

# SpaRC and SpaRP: Spatial Reasoning Characterization and Path Generation for Understanding Spatial Reasoning Capability of Large Language Models

Anonymous ACL submission

## Abstract

Spatial reasoning is a crucial component of both biological and artificial intelligence. In this work, we present a comprehensive study of the capability of current state-of-the-art large language models (LLMs) on spatial reasoning. To support our study, we created and contribute novel spatial characterization frameworks and datasets, the **Spatial Reasoning Characterization** (SpaRC), and **Spatial Reasoning Paths** (SpaRP), to enable an in-depth understanding of the spatial relations and compositions as well as the usefulness of spatial reasoning chains. We found that all state-of-the-art LLMs do not perform well on the datasets—their performances are consistently low across different setups. The spatial reasoning is an emergent capability as model sizes scale up. Finetuning both large language models (e.g., Llama-2-70B) and smaller ones (e.g., Llama-2-13B) can significantly improve their F1-scores by 7–32 absolute points. We also found that the top proprietary LLMs still significantly outperform their open-source counterparts in topological spatial understanding and reasoning.

## 1 Introduction

Spatial understanding and reasoning are a crucial component of both biological and artificial intelligence, essential for daily interactions and common tasks such as dialogues and conversations (Kruijff et al., 2007; Udagawa and Aizawa, 2019), navigation (Anderson et al., 2018; Chen et al., 2019; Zhang and Kordjamshidi, 2022), and robotics (Bisk et al., 2016; Venkatesh et al., 2021), among others. They require common reasoning steps such as identifying objects, determining other objects being involved, and aggregating multiple spatial relations to reach a conclusion. The advancement of the field has significantly benefited from many well-known tasks and datasets, including bAbI (Weston et al., 2016), SPARTQA (Mirzaee et al., 2021), SPARTUN and RESQ (Mirzaee and Kordjamshidi,

2022), and StepGame (Shi et al., 2022), among others.

Recently, Large Language Models (LLMs) have been shown to be capable of performing abstract, commonsense-based, and multi-hop reasoning (Wei et al., 2022b; Kojima et al., 2022; Wang et al., 2023). If such models are to be used as intelligent agents to answer questions, perform tasks, and collaborate with humans, whether they can understand the basic spatial relationships and perform corresponding reasoning would become critical to many real-life applications.

In this work, we present an extensive study on the state-of-the-art LLMs’ capability in spatial reasoning. The key components of spatial abilities include: (i) understanding spatial relations and composition, and (ii) developing reasoning chains to reach conclusions. Prior work (Mirzaee et al., 2021; Mirzaee and Kordjamshidi, 2022; Shi et al., 2022) has focused on the relations and spatial composition tied to a limited context setup, as will be detailed later in this paper. In our work, we propose a bottom-up approach that builds upon detailed spatial properties, providing fine control for constructing spatial rules and context setups. We formalize and propose **Spatial Reasoning Characterization** (SpaRC), a systematic framework in defining *spatial properties* of objects, relations, and contexts, as well as how they *characterize* spatial composition, which is inspired by the widely used benchmarks SPARTUN (Mirzaee and Kordjamshidi, 2022) and StepGame (Shi et al., 2022).

Reasoning paths are an integral part of the reasoning process and critical for analyzing and enhancing reasoning models. To the best of our knowledge, unlike other reasoning tasks such as mathematical reasoning, there exist no datasets with textual spatial reasoning paths. In this paper we develop deductively verified spatial reasoning paths by using spatial reasoners to generate step-by-step reasoning on SPARTUN and StepGame,

083 which is then verbalized to form textual chain-of- 133  
084 thoughts. We show that finetuning different sizes of 134  
085 LLMs (13B and 70B) on the reasoning paths signifi- 135  
086 cantly improve their spatial reasoning performance, 136  
087 which also highlights the poor performance of the 137  
088 generalist pretrained LLMs (without finetuning) on 138  
089 spatial reasoning. In summary, our contributions 139  
090 are as follows: 140

