# Commonsense Spatial Reasoning for Visually Intelligent Agents

**Agnese Chiatti**[1] , **Gianluca Bardaro**[1] , **Enrico Motta**[1] , **Enrico Daga**[1]

[1]Knowledge Media Institute, The Open University, United Kingdom

{name.surname}@open.ac.uk

## Abstract

Service robots are expected to reliably make sense of complex, fast-changing environments. From a cognitive standpoint, they need the appropriate reasoning capabilities and background knowledge required to exhibit human-like Visual Intelligence. In particular, our prior work has shown that the ability to reason about spatial relations between objects in the world is a key requirement for the development of Visually Intelligent Agents. In this paper, we present a framework for commonsense spatial reasoning which is tailored to real-world robotic applications. Differently from prior approaches to qualitative spatial reasoning, the proposed framework is robust to variations in the robot's viewpoint and object orientation. The spatial relations in the proposed framework are also mapped to the types of commonsense predicates used to describe typical object configurations in English. In addition, we also show how this formally-defined framework can be implemented in a concrete spatial database.

## 1 Introduction

In all cases where it is inconvenient or even dangerous for us to intervene, there is an incentive to delegate tasks to *service robots* - or *robot assistants*: e.g., under the extreme conditions imposed by space explorations (Nilsson et al. 2018), in hazardous manufacturing environments (Liu and Wang 2020), or whenever social distance needs to be maintained (Yang et al. 2020). Before delegating complex tasks to robots, however, we need to ensure that they can reliably *make sense* of the stimuli coming from their sensors. Autonomous sensemaking remains an open challenge, because it requires not only to reconcile the high-volume and diverse data collected from real-world settings (Alatise and Hancke 2020), but also to actually understand these data, going beyond mere pattern recognition (Lake et al. 2017; Davis and Marcus 2015).

From a vision perspective, the problem of robot sensemaking becomes one of enhancing the *Visual Intelligence* of service robots, i.e., their ability to make sense of the environment through their vision system and epistemic competences (Chiatti, Motta, and Daga 2020). Naturally, several epistemic competences are required to build *Visually Intelligent Agents (VIA)*. For instance, let us consider the case of HanS, a Health & Safety robot inspector. HanS is expected to autonomously detect potentially threatening situations, such as the fire hazard posed by a sweater left to dry

on top of an electric heater. To assess the risk associated with this situation, HanS first needs to recognise the sweater and the heater in question, i.e., it needs to exhibit robust *object recognition* capabilities. It also needs *spatial reasoning* capabilities, to infer that the sweater is touching the heater. Moreover, it also needs to know that sweaters are made of cloth and that a piece of cloth clogging an electric radiator can catch fire. The list goes on.

In (Chiatti, Motta, and Daga 2020), we identified a framework of *epistemic requirements*, i.e., knowledge properties and reasoning capabilities which are needed to develop *Visually Intelligent Agents (VIA)*. To form hypotheses on which epistemic requirements are more likely to significantly enhance the Visual Intelligence of a robot, we also mapped these epistemic ingredients to the types of object classification errors emerging from one of HanS' scouting routines. This error analysis highlighted that the majority of misclassifications could in principle have been avoided, if the robot was capable of considering: (i) the canonical size of objects, e.g., that mugs are generally smaller than bins, as well as (ii) the typical *Qualitative Spatial Relations (QSR)* between objects. For instance, a fire extinguisher may be mistaken for a bottle due to its shape. However, the proximity of a fire extinguisher sign is a strong indication that the observed object is in fact a fire extinguisher. This element of *typicality* relates to the broader objective of developing AI systems which can reason about what is *plausible* (Davis and Marcus 2015), i.e., which exhibit *common sense* (Levesque 2017) and *Intuitive Physics* reasoning abilities (Hayes 1988; Lake et al. 2017). Our most recent findings (Chiatti et al. 2021) confirmed that combining state-of-the-art Machine Learning methods with a component able to reason about object sizes improves the robot's object recognition performance. In this paper, we progress this line of research by characterising commonsense QSR between objects.

The problem of representing spatial relations has been actively researched for decades, producing many theoretical frameworks for autonomous spatial reasoning (Cohn and Renz 2008). In robotics, *semantic mapping* (Nüchter and Hertzberg 2008; Kostavelis and Gasteratos 2015) and *object anchoring* methods (Coradeschi and Saffiotti 2003) have enabled linking the robot sensor data and symbolic knowledge to the geometric maps modelling its environment. To combine the best of both worlds, a number

of approaches (Deeken, Wiemann, and Hertzberg 2018; Kunze et al. 2014; Young et al. 2017) have linked the spatial representations within semantic maps to the higher-level formal definitions provided by AI theories. In this paper, we propose a novel spatial reasoning framework which extends the work in (Deeken, Wiemann, and Hertzberg 2018; Borrmann and Rank 2010), to account for variations in the robot's viewpoint and in the relative orientation of objects. Moreover, we formally map the Qualitative Spatial Relations composing this framework to the type of linguistic predicates used to describe *commonsense spatial relations* in English, which are discussed within seminal theories of spatial cognition (Landau and Jackendoff 1993; Herskovits 1997). Finally, we show how the proposed framework can be implemented in state-of-the-art Geographic Information Systems (GIS), to support commonsense spatial reasoning in real-world robotic scenarios.

## 2 Related work

Broadly speaking, spatial relations can be represented qualitatively - e.g., A contains B - or quantitatively - e.g., the angle between A and B is $\theta$ (Thippur et al. 2015). Following (Borrmann and Rank 2010), Qualitative Spatial Relations (QSR) can be further characterised as (i) *metric*, i.e., based on the metric distance between objects (ii) *topological*, i.e., describing the neighbourhood of objects, and (iii) *directional*, i.e., relative to the axis directions in a reference coordinate system. The interested reader is referred to (Cohn and Renz 2008) for a foundational review of qualitative spatial representations. Compared to quantitative representations, qualitative representations are more similar to the types of spatial predicates involved in natural language discourse. As a result, qualitative spatial representations are easier to interpret and aid Human-Robot Interaction (Sarthou, Alami, and Clodic 2019; Sisbot and Connell 2019; Thippur, Stork, and Jensfelt 2017). Moreover, they are more similar to the types of linguistic predicates available within large-scale, general-purpose Knowledge Bases (KB), such as those surveyed in (Storks, Gao, and Chai 2019). Crucially, they are also aligned with the spatial predicates provided with benchmark image collections for visual reasoning tasks, such as Visual Genome (Krishna et al. 2017) and SpatialSense (Yang, Russakovsky, and Deng 2019). Thus, relying on qualitative representations has the potential to facilitate the repurposing of these resources in robotic contexts, especially given the paucity of comprehensive KBs for Visually Intelligent Agents (Chiatti, Motta, and Daga 2020).

Extensive efforts have been devoted to mapping the quantitative data collected through the robot's sensors to higher-level symbols describing a set of known object classes and their attributes (Nüchter and Hertzberg 2008; Coradeschi and Saffiotti 2003; Kostavelis and Gasteratos 2015). These efforts have produced intermediate representational models also known as *semantic maps*, i.e., maps that contain, "in addition to spatial information about the environment, assignments of mapped features to entities of known classes" (Nüchter and Hertzberg 2008). Further approaches have been proposed, where the content of semantic maps is also interpreted with respect to formal theories of qualitative spatial reasoning (Young et al. 2017; Kunze et al. 2014; Deeken, Wiemann, and Hertzberg 2018). In general, spatial relations are expressed between object pairs, where one of the two objects is considered as a *reference*, or *landmark*: e.g., bike near house. (Young et al. 2017) have used Ring Calculus to represent the closeness of objects. (Kunze et al. 2014) have relied on ternary point calculus (Moratz and Ragni 2008) to model directional relations with respect to both the robot's location and the location of the reference object. Thus, the 3D regions occupied by objects are reduced to point-like objects on the 2D plane. Moreover, (Kunze et al. 2014) assumed that the robot's location does not change over time, and is always defined with respect to a tabletop. Differently from (Kunze et al. 2014), (Deeken, Wiemann, and Hertzberg 2018) represented directional relations by comparing the 3D regions occupied by objects, through the halfspace-based model of (Borrmann and Rank 2010). However, this model is based on the assumption that the robot's viewpoint is always aligned both with the global coordinate system of the map and with the inherent orientation of the observed objects. Thus, it is not suitable to model mobile robots making sense of the environment during navigation. In real-world scenarios, as the robot moves, its viewpoint changes over time and the objects observed will be oriented differently. Thus, we propose to combine the robot's viewpoint and the orientation of the reference object within a *contextualised frame of reference*. This contextualised frame of reference allows us to define a contextualised 3D region, or *Contextualised Bounding Box*, which represents the location of the object with respect to both the robot's viewpoint and the frame of reference of a landmark. Crucially, the contextualised frame of reference and Bounding Box can be defined for any combination of robot and landmark location, thus ensuring that this framework can scale to many real-world robotic scenarios.

