calculated the understandability of these explanations by average reading grade level, as measured by the Coleman-Liau index (CLI) (Coleman and Liau, 1975). Since WHY’s goal is to produce explanations for non-experts, lower grade-level scores are more desirable. While one could manipulate the templates to improve these scores, CLI provides a method to compare the complexity of responses to one another.

Table 7 analyzes WHY’s answers to all 3000 (5 questions · 40 targets · 5 sequences · 3 worlds) questions. Its distinct natural explanations simulate people’s ability to vary explanations based on context (Malle, 1999). WHY averaged 10.4 msec to compute explanations for all five questions about each plan. WHY’s approach is also nuanced, with many unique responses per question. For example, WHY produced 49 unique responses to “Why does your plan go this way?” out of the 92 possible instantiations of the template. The CLI gauged them at about a sixth-grade reading level, readily understandable to a layperson.

8 Discussion

To capture useful spatial affordances for its world model, SemaFORR generalizes over its percepts,

the 660 distances to the nearest obstacle that its range finder reports 15 times per second. Each of SemaFORR’s planners generates paths in a graph biased by edge weights that represent that planner’s objective but share an underlying structure that facilitates plan comparison. Voting guarantees that any selected plan will be optimal with respect to at least one objective, and makes it likely that the plan will also perform well with respect to the others. This also facilitates contrastive explanations in natural spatial language for the robot’s planning objectives, alternative paths, and confidence.

How a robot control system represents knowledge is integral to natural communication between robots and people, especially in a spatial context. Misunderstandings between a robot and a human often arise from a discrepancy between their spatial mental models. This prompts questions about the robot’s underlying decision-making and reasoning mechanisms. WHY’s explanations rely on SemaFORR’s cognitive underpinnings. Language about the spatial model is readily understood because SemaFORR interprets its percepts much the way people do. SemaFORR’s freespace affordances were inspired by sketches after human subjects had actively explored complex virtual worlds (Chrastil and Warren, 2013). The planners’ objectives are also analogous to processes empirically identified in people (Hölscher et al., 2009). The results here demonstrate that natural language communication with robots benefits substantially when a robot’s control system and a human have similar cognitively-based spatial representations.

WHY’s templates flexibly and quickly produce many different explanations in natural language. The templates focus language generation on SemaFORR’s computational rationale rather than on linguistic structure and grammar. They also facilitate the introduction of new planners without the need to retrain a language generator for a new planning objective. For example, an objective that relied on landmarks could modify the cost graph to reduce costs near them, so that WHY might explain “I think my way is a lot better at following landmarks.” Although WHY assumes the human’s objective is the shortest path, it can easily substitute any objective representable in a cost graph with an admissible heuristic. SemaFORR could also incorporate a planning objective learned from external demonstration (e.g., inverse reinforcement learning) if that objective were representable as

increments to the cost graph’s weights.

Whenever SemaFORR selects FASTP’s plan here, it assumes that it shares the human’s objective. Any questions about the robot’s plan necessarily challenge that assumption. Presumably, the person asks because they do not recognize their objective there. WHY responds by agreement that the person’s plan is the correct way to go (e.g., “Actually, I agree that we should take the shortest route.”), even though the question should not have arisen. Another way to address this would be to offer an alternative plan when FASTP is selected.

Our current work examines how well human subjects understand and feel comfortable with WHY. Although SemaFORR’s parameters for intervals (e.g., in Table 3) were chosen for G5 and also worked well in other worlds, human subject evaluation will allow us to confirm or reassess these values. Human-subject studies could also help refine WHY’s explanations and incorporate psychophysics and proxemics.

Future work could extend WHY for dialogue (e.g., to clarify confusion or guide navigation (Roman et al., 2020)). This could incorporate natural language generation with deep learning and facilitate queries to the person. WHY presumes that questions arise from a difference between the human’s and the robot’s objectives, but they could also stem from a violation of the shared target assumption. A broader system for human-robot collaboration would seek the cause of such a mismatch, use plan explanations to resolve it, and then allow the robot to adjust its responses based on feedback from its human partner. For example, given a plan P from a person or an unspecified heuristic planner, WHY could use the individual objectives in its repertoire to tease apart and then characterize how P weighted its objectives (e.g., “So distance is more important than travel time?”).

Meanwhile, SemaFORR’s cognitively-based spatial model supports important path planning objectives and human-friendly explanations of its behavior, intentions, and confidence. Empirical results in three large, complex, realistic worlds show that our approach produces diverse, understandable contrastive explanations in natural language.