- 091 • We present a comprehensive study on the spa- 141  
092 tial reasoning capabilities of the state-of-the- 142  
093 art LLMs, under extensive setups: compre- 143  
094 hensive spatial characterizations, different parame- 144  
095 ter scales, pretrained vs. finetuned models, and 145  
096 different decoding strategies. We show that the 146  
097 current LLMs do not perform well on the spatial 147  
098 reasoning tasks. We observe that spatial reason- 148  
099 ing is an emergent capability as model sizes scale 149  
100 up. Top proprietary LLMs still significantly out- 150  
101 perform their open-source counterparts in topo- 151  
102 logical spatial reasoning. 152
- 103 • To support an in-depth study, we present 153  
104 the **Spatial Reasoning Characterization (SpaRC)** 154  
105 framework, a systematic bottom-up approach 155  
106 that shifts the focus towards spatial properties, 156  
107 providing a fine and flexible control on the spa- 157  
108 tial composition rules and context setups. We 158  
109 characterize and extend the widely used bench- 159  
110 mark datasets SPARTUN and StepGame under 160  
111 the SpaRC framework. 161
- 112 • We develop **Spatial Reasoning Paths (SpaRP)** by 162  
113 generating reasoning steps using symbolic spa- 163  
114 tial reasoners and verbalizing them in a deductive 164  
115 step-by-step process. We demonstrate that fine- 165  
116 tuning large language models on our reasoning 166  
117 paths can consistently improve their spatial rea- 167  
118 soning abilities. 168

## 119 2 Related Work 169

120 **Text-based Spatial Reasoning.** Textual spatial 170  
121 reasoning datasets present the task as question- 171  
122 answering (SRQA) over a textual spatial context. 172  
123 [Weston et al. \(2016\)](#) introduced bAbI containing 173  
124 two datasets focused on positional (Task 17) and 174  
125 navigational (Task 19) reasoning. Their simplistic 175  
126 nature and small size prompted subsequent works 176  
127 to create new and challenging datasets. [Mirzaee](#) 177  
128 [et al. \(2021\)](#) designed reasoning rules, and created 178  
129 human-generated and synthetic context-question- 179  
130 answer tuples from spatial description of visual 180  
131 scenes (SPARTQA) to train and evaluate spatial 181  
132 reasoning of neural language models. [Mirzaee](#)

and [Kordjamshidi \(2022\)](#) further extended the spa- 133  
tial rules to cover 16 spatial relations over mul- 134  
tiple formalism in 3D in their synthetic SPAR- 135  
TUN dataset, and commonsense spatial reason- 136  
ing in human-generated RESQ dataset. StepGame 137  
([Shi et al., 2022](#)) was introduced to assess robust 138  
positional multi-hop spatial reasoning in 2D. Our 139  
SpaRC framework builds on top of SPARTUN and 140  
StepGame as they provide a broad coverage over 141  
the number of hops and relations for abstract spatial 142  
reasoning. 143

## 144 Reasoning Abilities of Large Language Models. 144

145 Certain reasoning capabilities have been shown to 145  
146 be emergent abilities of LLMs ([Wei et al., 2022a](#)), 146  
147 which are further elicited by various chain-of- 147  
148 thought prompting techniques ([Wei et al., 2022b](#); 148  
149 [Kojima et al., 2022](#); [Yao et al., 2023](#); [Hao et al.,](#) 149  
150 [2023](#)). On logic-based tasks, including spatial rea- 150  
151 soning, they however lag significantly when com- 151  
152 pared to neuro-symbolic methods ([Mirzaee and](#) 152  
153 [Kordjamshidi, 2023](#); [Yang et al., 2023](#)). 153