## 3 Proposed Framework

To define a spatial reasoning framework which satisfies the concrete requirements of robot sensemaking, we extend the formal theory of spatial reasoning by (Borrmann and Rank 2010). Moreover, we map the obtained spatial relations to the commonsense predicates used to describe spatial relations between objects in English. These predicates are gathered from cognitive theories (Landau and Jackendoff 1993). By making an explicit link between formal AI theories and informal linguistic representations, we obtain a framework for commonsense spatial reasoning in robotic scenarios.

**Notation** In what follows, we model definitions as First Order Logic (FOL) statements. We represent logic variables through lowercase letters and constants through uppercase letters. We also use lowercase initials to denote functions, while uppercase initials symbolise predicates. For instance, *sReg* is a function, whereas *Above* is a predicate. Unless otherwise stated, free variables are universally quantified. Finally, we use the standard notation (X, Y, Z) to denote reference axes, while $x, y, z$ are used to refer to the spatial coordinates with respect to those axes.

**Spatial primitives**  Our domain of discourse $\mathbb{D}$ is that of *spatial objects*, i.e., physical objects, "which have spatial extensions" (Cohn and Renz 2008). From this perspective, a spatial object is represented in terms of the associated *spatial region*. In particular, our spatial primitive is the concept of *spatial point*. Thus, spatial regions are represented as sets of spatial points, $p$. Let $P$ be the set of all spatial points, then, for each spatial object $o \in \mathbb{D}$, we assume the existence of a function $sReg$ which, given $o$, returns the subset of $P$ which includes all the points in the spatial region of $o$.

$$SpatialObj(o) \Rightarrow sReg(o) \subseteq P \tag{1}$$

$$SpatialObj(o) \Rightarrow sReg(o) \neq \varnothing \tag{2}$$

In particular, our focus is not on arbitrary collections of spatial points, but rather on one-piece regions (Cohn and Renz 2008), i.e., on sets of internally connected points:

$$SpatialObj(o) \Rightarrow ProperSR(sReg(o)) \tag{3}$$

To provide a formal definition of the concept of *proper spatial region*, we need first to establish a *spatial frame of reference*.

**Spatial Frame of Reference**  A spatial object is characterised not only with respect to a spatial region but also in terms of a reference coordinate system, also known as *frame of reference*. A frame of reference consists of an origin point and of a set of directed axes intersecting at the origin, $O$. In particular, modelling the 3D space requires three reference axes $X, Y, Z$. Although spatial points and spatial regions exist independently of the frame of reference, the interpretation of these spatial primitives only makes sense in the context of a frame of reference. Once we have defined a reference frame, we can interpret spatial points as *geometrical points*, i.e., as coordinate triples in $\mathbb{R}^3$. Let $GP$ be the set of all geometrical points in the considered space:

$$GP = \{p | p = (x, y, z) \in \mathbb{R}^3\} \tag{4}$$

The identified frame of reference also has an associated *granularity*, i.e., an infinitesimally small constant $D > 0$ in $\mathbb{R}$, which defines the minimum distance for two geometrical points to be considered as distinct entities. Two geometrical points are then said to be *adjacent* iff their geometrical distance is equal to $D$. To compute the distance between two geometrical points, they have to be in the same frame of reference. Let $d(gp, gp')$ be a function which returns a real number indicating the geometric distance between points $gp$ and $gp'$. Then:

$$Adj(gp, gp') \Leftrightarrow d(gp, gp') = D \tag{5}$$

The definition of proper spatial region then follows from the defined notion of adjacency:

$$ProperSR(sr) \Leftrightarrow \forall gp[gp \in sr \Rightarrow Conn(gp, sr)] \tag{6}$$

$$Conn(gp, sr) \Leftrightarrow \forall gp'[gp' \in sr \wedge$$
$$gp' \neq gp] \Rightarrow ConnP(gp, gp') \tag{7}$$

$$ConnP(gp_1, gp_2) \Leftrightarrow Adj(gp_1, gp_2) \vee$$
$$\exists gp_3[Adj(gp_1, gp_3) \wedge ConnP(gp_3, gp_2)] \tag{8}$$

In our model, we assume that the *global spatial region*, $GP$, is a fully-connected set of points. Moreover, we assume that spatial regions can be approximated through 3D boxes[1]. This simplifying assumption is consistent with standard practice in the literature (Deeken, Wiemann, and Hertzberg 2018; Borrmann and Rank 2010). Bounding boxes can have an arbitrary orientation around the Z axis aligned with gravity, but their base is always parallel to the XY plane. In particular, we consider the minimum bounding box which best approximates the real volume occupied by an object and which is aligned with its *natural orientation* (Chiatti, Motta, and Daga 2020). Let $b$ be a set of geometrical points which contains the spatial region of $o$:

$$BoundBox(b, o) \Leftrightarrow sReg(o) \subseteq b \wedge b \subseteq GP \tag{9}$$

$$MinBoundBox(b, o) \Leftrightarrow BoundBox(b, o) \wedge$$
$$\neg \exists b'[BoundBox(b', o) \wedge b' \subset b] \tag{10}$$

In this scenario, the environment navigated by a robot can also be modelled as a spatial region including an arbitrary number of objects, i.e., as a global spatial region. Consequently, the outer region of a spatial region, $sr$, is:

$$outReg(sr) = \{gp | gp \in GP \wedge gp \notin sr\} \tag{11}$$

The frame of reference of the global region, $F_g$, is *extrinsic*, i.e., based on a reference point which is external to both an object and an observer. $F_g$ remains fixed as the robot navigates the environment. Conversely, the robot's frame of reference, $F_r$, changes as the robot moves. Thus, it is *deitic*, relative to the observer's position. Consequently, the location of objects at each point in time can be interpreted differently, based on which frame of reference is considered. Within the formal spatial reasoning frameworks of (Borrmann and Rank 2010; Deeken, Wiemann, and Hertzberg 2018), all the spatial relations between objects are defined according to the same pre-defined frame of reference, whether it is an extrinsic, deitic or intrinsic one, i.e., inherent to a specific object.

Differently from the latter relations, linguistic spatial predicates implicitly refer both to (i) the location of whichever object is considered as reference, and to (ii) the observer's point of view (Landau and Jackendoff 1993). Similarly, a robot would conclude that "A is on the left of B" based not only on the location of objects A and B within $F_g$, but also on $F_r$. From a different standpoint, A might appear on the right of B, for instance, or in front of it. To model such cases, we introduce the notion of *robot's viewpoint*, $F_{r'}$. Let $C_o$ be the centroid of the spatial region representing object $o$. Then, $F_{r'}$ is obtained by rotating $F_r$ along $Z_r$, by an angle $\alpha$. Specifically, $\alpha$ is the angle between $X_r$ and the imaginary line connecting the origin of $F_r$ with $C_o$. Let $F_o$ of origin $C_o$ and axes $X_o, Y_o, Z_o$ be the intrinsic frame of reference of $o$, i.e., the frame of reference which is aligned with the orientation of $o$. Then, the contextualised frame of reference of the object, $F_c$, is the frame of reference of origin $C_o$

---

[1] In the case of large regions of negligible thickness, such as floors, walls and ceilings, the spatial region reduces to a 2D surface.

whose axes have the same orientation of the axes defining the robot's viewpoint, $F_{r'}$ (Figure 1).