Acknowledgments

This work was supported in part by The National Science Foundation under CNS-1625843. The authors thank Anoop Aroor for MengeROS.

References

- Anoop Aroor, Susan L Epstein, and Raj Korpan. 2017. [MengeROS: A Crowd Simulation Tool for Autonomous Robot Navigation](#). In *Proceedings of AAAI Fall Symposium on Artificial Intelligence for Human-Robot Interaction*, pages 123–125. AAAI.
- Daniel Paul Barrett, Scott Alan Bronikowski, Haonan Yu, and Jeffrey Mark Siskind. 2017. [Driving Under the Influence \(of Language\)](#). *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–16.
- Tathagata Chakraborti, Sarath Sreedharan, Sachin Grover, and Subbarao Kambhampati. 2019. [Plan explanations as model reconciliation—an empirical study](#). In *2019 14th ACM/IEEE International Conference on Human-Robot Interaction (HRI)*, pages 258–266. IEEE.
- Elizabeth R Chrastil and William H Warren. 2013. [Active and passive spatial learning in human navigation: Acquisition of survey knowledge](#). *Journal of experimental psychology: learning, memory, and cognition*, 39(5):1520.
- Elizabeth R Chrastil and William H Warren. 2014. [From cognitive maps to cognitive graphs](#). *PLoS one*, 9(11):e112544.
- Meri Coleman and Ta Lin Liao. 1975. [A Computer Readability Formula Designed for Machine Scoring](#). *Journal of Applied Psychology*, 60(2):283–284.
- Susan L Epstein, Anoop Aroor, Matthew Evanusa, Elizabeth I Sklar, and Simon Parsons. 2015. [Learning spatial models for navigation](#). In *International Conference on Spatial Information Theory*, pages 403–425. Springer.
- Susan L. Epstein and Raj Korpan. 2019. [Planning and explanations with a learned spatial model](#). In *International Conference on Spatial Information Theory*, volume 142 of *LIPICS*, pages 22:1–22:20. Schloss Dagstuhl.
- Brett R Fajen and Flip Phillips. 2013. [Spatial perception and action](#). In *Handbook of spatial cognition*. American Psychological Association.
- Patrick Foo, William H Warren, Andrew Duchon, and Michael J Tarr. 2005. [Do humans integrate routes into a cognitive map? map- versus landmark-based navigation of novel shortcuts](#). *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 31(2):195–215.
- Maria Fox, Derek Long, and Daniele Magazzeni. 2017. [Explainable planning](#). In *IJCAI-17 Workshop on Explainable AI (XAI)*, page 24.
- Charles R Gallistel. 1990. *The organization of learning*. The MIT Press.
- James J Gibson. 1977. [The theory of affordances](#). *Perceiving, Acting, and Knowing: Toward an Ecological Psychology*, pages 67–82.
- Reginald G Golledge. 1999. Human wayfinding and cognitive maps. *Wayfinding behavior: Cognitive mapping and other spatial processes*, pages 5–45.
- Antoine Grea, Laëtitia Matignon, and Samir Aknine. 2018. [How explainable plans can make planning faster](#). In *Workshop on Explainable Artificial Intelligence*, pages 58–64.
- P. E. Hart, N. J. Nilsson, and B. Raphael. 1968. [A formal basis for the heuristic determination of minimum cost paths](#). *IEEE Transactions on Systems Science and Cybernetics*, 4(2):100–107.
- Jörg Hoffmann and Daniele Magazzeni. 2019. [Explainable AI planning \(XAIP\): overview and the case of contrastive explanation](#). *Reasoning Web. Explainable Artificial Intelligence*, pages 277–282.
- Christoph Hölscher, Simon J Büchner, Tobias Meilinger, and Gerhard Strube. 2009. [Adaptivity of wayfinding strategies in a multi-building ensemble: The effects of spatial structure, task requirements, and metric information](#). *Journal of Environmental Psychology*, 29(2):208–219.
- Sandy H Huang, David Held, Pieter Abbeel, and Anca D Dragan. 2019. [Enabling robots to communicate their objectives](#). *Autonomous Robots*, 43(2):309–326.
- Ioannis Kostavelis and Antonios Gasteratos. 2015. [Semantic mapping for mobile robotics tasks: A survey](#). *Robotics and Autonomous Systems*, 66:86–103.
- Benjamin Krarup, Michael Cashmore, Daniele Magazzeni, and Tim Miller. 2019. [Model-based contrastive explanations for explainable planning](#). In *ICAPS 2019 Workshop on Explainable AI Planning (XAIP)*.
- Benjamin Kuipers. 2000. [The spatial semantic hierarchy](#). *Artificial intelligence*, 119(1-2):191–233.
- Christian Landsiedel, Verena Rieser, Matthew Walter, and Dirk Wollherr. 2017. [A Review of Spatial Reasoning and Interaction for Real-World Robotics](#). *Advanced Robotics*, 31(5):222–242.
- Maurício Cecílio Magnaguagno, Ramon Fraga Pereira, Martin Duarte Móre, and Felipe Rech Meneguzzi. 2017. [Web planner: A tool to develop classical planning domains and visualize heuristic state-space search](#). In *2017 Workshop on User Interfaces and Scheduling and Planning*.
- Bertram F Malle. 1999. [How People Explain Behavior: A New Theoretical Framework](#). *Personality and Social Psychology Review*, 3(1):23–48.
- Matthew Marge, Carol Espy-Wilson, and Nigel Ward. 2020. [Spoken language interaction with robots: Research issues and recommendations](#). *Report from the NSF Future Directions Workshop*.