154 To understand spatial reasoning abilities, [Bang](#) 154  
155 [et al. \(2023\)](#) provided a preliminary probing anal- 155  
156 ysis on ChatGPT using a very small dataset (60 156  
157 examples from each of StepGame and SPARTQA). 157  
158 [Yang et al. \(2023\)](#) evaluated the performance of 158  
159 GPT-3 on StepGame; [Mirzaee and Kordjamshidi](#) 159  
160 ([2023](#)) reported the performance of GPT-3 on 160  
161 SPARTQA, SPARTUN, and RESQ datasets. How- 161  
162 ever, these work are limited in terms of evalua- 162  
163 tion metric, qualitative analysis, past generation of 163  
164 LLMs, pretrained LLMs, or generation strategies. 164  
165 To the best of our knowledge, our work is the first 165  
166 attempt at a comprehensive evaluation of spatial 166  
167 reasoning of LLMs under these settings. 167

## 168 3 The Spatial Reasoning Characterization 168 169 (SpaRC) Framework 169

170 The steps to identify and compose spatial relations 170  
171 between entities distinguish spatial reasoning from 171  
172 other reasoning tasks. Prior work e.g. SPARTUN 172  
173 ([Mirzaee and Kordjamshidi, 2022](#)) and StepGame 173  
174 ([Shi et al., 2022](#)), have focused directly on the spa- 174  
175 tial composition rules coupled with the contexts, 175  
176 which can lead to different conclusions even for 176  
177 the same set of relations. For example, for the 177  
178 same context “A is left of B and B is above C”, 178  
179 applying the spatial composition of StepGame 179  
180 concludes that A is to the left and above C, while no 180  
181 directional relation between A and C can be con- 181  
182 cluded at all by applying the spatial rules of SPAR-

TUN. The conclusions are completely different *but* equally valid. This difference can be reconciled by examining the underlying spatial properties of the objects and relations, specifically the treatment of objects as points vs extended, and completeness of the knowledge of relations in the context. We, therefore, advocate for an extendable bottom-up approach starting from a more granular level and introduce the **Spatial Reasoning Characterization** (SpaRC) framework. SpaRC prioritizes spatial properties over spatial composition rules. Consequently, it offers finer control in creating contexts and facilitates a deeper and systematic examination of the spatial reasoning capabilities.

To keep our work closer and comparable to the widely used existing benchmarks, SPARTUN (Mirzaee and Kordjamshidi, 2022) and StepGame (Shi et al., 2022), we identify *six* properties that cover and characterize these datasets by two *distinct and mutually exclusive* sets of *three* properties each. With SpaRC, we further explore two *properties sets* (PS) with properties in common to these existing benchmarks.

$\mathcal{F}$	Sub-Type	Relations ( $\mathcal{R}$ )	Textual Label ( $\mathcal{L}$ )
Topological	$\mathcal{T}_R$ (RCC8)	DC	outside
		EC	outside and touching
		PO	partially overlapping
		EQ	overlapping
		TPP	inside and touching
Directional	$\mathcal{D}_R$ (Relative)	NTPP	inside
		TPPI	contains and touches
		NTPPI	contains
		LEFT	left
		RIGHT	right
	$\mathcal{D}_C$ (Cardinal)	ABOVE	above
		BELOW	below
		FRONT	front
		BEHIND	behind
		$\mathcal{D}_T$ (Clock)	12 o'clock
3 o'clock	right		
6 o'clock	below		
9 o'clock	left		
Distance	$\mathcal{S}_Q$ (Qualitative)	NEAR FAR	near far
	$\mathcal{S}_U$ (Quantitative)	–	–

Table 1: Formalisms ( $\mathcal{F}$ ) and their sub-types, relations ( $\mathcal{R}$ ) in the datasets and their labels ( $\mathcal{L}$ ). Labels are presented in natural language to work with language models. Composite relations e.g. lower-left are considered in a multi-label setting in the present work.