Based on $F_c$, we can construct a *Contextualised Bounding Box (CBB)*, which is obtained by aligning the minimum bounding box with $F_c$. Let $rotZ(b, \theta)$ be a function which returns the spatial region, $sr$, obtained by rotating an input bounding box along $Z$ by an angle $\theta$. Given a frame of reference $F_c$, then $yaw(sr, F_c)$ returns the angle between the intrinsic frame of reference of $sr$ and $F_c$, along $Z$. Then, given $\pi/2$:

$$IsCBB(rotZ(b, \theta), o) \Leftrightarrow MinBoundBox(b, o) \wedge$$
$$\exists\theta[mod(yaw(rotZ(b, \theta), F_c), \pi/2) = 0 \wedge \neg\exists\theta'$$
$$[mod(yaw(rotZ(b, \theta'), F_c), \pi/2) = 0 \wedge \theta' < \theta]] \quad (12)$$

Namely, to construct CBB, we select the minimum angle $\theta$ so that the value returned by the *yaw* function is divisible by $\pi/2$, i.e., the remainder of their division, *mod*, is zero. There are always four possible alignments of a bounding box, $b$, for which *mod* is zero. Thus, by selecting the minimum angle among these four, we apply the transformation which is least disruptive of the natural orientation of the object. Thanks to these newly-defined spatial concepts, we can now map the metric, topological and directional relations in (Borrmann and Rank 2010; Deeken, Wiemann, and Hertzberg 2018) to commonsense predicates expressed in natural language.

**Metric spatial relations** Given two spatial objects $o_1, o_2$ and two geometrical points $gp_1, gp_2$ where $gp_1 \in o_1$ and $gp_2 \in o_2$, we define the distance between two geometrical points as the their Euclidean distance:

$$d(gp_1, gp_2) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2} \quad (13)$$

Then, the distance between two spatial objects is defined as the global minimum of the pointwise distance function, $d$:

$$[distance(o_1, o_2) = d(gp_1, gp_2)] \Leftrightarrow gp_1 \in o_1 \wedge gp_2 \in o_2 \wedge$$
$$\forall gp_3, gp_4[gp_3 \in o_1 \wedge gp_4 \in o_2] \Rightarrow d(gp_3, gp_4) >$$
$$d(gp_1, gp_2) \quad (14)$$

A distance threshold, $T$, can be then introduced, to represent closeness between objects. That is, for a $T$ greater than or equal to the frame granularity $D$ defined earlier:

$$IsClose(o_1, o_2) \Leftrightarrow distance(o_1, o_2) \leqslant T \quad (15)$$

In particular, if the minimum distance between two objects equals $D$, then the two objects touch:

$$Touches(o_1, o_2) \Leftrightarrow distance(o_1, o_2) = D \quad (16)$$

**Topological spatial relations** Topological relations are spatial relations which are invariant under a topological isomorphism, i.e., a function $f : X \rightarrow Y$ which preserves neighbourhood relationships while mapping $X$ to $Y$. Although several different qualitative representations of topological relations have been proposed (Cohn and Renz 2008), here we focus on a subset of topological relations, namely on the intersection and containment relations. As shown in the remainder of this Section, this minimal subset of relations, combined with metric and directional relations, is sufficient

to cover all the commonsense spatial relations required in the scenario of interest. First, based on our prior definitions, two spatial regions, $sr, sr'$ intersect iff they have at least one geometrical point in common:

$$Int(sr, sr') \Leftrightarrow \exists gp[gp \in sr \wedge gp \in sr'] \quad (17)$$

We also define the spatial region representing the intersection between two objects (i.e., the intersection between the associated spatial regions) as follows:

$$inter(o1, o2) = \{gp | gp \in sReg(o1) \wedge gp \in sReg(o2)\} \quad (18)$$

Then, a special case of the intersection relation is the case where one spatial region completely contains the other:

$$ComplCont(sr, sr') \Leftrightarrow \forall gp[gp \in sr' \Rightarrow gp \in sr] \quad (19)$$

Semantically, $o$ contains $o'$ completely iff all the geometrical points in the spatial region of $o'$ are also members of the spatial region of $o$.

**Directional spatial relations** Differently from metric and topological relations, directional spatial relations are interpreted differently based on the considered frame of reference. (Borrmann and Rank 2010) have proposed a qualitative representation for directional relations where the region outside a 3D bounding box is partitioned into six halfspaces, i.e., one halfspace for each semi-axis of $X, Y, Z$. Because we have defined a global spatial region containing all geometrical points in the robot's environment, these halfspaces are also proper spatial regions, which can be approximated through 3D bounding boxes. In particular, as in (Deeken, Wiemann, and Hertzberg 2018), they can be modelled as 3D extrusions, obtained by multiplying the extent of the object spatial region by a scaling factor $s \in \mathbb{R}$.

The coordinates of all geometrical points in the minimum bounding box are bound to a minimum and maximum value, e.g., $x_{min}$ and $x_{max}$. Let $X_o^+$ and $X_o^-$ be the positive and negative semi-axes of $X_o$ in $F_o$. Then, we define a function, $hs$, which returns the halfspace of an input bounding box, given semi-axis, $X_o^+$, and frame of reference, $F_o$:

$$MinBoundBox(mb_1, o_1) \Rightarrow hs(mb_1, X_o^+, F_o) =$$
$$\{gp \in outReg(mb_1) | gp = (x, y, z) \text{ w.r.t } F_o,$$
$$x_{max} \leqslant x \leqslant x_{max} + x_{max} \cdot s,$$
$$y_{min} \leqslant y \leqslant y_{max},$$
$$z_{min} \leqslant z \leqslant z_{max}\}$$
$$(20)$$

Additional halfspaces can be similarly derived for the other semi-axes in $F_o$, as further documented in the supplementary materials. Once these halfspaces have been defined, one can test whether a second object $o_2$ lies within any of the halfspaces of $o_1$. In particular, (Borrmann and Rank 2010) differentiate between "relaxed" (_r) and "strict" (_s) spatial operators, based on whether $o_2$ intersects or is completely contained in the halfspaces of $o_1$. In the following, we represent predicates symbolising cardinal directions *East*, *West*, *North*, *South*, *Above* and *Below* through their capital initial.
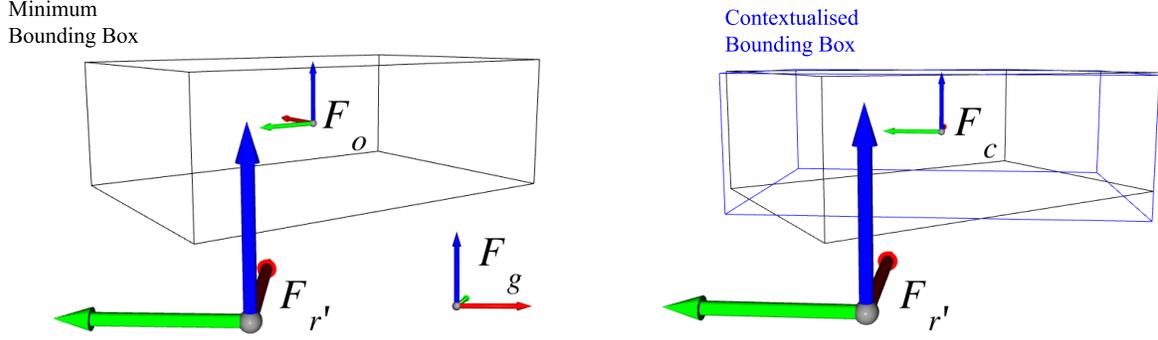
Figure 1: The robot's viewpoint, $F_{r'}$, consists of an origin and three axes $X_{r'}$ (in red), $Y_{r'}$ (in green) and $Z_{r'}$ (in blue). $F_{r'}$ may not coincide with the frame of reference characterising the global map, $F_g$, nor with the intrinsic frame of reference of a certain object, $F_o$. As shown on the left-hand side of the Figure, a spatial object is first modelled as the minimum 3D box bounding the object. Then, $F_{r'}$ is translated to the object's centroid to define a contextualised frame of reference $F_c$. Moreover, a Contextualised Bounding Box (highlighted in blue) is generated, i.e., the bounding box which requires the minimum rotation along the $Z$ axis to align the minimum bounding box with $F_c$.