- Raija Markkanen and Hartmut Schröder. 1997. *Hedging and Discourse: Approaches to the Analysis of a Pragmatic Phenomenon in Academic Texts*, volume 24. Walter de Gruyter.
- Tim Miller. 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artificial intelligence*, 267:1–38.
- Van Nguyen, Stylianos Loukas Vasileiou, Tran Cao Son, and William Yeoh. 2020. Explainable planning using answer set programming. In *Proceedings of the International Conference on Principles of Knowledge Representation and Reasoning*, volume 17, pages 662–666.
- Brandon S Perelman, Arthur W Evans III, and Kristin E Schaefer. 2020. Where do you think you’re going? characterizing spatial mental models from planned routes. *ACM Transactions on Human-Robot Interaction (THRI)*, 9(4):1–55.
- Vittorio Perera, Sai P Selveraj, Stephanie Rosenthal, and Manuela Veloso. 2016. Dynamic Generation and Refinement of Robot Verbalization. In *25th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*, pages 212–218. IEEE.
- Thomas Pouncy, Pedro Tsividis, and Samuel J Gershman. 2021. What is the model in model-based planning? *Cognitive Science*, 45(1):e12928.
- Morgan Quigley, Ken Conley, Brian P. Gerkey, Josh Faust, Tully Foote, Jeremy Leibs, Rob Wheeler, and Andrew Y. Ng. 2009. ROS: an open-source Robot Operating System. In *ICRA Workshop on Open Source Software*.
- Homero Roman Roman, Yonatan Bisk, Jesse Thomason, Asli Celikyilmaz, and Jianfeng Gao. 2020. Rmm: A recursive mental model for dialog navigation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 1732–1745.
- Avi Rosenfeld and Ariella Richardson. 2019. Explainability in human-agent systems. *Autonomous Agents and Multi-Agent Systems*, 33(6):673–705.
- Stephanie Rosenthal, Sai P Selvaraj, and Manuela Veloso. 2016. Verbalization: Narration of Autonomous Mobile Robot Experience. In *Proceedings of the International Joint Conference on Artificial Intelligence*, volume 16, pages 862–868.
- Rosario Scalise, Stephanie Rosenthal, and Siddhartha Srinivasa. 2017. Natural Language Explanations in Human-Collaborative Systems. In *Proceedings of the 2017 ACM/IEEE International Conference on Human-Robot Interaction*, pages 377–378. ACM.
- Bastian Seegebarth, Felix Müller, Bernd Schattenberg, and Susanne Biundo. 2012. Making hybrid plans more clear to human users—a formal approach for generating sound explanations. In *Proceedings of the International Conference on Automated Planning and Scheduling*, volume 22.
- Avinash Kumar Singh, Neha Baranwal, Kai-Florian Richter, Thomas Hellström, and Suna Bensch. 2021. Verbal explanations by collaborating robot teams. *Paladyn, Journal of Behavioral Robotics*, 12(1):47–57.
- Roykronk Sukkerd, Reid Simmons, and David Garlan. 2020. Tradeoff-focused contrastive explanation for MDP planning. In *2020 29th IEEE International Conference on Robot and Human Interactive Communication*, pages 1041–1048. IEEE.
- Barbara Tversky. 1993. Cognitive maps, cognitive collages, and spatial mental models. In *European Conference on Spatial Information Theory*, pages 14–24. Springer.
- William H Warren, Daniel B Rothman, Benjamin H Schnapp, and Jonathan D Ericson. 2017. Wormholes in virtual space: From cognitive maps to cognitive graphs. *Cognition*, 166:152–163.