### 3.1 Principle and Design of SpaRC

We focus on a set of binary spatial relations  $\mathcal{R}$  (Table 1) by following the previous work (Mirzaee and Kordjamshidi, 2022; Shi et al., 2022). The relations cover three formalism ( $\mathcal{F}$ )—topological  $\mathcal{T}$ , directional  $\mathcal{D}$ , and distance  $\mathcal{S}$ , divided into sub-types—

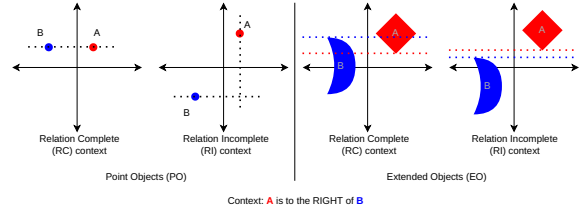


Figure 1: Visualization of Relation Complete (RC) and Relation Incomplete (RI) contexts for the RIGHT relation for Point Objects (PO) and Extended Objects (EO).

region connection calculus (RCC8)  $\mathcal{T}_R$ , relative directions  $\mathcal{D}_R$ , cardinal directions  $\mathcal{D}_C$ , clock-face directions  $\mathcal{D}_T$ , qualitative distance  $\mathcal{S}_Q$ , and quantitative distance  $\mathcal{S}_U$ .

For the relations set  $\mathcal{R}$  and a given set of entities  $\mathcal{E}$ , we denote a context  $\mathcal{C} = \{(h, r, t)_i\}_{i=1}^N$  as a set of  $(h, r, t)$  tuples, where  $h \in \mathcal{E}$  is a head entity,  $t \in \mathcal{E}$  is the tail entity, and  $r \in \mathcal{R}$  is the binary relation. Without loss of generality, objects are considered to be in a 2D space with  $(x_s, y_s)$  and  $(x_e, y_e)$  as the start and end positions. We now identify and describe six spatial properties of the objects, contexts, and relations that are crucial in determining their spatial composition rules. Refer to Appendix A for a more detailed discussion.

#### Fixed Orientation or Point of View (FPoV).

The directional relations are considered to be axis-aligned from a fixed orientation or point of view, i.e., fixed axes in a 2D or 3D space. A fixed mapping across the relative, cardinal, and clock-face directions is usually chosen. Consistent with the prior work, we map and canonicalize cardinal  $\mathcal{D}_C$  and clock-face  $\mathcal{D}_T$  relations to four relative directions  $\mathcal{D}_R$  (Table 1), *only* for their label representations  $\mathcal{L}$ . We denote the 2D subset of directions as  ${}^{2D}\mathcal{D} = \mathcal{D} \setminus \{\text{FRONT}, \text{BEHIND}\}$ .

#### Point Objects (PO).

A point object satisfies  $x_s = x_e \wedge y_s = y_e$ . As they are dimensionless, point objects have reduced set of relations with reference to other point objects. Real objects can be treated as point objects in practical contexts when their sizes are negligible.

#### Extended Objects (EO).

An object is said to be an extended object if  $x_s \neq x_e \vee y_s \neq y_e$ . In SpaRC, we extend StepGame by considering extended objects in addition to point objects. We further study additional composition rules for extended objects than presented in SPARTUN, as will be detailed later in Section 3.2.

Relation	Point Objects (PO)	Extended Objects (EO)
<i>Incomplete (RI):</i>		
RIGHT(A,B)	$x_A > x_B$	$x_s^A \geq x_e^B$
BELOW(A,B)	$y_A < y_B$	$y_e^A \leq y_s^B$
<i>Complete (RC):</i>		
RIGHT(A,B)	$x_A > x_B \wedge y_A = y_B$	$x_s^A \geq x_e^B \wedge y_s^B \leq y_e^A \wedge y_e^B \geq y_s^A$
BELOW(A,B)	$y_A < y_B \wedge x_A = x_B$	$y_e^A \leq y_s^B \wedge x_s^A \leq x_e^B \wedge x_e^A \geq x_s^B$
RIGHT(A,B) $\wedge$ BELOW(A,B)	$x_A > x_B \wedge y_A < y_B$	$x_s^A \geq x_e^B \wedge y_e^A \leq y_s^B$