Given a $F_o$ which coincides with $F_g$, the strict and relaxed definitions of the relation $East(o_2, o_1)$ are:

$$E\_s(o_2, o_1, F_o) \Leftrightarrow MinBoundBox(mb_1, o_1) \wedge$$
$$ComplCont(hs(mb_1, X_o^+, F_o), sReg(o_2)) \quad (21)$$

$$E\_r(o_2, o_1, F_o) \Leftrightarrow MinBoundBox(mb_1, o_1) \wedge$$
$$Int(hs(mb_1, X_o^+, F_o), sReg(o_2)) \quad (22)$$

Based on our prior definitions, we can model directional relations with respect to a given $F_o$ as follows:

$$W\_r(o_2, o_1, F_o) \Leftrightarrow Int(sReg(o_2), hs(mb_1, X_o^-, F_o)) \quad (23)$$

$$N\_r(o_2, o_1, F_o) \Leftrightarrow Int(sReg(o_2), hs(mb_1, Y_o^+, F_o)) \quad (24)$$

$$S\_r(o_2, o_1, F_o) \Leftrightarrow Int(sReg(o_2), hs(mb_1, Y_o^-, F_o)) \quad (25)$$

$$A\_r(o_2, o_1, F_o) \Leftrightarrow Int(sReg(o_2), hs(mb_1, Z_o^+, F_o)) \quad (26)$$

$$B\_r(o_2, o_1, F_o) \Leftrightarrow Int(sReg(o_2), hs(mb_1, Z_o^-, F_o)) \quad (27)$$

For brevity, we have omitted the predicate $MinBoundBox(mb_1, o_1)$ from axioms 23-27. The full definition is in the supplementary materials.

(Borrmann and Rank 2010) defined the aforementioned relations under the assumption that $F_o$ is always aligned with $F_g$. However, this assumption does not hold in the case of mobile robot sensemaking. Indeed, the frame of reference of the robot, $F_r$ is mobile, i.e., its origin and orientation change over time. Moreover, the natural orientation of objects may not be aligned with $F_g$. Thus, to produce a representational model which suits the case of robot sensemaking, we need to map axioms 22-27 to the contextualised frame of reference, $F_c$, defined earlier. In a typical robotic setting, the robot always faces towards $X_r^+$, and $Z_r^+$ is directed upwards, i.e., opposite to the direction of gravity. Then, the orientation of $Y_r$ is given by applying the right hand rule (Figure 1). Based on these premises, all $Z$ axes always share the same orientation. Namely, the top and bottom halfspaces of an object do not change with the robot's reference frame (Figures 1,2). Therefore, the $A\_r$ and $B\_r$ predicates based on the minimum oriented bounding box of an object w.r.t. a given $F_o$ can be directly reused to define the *Above* and *Below* relations w.r.t. $F_c$:

$$Above(o_2, o_1, F_c) \Leftrightarrow A\_r(o_2, o_1, F_o) \quad (28)$$

$$Below(o_2, o_1, F_c) \Leftrightarrow B\_r(o_2, o_1, F_o) \quad (29)$$

Nonetheless, to model relations such as *RightOf* or *LeftOf*, we need to account for the robot's viewpoint. Thus, we apply the halfspace-based model to the Contextualised Bounding Box we have defined earlier (see Figure 2). By definition, CBB is aligned with the contextualised frame of reference, $F_c$, so the front halfspace of CBB, for instance, can be defined w.r.t. a given $F_c$ as follows:

$$IsCBB(cbb_1, o_1) \Rightarrow hs(cbb_1, X_c^-, F_c) =$$
$$\{gp \in outReg(cbb_1) | gp = (x, y, z) \ w.r.t. F_c,$$
$$x'_{min} - x'_{min} \cdot s \leqslant x \leqslant x'_{min},$$
$$y'_{min} \leqslant y \leqslant y'_{max},$$
$$z'_{min} \leqslant z \leqslant z'_{max}\}$$
$$(30)$$

Capitalising on these spatial constructs, we can define the remaining directional relations:

$$RightOf(o_2, o_1, F_c) \Leftrightarrow Int(sr_2, hs(cbb_1, Y_c^-, F_c)) \quad (31)$$

$$LeftOf(o_2, o_1, F_c) \Leftrightarrow Int(sr_2, hs(cbb_1, Y_c^+, F_c)) \quad (32)$$

$$InFrontOf(o_2, o_1, F_c) \Leftrightarrow Int(sr_2, hs(cbb_1, X_c^-, F_c)) \quad (33)$$

$$Behind(o_2, o_1, F_c) \Leftrightarrow Int(sr_2, hs(cbb_1, X_c^+, F_c)) \quad (34)$$

For brevity, in axioms 31-34 we have omitted the predicate $IsCBB(cbb_1, o_1)$, which is always valid. The full definition is given in the supplementary materials. Next, we need to specify how the qualitative spatial relations we have identified align with commonsense spatial predicates.

**Commonsense spatial relations** In English, objects are represented by nouns while the spatial relationships between objects are mainly represented through prepositions - e.g.,
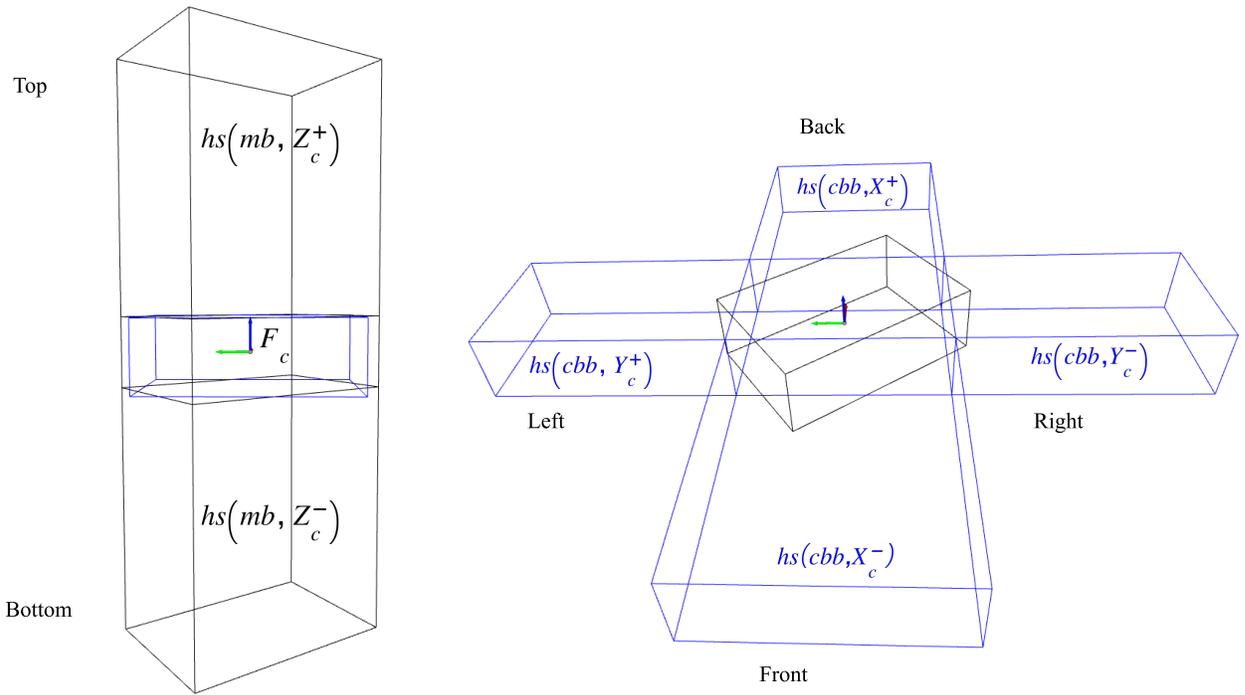
Figure 2: Halfspaces are generated by extruding the 3D bounding boxes along the axis direction in the frame of reference. Specifically, the top and bottom halfspaces, i.e., extruded along the $Z$ axis, are derived from the minimum oriented bounding box (left-hand side of the Figure). The left, right, front and bottom halfspaces are instead extruded from the Contextualised Bounding Box (right-hand side of the Figure).

on, next to, behind (Landau and Jackendoff 1993). Spatial relations are also implied by using certain verbs (e.g., person wears shirt). However, almost invariably, these verbs can be reduced to a simplified form, followed by a preposition (e.g., person has shirt on). Hence, the canonical structure of a spatial sentence consists of three elements: (i) a *reference* object and (ii) a *figure* object, both expressed as noun phrases, as well as (iii) a spatial preposition. The reference object and the preposition, together, define the spatial region occupied by the figure object.

As pointed out in (Landau and Jackendoff 1993), a set of reference axes is needed to differentiate the *front*, *back*, *top*, *bottom* and *sides* of an object. Specifically, the object's *top* and *bottom* are defined as "the regions at the ends of whichever axis is vertical in the object's normal orientation" (Landau and Jackendoff 1993). Thus, they are conceptually equivalent to the notion of top and bottom halfspaces we defined for the minimum oriented bounding box. Moreover, the object *front* is defined as the region at the end of the object's horizontal axis which also faces the observer. Conversely, the object *back* region is located opposite to the observer along the same axis. Finally, the region at the end of any other horizontal axis can be called a *side*. Thus, the four halfspaces we have defined based on the CBB cover these concepts.

Consequently, the directional relations defined at statements 28-34 are fit to model the directional predicates in (Landau and Jackendoff 1993). Moreover, the *LeftOf* and *RightOf*

relations can be combined so that, given $F_c$:

$$Beside(o_2, o_1, F_c) \Leftrightarrow RightOf(o_2, o_1, F_c) \vee \\ LeftOf(o_2, o_1, F_c) \tag{35}$$

An interesting case is that of the "on" preposition. One of the senses of "on" is semantically related to "above". However, while "above" typically implies absence of contact between the two objects, "on" strongly favours a contact reading (Landau and Jackendoff 1993). Formally, we make this distinction by defining:

$$OnTopOf(o_2, o_1, F_c) \Leftrightarrow Above(o_2, o_1, F_c) \wedge \\ Touches(o_2, o_1) \tag{36}$$

Nonetheless, the "on" preposition can also be used to denote that the figure object is supported by the reference object. For instance, we say that a "clock is on the wall" although the two objects overlap horizontally. The phrase "clock on wall" also implies that the wall is adequately stable to support the clock. Indeed, if two objects differ in terms of size and mobility, we typically prefer to consider the larger and more stable object as reference (Landau and Jackendoff 1993). To disambiguate these additional uses of "on", we define, for a given $F_c$:

$$LeansOn(o_2, o_1, F_c) \Leftrightarrow Touches(o_2, o_1) \wedge \\ \neg Above(o_2, o_1, F_c) \wedge \neg Below(o_2, o_1, F_c) \wedge \\ \exists o_3 [Touches(o_2, o_3) \wedge Below(o_3, o_2, F_c)] \tag{37}$$

$$Touches(o_2, o_1) \wedge \neg Above(o_2, o_1, F_c) \wedge$$
$$\neg \exists o_3 Touches(o_3, o_2) \Rightarrow AffixedOn(o_2, o_1, F_c) \quad (38)$$

Namely, whenever $o_2$ is supported by a reference object $o_1$ along the horizontal direction, it is typically said to be "leaning against" $o_1$: e.g., a ladder leaning against a wall. Furthermore, if the reference object $o_1$ provides the only support surface for $o_2$, $o_2$ is typically said to be *AffixedOn* $o_1$: e.g., a ladder which is affixed on the wall, above ground. Nonetheless, there may be cases where an object, $o_2$, is physically affixed to a surface, $o_1$, even though $o_1$ is not the only surface in contact with $o_2$: e.g., a ladder affixed at ground level. Hence, we used a single logic implication in Statement 37.

Similar considerations apply for the spatial preposition "in", which is also polysemous. First, "in" is generally used to imply that one object is "inside" another, or, based on our prior topological definitions, that one object is completely contained in the other (axiom 19). However, "in" is also used in cases where two objects only partially compenetrate each other. For instance, we would say that "a cat is in the box" even when the cat's tail is peeping from the box. To define this notion of partial containment we need first to define a function, *adjSRCard*, which, given two spatial regions, sr and sr', returns the cardinality of the set of points in sr' that are adjacent to points in sr:

$$adjSRCard(sr, sr') = |\{gp'|gp' \in outReg(sr) \wedge gp' \in sr' \wedge$$
$$\exists gp[gp \in sr \wedge Adj(gp, gp')]\}| \quad (39)$$

Hence, we can now define partial containment as follows:

$$PartIn(o_1, o_2) \Leftrightarrow sr = inter(sReg(o_1), sReg(o_2)) \wedge$$
$$adjSRCard(sr, sReg(o_1)) < adjSRCard(sr, sReg(o_2))] \quad (40)$$

Namely, $o_1$ is partially contained in $o_2$ iff the number of points in $o_1$ that are adjacent to the intersection region of $o_1$ and $o_2$ is strictly smaller than the number of points in $o_2$ that are adjacent to the same intersection region.

Lastly, an object is said to be "near" another object if it is located in a region "extending up to some critical distance" (Landau and Jackendoff 1993). This notion corresponds exactly to our definition of predicate *IsClose* (axiom 15).

## 4 Framework Applicability

In this Section, we assess the applicability of the proposed logic framework in concrete robotic scenarios. First, we illustrate how the raw sensor data can be opportunely processed to populate a spatial database (Section 4.1). Then, in Section 4.2, we evaluate the extent to which state-of-the-art GIS operators can cover the proposed set of Qualitative Spatial Relations. Furthermore, we show that, once a set of basic spatial concepts has been derived through GIS operators, our framework provides a method to combine these basic spatial concepts to model the commonsense spatial predicates of (Landau and Jackendoff 1993).

### 4.1 Populating the Spatial Database

Figure 3a shows an example of RGB-Depth (RGB-D) data collected through HanS' Orbbec Astra Pro monocular camera. At each time frame, $t$, the distance between the robot's

pose and the surfaces reached by the laser in the depth sensor is measured. These data are also known as *depth images*, and can be converted to collections of 3D geometrical points in the considered frame of reference, i.e., *PointClouds*. As in (Chiatti et al. 2021), RGB images are autonomously classified through Machine Learning (ML), i.e., based on the multi-branch Network of (Zeng et al. 2018). We then project the object regions annotated on RGB images on the PointCloud representing the observed scene, to obtain a segmented 3D region for each annotated object. In addition to the annotated object regions, we extract the planar surfaces representing the wall and floor areas. Specifically, we use the PCL library (Rusu and Cousins 2011) to reproduce the RANSAC planar segmentation algorithm. Then, we differentiate walls from floors based on the orientation of the plane normal. Because the robot pose at each $t$ is known, the depth values within each 3D object region can be converted to coordinate triples in the global frame of reference, $F_g$. In sum, we have produced a semantic map of the robot's environment.

Consistently with (Deeken, Wiemann, and Hertzberg 2018), we store the object regions and labels in the semantic map within a spatial database, implemented in PostgreSQL[2]. By linking these data to a spatial database, we can capitalise on the PostGIS engine[3], which provides a series of query operators for spatial reasoning over PostgreSQL databases. In particular, we rely on the SFCGAL backend[4], which extends PostGIS by supporting more advanced 3D operations. Objects are stored in the PostGIS database using a minimum oriented polyhedron derived by applying the convex hull algorithm on the segmented PointCloud. Planar surfaces are instead stored as 2D polygons. To populate the spatial database, for each 3D solid or 2D polygon, a new database record is added, which includes: (i) a unique identifier, obtained by concatenating the data collection timestamp with an incremental digit; (ii) the robot's heading and $x, y, z$ coordinates w.r.t. $F_g$; (iii) the top-5 object labels and related confidence scores, as predicted through ML; as well as (iv) a set of Bounding Box representations of the 3D solid, as further detailed in the next Section.

### 4.2 Coverage Study

The mapping of spatial concepts in our framework to GIS operators is summarised in Table 1a. In the following, we also explain the operational steps applied to obtain the output of each row in Table 1a.

**Minimum Oriented Bounding Box** To compute the minimum oriented bounding box, we input the 2D projection of the solid on the XY plane to the *ST_OrientedEnvelope* operator. Second, the *ST_ZMin* and *ST_ZMax* functions can be used to find the minimum and maximum coordinate of the 3D solid with respect to the vertical axis. The absolute difference between these two coordinates yields the height of the target bounding box, $h$. Then, the *ST_Extrude* operator can be used to extrude the 2D envelope along $Z$ by $h$.

---

[2]https://www.postgresql.org/
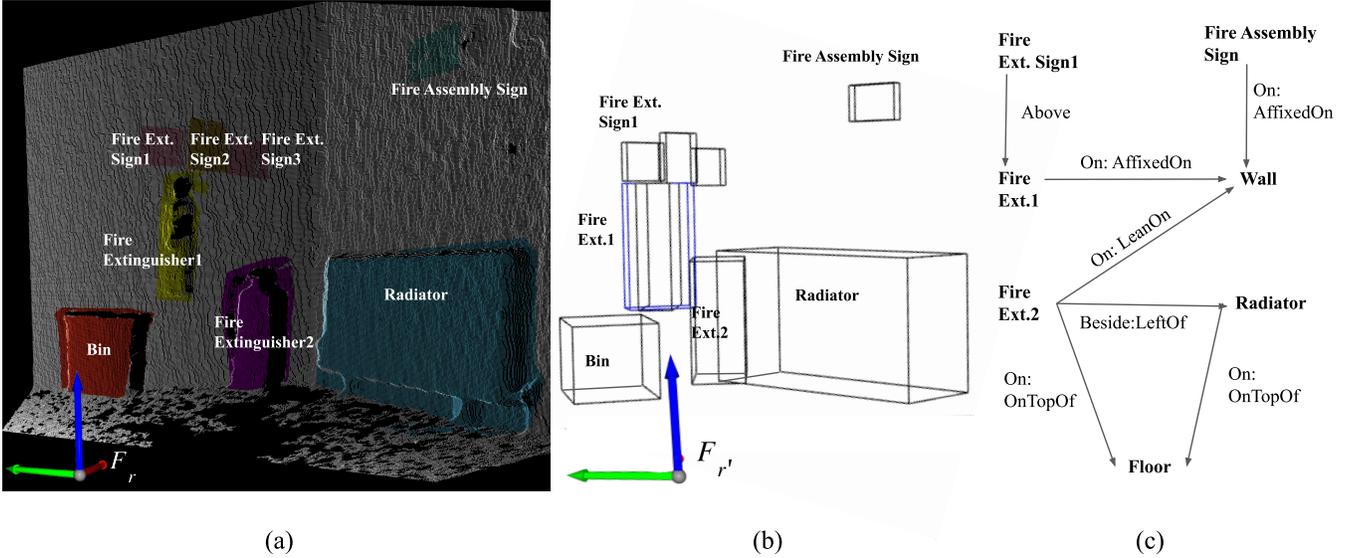[3]https://postgis.net/
[4]http://www.sfcgal.org/

Figure 3: Example of operational workflow: (a) the PointCloud representing the observed scene is first segmented and annotated with object categories. Then, (b) the minimum oriented bounding boxes and CBBs (in blue) are constructed. Lastly, (c) a set of QSR in figure-reference form is derived. In Figures 3b, 3c we show a subset of the bounding boxes and QSR representing the scene, for readability.

**Contextualised Bounding Box** The Contextualised Bounding Box, CBB, is obtained by rotating the minimum oriented bounding box by an angle $\theta$, to align it with $F_c$ (axiom 12). This operation can be achieved through the *ST_Rotate* operator. *ST_Rotate* requires, as input, a geometry, a rotation origin and an angle. We use *ST_Angle* to compute the angle between the heading of the robot and the bounding box, i.e., $\theta$. To derive the rotation origin, we can exploit *ST_Centroid*. Since *ST_Centroid* is a 2D operator, it is applied to the minimum oriented 2D rectangle returned by *ST_OrientedEnvelope*. Lastly, the transformation is applied to the line passing through the rotation origin and parallel to the $Z$ axis.

**Object Halfspaces** To derive the six object halfspaces, we can rely on the *ST_Extrude* operator. As emphasised in Table 1a, we generate the top and bottom halfspaces by applying *ST_Extrude* to the minimum oriented bounding box. However, to account for the orientation of both the object and the robot, the remaining four halfspaces are generated by extruding the Contextualised Bounding Box.

**Identifying the reference objects** Commonsense spatial relations are expressed with respect to a reference object and are asymmetric (Landau and Jackendoff 1993). For instance, we would say that *a plant is on the floor* (figure-reference form), but we would never say that *the floor is under a plant* (reference-figure form). Therefore, to extract only the QSR which are in figure-reference form, we need first to identify the set of reference objects in each scene. Reference objects are usually the largest and most stable among the observed objects (Landau and Jackendoff 1993). Hence, by definition, we consider as reference the objects of large area and negligible thickness, which we have modelled as 2D planes, i.e., walls and floors. Moreover, the *ST_Volume* operator pro-vides a way to measure the size of the 3D bounding boxes. To compute only QSR which are in figure-reference form, we sort objects by volume in descending order. Then, we only compute the QSR between one object and the nearby objects which are smaller than it, if any is found. As a result, in the example depicted in Figure 3, QSR such as *wall behind fire extinguisher1* or *floor under fire extinguisher2* would not be extracted. Thus, this design choice also reduces the computational load of extracting QSR for all pair-wise combinations of objects.

**Metric relations** To identify the set of objects which lie nearby a reference object, we can use the *ST_3DDWithin* operator. This operator returns true if the minimum 3D Euclidean distance between two objects is within a specified threshold. Thus, it is equivalent to our definition of object closeness (axiom 15). Hence, for all object pairs $o'$ and $o$ for which *ST_3DDWithin* returns true for a specified $T > 0$, the relation $IsClose(o, o')$ also holds. A special case is that of the $Touches(o, o')$ relation, where $T = 0$.

**Topological relations** The intersection relation we have defined at axiom 17 is directly covered by the *ST_3DIntersects* operator. Similarly, *ST_3DIntersection*, is equivalent to the aforementioned *inter* function (axiom 18). Neither PostGIS nor SFCGAL support 3D containment tests. To circumvent this limitation, we derive containment relations by comparing the volume of objects with the volume of their intersection region, through *ST_Volume*. Namely, if the volume of the intersection region equals the volume of the smaller object, e.g., $o'$, then $ComplCont(o, o')$.

**Directional relations** To derive directional QSR, the *ST_3DIntersects* operator can be applied to the object half-spaces constructed earlier. For instance, in Figure 3c, *fire*

| Input Geometry/ies | GIS operators applied | Output |
| --- | --- | --- |
| Convex Hull | *ST_OrientedEnvelope, ST_ZMin, ST_ZMax, ST_Extrude* | Min Oriented BBox |
| Min Oriented BBox, Robot heading | *ST_Rotate, ST_Angle, ST_Centroid* | CBB |
| Min Oriented BBox, $s$ | *ST_Extrude* | Top/Bottom Halfspaces |
| CBB, $s$ | *ST_Extrude* | L/R/Front/Back Halfspaces |
| Min Oriented BBoxes | *ST_Volume* | Reference object set |
| Min Oriented BBoxes | *ST_3DDWithin* | *IsClose, Touches* |
| Min Oriented BBoxes | *ST_3DIntersects* | *Intersects (Int)* |
| Min Oriented BBoxes | *ST_3DIntersection* | *inter* |
| Min Oriented BBoxes | *ST_3DIntersection, ST_Volume* | *CompletelyContains (ComplCont)* |
| Min Oriented BBox Halfspaces | *ST_3DIntersects* | *Left/RightOf Above, Below InFrontOf, Behind* |
| Min Oriented BBoxes | *ST_Scale, ST_Volume ST_Intersection* | *adjSRCard* |

(a)

| QSR | Follows from |
| --- | --- |
| *Beside* | *RightOf, LeftOf* |
| *OnTopOf* | *Touches, Above* |
| *LeansOn* | *Touches, Above Below* |
| *AffixedOn* | *Touches, Above* |
| *Inside* | *ComplCont* |
| *PartIn* | *inter, adjSRCard* |
| *Near* | *isClose* |

(b)

Table 1: (a) Coverage of spatial notions through PostGIS operators. (b) The basic spatial relations covered by PostGIS are combined to derive more complex QSR.

*extinguisher2* is on the left of the *radiator*, because it intersects the left halfspace of the *radiator*. Differently from (Deeken, Wiemann, and Hertzberg 2018), however, the left halfspace was here defined on a Contextualised Bounding Box, so that the robot's viewpoint is also accounted for.

**Commonsense spatial relations** As shown in Table 1a, PostGIS ensures a full coverage of the basic building blocks of our spatial framework. Then, the commonsense relations defined in Section 3 can be seen as a combination of these building blocks (Table 1b). For instance, having mapped the *LeftOf* and *RightOf* relation to GIS operators (Table 1a) we can also conclude whether $o, o'$ are *Beside* one another (Table 1b). In the case of *PartIn*, we first apply *ST_Intersection* to obtain the intersection region of $o, o'$, i.e., $inter(o, o')$. Then, to approximate the frontier of points which are adjacent to the intersection region, we scale $inter(o, o')$ by a $D$, via *ST_Scale*. In other words, we obtain a region which is infinitesimally larger than the intersection region. Thus, we can use this region to test, again through *ST_Intersection*, which sets of points overlap $o$ and $o'$. Specifically, since we are dealing with geometric regions, the cardinality of each point set (axiom 39) is given by the region volume, i.e., via *ST_Volume*. In sum, the aforementioned operators, also listed in Table 1, cover the logic functions needed to evaluate $PartIn(o, o')$.

Crucially, the introduction of commonsense QSR allows us to disambiguate polysemous spatial prepositions, such as "in" or "on". For instance, in Figure 3, both fire extinguishers are generically "on the wall". However, *fire extinguisher 1*, is affixed on the wall, whereas *fire extinguisher 2*, is also supported by the floor and thus leans on the wall, i.e., it may or may not be affixed on the wall. Thus, compared to the "on" preposition, the introduced QSR more precisely express the intuitive spatial and physics relations at play.

## 5 Conclusion and Future Work

In this paper, we have identified a framework for commonsense spatial reasoning which satisfies the concrete requirements of robot sensemaking in real-world scenarios. Differently from prior approaches to qualitative spatial reasoning in robotics, this framework is robust to variations in the robot's viewpoint and object orientation, thus ensuring scalability to many application scenarios. As highlighted by our coverage study, the proposed framework can be fully implemented by capitalising on state-of-the-art GIS technologies. Moreover, the proposed framework contributes a cognitively-inspired conceptual layer on top of these basic spatial operators, to model commonsense spatial predicates. The resulting linguistic predicates facilitate Human-Robot Interaction, as well as the integration of background spatial knowledge from external resources. As such, the proposed framework contributes to the broader objective of developing Visually Intelligent Agents (VIA), which can reliably assist us with our daily tasks.

Our future efforts will be focused on evaluating the performance impacts of augmenting ML-based object recognition methods with the proposed commonsense spatial reasoner. In this context, we will also assess the effects of combining commonsense spatial reasoning with the other types of reasoners contributing to the Visual Intelligence of a robot (Chiatti, Motta, and Daga 2020): e.g., size-based (Chiatti et al. 2021), motion-based and others.

# A Supplementary Materials

## A.1 Halfspace Definitions

The *hs* function we introduced at axiom 20 of the paper can be similarly applied to the other semi-axes in a given $F_o$. Specifically, given an input bounding box, and frame of reference, $F_o$:

$$MinBoundBox(mb_1, o_1) \Rightarrow hs(mb_1, X_o^-, F_o) =$$
$$\{gp \in outReg(mb_1) | gp = (x, y, z) \text{ w.r.t } F_o,$$
$$x_{min} - x_{min} \cdot s \leqslant x \leqslant x_{min},$$
$$y_{min} \leqslant y \leqslant y_{max},$$
$$z_{min} \leqslant z \leqslant z_{max}\}$$
$$(41)$$

$$MinBoundBox(mb_1, o_1) \Rightarrow hs(mb_1, Y_o^+, F_o) =$$
$$\{gp \in outReg(mb_1) | gp = (x, y, z) \text{ w.r.t } F_o,$$
$$x_{min} \leqslant x \leqslant x_{max},$$
$$y_{max} \leqslant y \leqslant y_{max} + y_{max} \cdot s,$$
$$z_{min} \leqslant z \leqslant z_{max}\}$$
$$(42)$$

$$MinBoundBox(mb_1, o_1) \Rightarrow hs(mb_1, Y_o^-, F_o) =$$
$$\{gp \in outReg(mb_1) | gp = (x, y, z) \text{ w.r.t } F_o,$$
$$x_{min} \leqslant x \leqslant x_{max},$$
$$y_{min} - y_{min} \cdot s \leqslant y \leqslant y_{min},$$
$$z_{min} \leqslant z \leqslant z_{max}\}$$
$$(43)$$

$$MinBoundBox(mb_1, o_1) \Rightarrow hs(mb_1, Z_o^+, F_o) =$$
$$\{gp \in outReg(mb_1) | gp = (x, y, z) \text{ w.r.t } F_o,$$
$$x_{min} \leqslant x \leqslant x_{max},$$
$$y_{min} \leqslant y \leqslant y_{max},$$
$$z_{max} \leqslant z \leqslant z_{max} + z_{max} \cdot s\}$$
$$(44)$$

$$MinBoundBox(mb_1, o_1) \Rightarrow hs(mb_1, Z_o^-, F_o) =$$
$$\{gp \in outReg(mb_1) | gp = (x, y, z) \text{ w.r.t } F_o,$$
$$x_{min} \leqslant x \leqslant x_{max},$$
$$y_{min} \leqslant y \leqslant y_{max},$$
$$z_{min} - z_{min} \cdot s \leqslant z \leqslant z_{min}\}$$
$$(45)$$

Similarly, to complete axiom 30 in the main paper, the back, left and right halfspaces are defined as follows, given $F_c$:

$$IsCBB(cbb_1, o_1) \Rightarrow hs(cbb_1, X_c^+, F_c) =$$
$$\{gp \in outReg(cbb_1) | gp = (x, y, z) \text{ w.r.t.} F_c,$$
$$x'_{max} \leqslant x \leqslant x'_{max} + x'_{max} \cdot s,$$
$$y'_{min} \leqslant y \leqslant y'_{max},$$
$$z'_{min} \leqslant z \leqslant z'_{max}\}$$
$$(46)$$

$$IsCBB(cbb_1, o_1) \Rightarrow hs(cbb_1, Y_c^+, F_c) =$$
$$\{gp \in outReg(cbb_1) | gp = (x, y, z) \text{ w.r.t.} F_c,$$
$$x'_{min} \leqslant x \leqslant x'_{max},$$
$$y'_{max} \leqslant y \leqslant y'_{max} + y'_{max} \cdot s,$$
$$z'_{min} \leqslant z \leqslant z'_{max}\}$$
$$(47)$$

$$IsCBB(cbb_1, o_1) \Rightarrow hs(cbb_1, Y_c^-, F_c) =$$
$$\{gp \in outReg(cbb_1) | gp = (x, y, z) \text{ w.r.t.} F_c,$$
$$x'_{min} \leqslant x \leqslant x'_{max},$$
$$y'_{min} - y'_{min} \cdot s \leqslant y \leqslant y'_{min},$$
$$z'_{min} \leqslant z \leqslant z'_{max}\}$$
$$(48)$$

## A.2 Directional Relations

The extended definitions of axioms 23-27 in the main paper for a given $F_o$ are:

$$W\_r(o_2, o_1, F_o) \Leftrightarrow MinBoundBox(mb_1, o_1) \wedge$$
$$Int(sReg(o_2), hs(mb_1, X_o^-, F_o)) \quad (49)$$

$$N\_r(o_2, o_1, F_o) \Leftrightarrow MinBoundBox(mb_1, o_1) \wedge$$
$$Int(sReg(o_2), hs(mb_1, Y_o^+, F_o)) \quad (50)$$

$$S\_r(o_2, o_1, F_o) \Leftrightarrow MinBoundBox(mb_1, o_1) \wedge$$
$$Int(sReg(o_2), hs(mb_1, Y_o^-, F_o)) \quad (51)$$

$$A\_r(o_2, o_1, F_o) \Leftrightarrow MinBoundBox(mb_1, o_1) \wedge$$
$$Int(sReg(o_2), hs(mb_1, Z_o^+, F_o)) \quad (52)$$

$$B\_r(o_2, o_1, F_o) \Leftrightarrow MinBoundBox(mb_1, o_1) \wedge$$
$$Int(sReg(o_2), hs(mb_1, Z_o^-, F_o)) \quad (53)$$

Moreover, the full definitions of axioms 31-34 in the main paper, given $F_c$, are:

$$RightOf(o_2, o_1, F_c) \Leftrightarrow IsCBB(cbb_1, o_1) \wedge$$
$$Int(sr_2, hs(cbb_1, Y_c^-, F_c)) \quad (54)$$

$$LeftOf(o_2, o_1, F_c) \Leftrightarrow IsCBB(cbb_1, o_1) \wedge$$
$$Int(sr_2, hs(cbb_1, Y_c^+, F_c)) \quad (55)$$

$$InFrontOf(o_2, o_1, F_c) \Leftrightarrow IsCBB(cbb_1, o_1) \wedge$$
$$Int(sr_2, hs(cbb_1, X_c^-, F_c)) \quad (56)$$

$$Behind(o_2, o_1, F_c) \Leftrightarrow IsCBB(cbb_1, o_1) \wedge$$
$$Int(sr_2, hs(cbb_1, X_c^+, F_c)) \quad (57)$$

# References

Alatise, M. B., and Hancke, G. P. 2020. A review on challenges of autonomous mobile robot and sensor fusion methods. *IEEE Access* 8:39830–39846.

Borrmann, A., and Rank, E. 2010. Query Support for BIMs using Semantic and Spatial Conditions. In *Handbook of Research on Building Information Modeling and Construction Informatics: Concepts and Technologies*. IGI Global. 405–450.

Chiatti, A.; Motta, E.; Daga, E.; and Bardaro, G. 2021. Fit to measure: Reasoning about sizes for robust object recognition. In *To appear in proceedings of the AAAI2021 Spring Symposium on Combining Machine Learning and Knowledge Engineering (AAAI-MAKE 2021)*.

Chiatti, A.; Motta, E.; and Daga, E. 2020. Towards a Framework for Visual Intelligence in Service Robotics: Epistemic Requirements and Gap Analysis. In *Proceedings of KR 2020- Special session on KR & Robotics*, 905–916. IJCAI.

Cohn, A. G., and Renz, J. 2008. Qualitative Spatial Representation and Reasoning. In van Harmelen, F.; Lifschitz, V.; and Porter, B., eds., *Foundations of Artificial Intelligence*, volume 3 of *Handbook of Knowledge Representation*. Elsevier. 551–596.

Coradeschi, S., and Saffiotti, A. 2003. An introduction to the anchoring problem. *Robotics and autonomous systems* 43(2-3):85–96.

Davis, E., and Marcus, G. 2015. Commonsense reasoning and commonsense knowledge in artificial intelligence. *Communications of the ACM* 58(9):92–103.

Deeken, H.; Wiemann, T.; and Hertzberg, J. 2018. Grounding semantic maps in spatial databases. *Robotics and Autonomous Systems* 105:146–165.

Hayes, P. J. 1988. *The Second Naive Physics Manifesto. Formal theories of the common sense world*. Ablex Publishing Corporation.

Herskovits, A. 1997. Language, Spatial Cognition, and Vision. In Stock, O., ed., *Spatial and Temporal Reasoning*. Dordrecht: Springer Netherlands. 155–202.

Kostavelis, I., and Gasteratos, A. 2015. Semantic mapping for mobile robotics tasks: A survey. *Robotics and Autonomous Systems* 66:86–103.

Krishna, R.; Zhu, Y.; Groth, O.; Johnson, J.; and et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision* 123(1):32–73.

Kunze, L.; Burbridge, C.; Alberti, M.; Thippur, A.; Folkesson, J.; Jensfelt, P.; and Hawes, N. 2014. Combining top-down spatial reasoning and bottom-up object class recognition for scene understanding. In *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2910–2915.

Lake, B. M.; Ullman, T. D.; Tenenbaum, J. B.; and Gershman, S. J. 2017. Building machines that learn and think like people. *Behavioral and Brain Sciences* 40.

Landau, B., and Jackendoff, R. 1993. "What" and "where" in spatial language and spatial cognition. *Behavioral and Brain Sciences* 16(2):217–238.

Levesque, H. 2017. *Common Sense, the Turing Test, and the Quest for Real AI*. The MIT Press.

Liu, H., and Wang, L. 2020. Remote human–robot collaboration: A cyber–physical system application for hazard manufacturing environment. *Journal of manufacturing systems* 54:24–34.

Moratz, R., and Ragni, M. 2008. Qualitative spatial reasoning about relative point position. *Journal of Visual Languages & Computing* 19(1):75–98.

Nilsson, P.; Haesaert, S.; Thakker, R.; Otsu, K.; Vasile, C.-I.; Agha-Mohammadi, A.-A.; Murray, R. M.; and Ames, A. D. 2018. Toward specification-guided active mars exploration for cooperative robot teams. *Robotics: Science and Systems (RSS)*.

Nüchter, A., and Hertzberg, J. 2008. Towards semantic maps for mobile robots. *Robotics and Autonomous Systems* 56(11):915–926.

Rusu, R. B., and Cousins, S. 2011. 3d is here: Point cloud library (pcl). In *2011 IEEE International Conference on Robotics and Automation (ICRA)*, 1–4. IEEE.

Sarthou, G.; Alami, R.; and Clodic, A. 2019. Semantic Spatial Representation: a unique representation of an environment based on an ontology for robotic applications. In *Combined Workshop on Spatial Language Understanding (SpLU) and Grounded Communication for Robotics (RoboNLP).*, 12.

Sisbot, E. A., and Connell, J. H. 2019. Where is My Stuff? An Interactive System for Spatial Relations. *arXiv:1909.06331 [cs]*. arXiv: 1909.06331.

Storks, S.; Gao, Q.; and Chai, J. Y. 2019. Recent advances in natural language inference: A survey of benchmarks, resources, and approaches. *arXiv preprint arXiv:1904.01172*.

Thippur, A.; Burbridge, C.; Kunze, L.; Alberti, M.; Folkesson, J.; Jensfelt, P.; and Hawes, N. 2015. A comparison of qualitative and metric spatial relation models for scene understanding. In *29th AAAI Conference and the 27th Innovative Applications of Artificial Intelligence Conference (IAAI)*, volume 2, 1632–1640. AI Access Foundation.

Thippur, A.; Stork, J. A.; and Jensfelt, P. 2017. Non-parametric spatial context structure learning for autonomous understanding of human environments. In *2017 26th IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*, 1317–1324. ISSN: 1944-9437.

Yang, G.; Lv, H.; Zhang, Z.; Yang, L.; Deng, J.; You, S.; Du, J.; and Yang, H. 2020. Keep healthcare workers safe: application of teleoperated robot in isolation ward for covid-19 prevention and control. *Chinese Journal of Mechanical Engineering* 33(1):1–4.

Yang, K.; Russakovsky, O.; and Deng, J. 2019. SpatialSense: An Adversarially Crowdsourced Benchmark for Spatial Relation Recognition. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2051–2060. Seoul, Korea (South): IEEE.

Young, J.; Kunze, L.; Basile, V.; Cabrio, E.; Hawes, N.; and Caputo, B. 2017. Semantic web-mining and deep vision for lifelong object discovery. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*, 2774–2779. IEEE.

Zeng, A.; Song, S.; Yu, K.-T.; Donlon, E.; Hogan, F. R.; Bauza, M.; Ma, D.; Taylor, O.; Liu, M.; Romo, E.; et al. 2018. Robotic pick-and-place of novel objects in clutter with multi-affordance grasping and cross-domain image matching. In *2018 IEEE ICRA Conference*, 1–8. IEEE.