Table 2: Mathematical descriptions of Relation Incomplete (RI) and Relation Complete (RC) contexts for the relations RIGHT, BELOW, and their combination in terms of entity positions  $(x, y)$  for Point Objects (PO) or entity boundaries  $(x_s, x_e, y_s, y_e)$  for Extended Objects (EO).

**Relation Incomplete (RI).** We introduce the term relation incomplete (RI) for a context  $\mathcal{C}$  between a head  $h$  and a tail  $t$  entity if *not all* the relations  $r \in \mathcal{R}$  between these entities are *considered* to be known and expressed in the context. Thus, the knowledge for the expressed relations should be *treated* as incomplete or partial for spatial composition. For example, “Ron is to the right of Hermione” as a RI context means that the direction orthogonal to the RIGHT could be ABOVE or BELOW as well. The state of positions or boundaries of objects on the *orthogonal* axes cannot be assumed. Table 2 and Figure 1 exemplifies and visualizes this for a few scenarios.

**Relation Complete (RC).** We introduce the term relation complete (RC) for a context  $\mathcal{C}$  between  $h$  and  $t$  if *all* the relations  $r \in \mathcal{R}$  between these entities are *considered* to be known and expressed in the context, and treated as such for spatial compositions. For the previous example “Ron is to the right of Hermione” to be considered as RC, the context should mean that Ron is only to the RIGHT of Hermione, and not to her lower-right or upper-right side. The positions or boundaries of objects on the *orthogonal* direction axes should coincide or overlap. Table 2 and Figure 1 exemplifies and visualizes this for a few scenarios. In SpaRC, we further consider this property in conjunction with other properties, such as extended objects, to design composition rules that are not present in StepGame, as discussed later in Section 3.2.

**Quantitatively Specified (QS).** A relation which is stated in terms of a unit of measurement is said to be quantitatively specified in the given context. Quantitatively specified relations that are *reverse*

of each other, e.g. {LEFT, RIGHT}, can readily be composed.

**Quantitatively Unspecified (QU).** A relation which can be stated in terms of a unit of measurement but is not stated as such in a given context is said to be quantitatively unspecified. Quantitatively unspecified relations that are *reverse* of each other, e.g. {LEFT, RIGHT}, cannot be composed unless they are quantified. In SpaRC, we design and study the reasoning abilities for this property in conjunction with other properties, such as point objects, that are not present in SPARTUN and StepGame, as discussed later in Section 3.2.

We restrict our study to the above 6 properties to keep it closer and comparable to the existing benchmarks, SPARTUN and StepGame. These properties form 3 *mutually exclusive* pairs—{EO, PO}, {RI, RC}, {QS, QU}, leading to 8 possible sets. SpaRC can be extended with additional properties, however, we note that the number of possible characterizations increases exponentially with the number of properties.

### 3.2 Creation of The SpaRC Dataset

We *identify* the properties set PS for the existing benchmarks, as formalized in the previous section, based on the generation process of the context and the spatial composition rules. More concretely, we identify that SPARTUN is characterized by the properties set PS1 = {EO, RI, QU}, while StepGame is characterized by the properties set PS2 = {PO, RC, QS}. These properties sets are mutually exclusive with PS2 supporting *stronger* composition rules than PS1 for a given context, e.g. “A is left of B and B is above C” as discussed earlier. Refer to Appendix B for more details.

In the SpaRC framework, we construct two additional datasets by relaxing the properties of StepGame from PO to EO, and QS to QU. We chose to extend StepGame as it is simple with fewer relations (only directional which is common across datasets and benchmarks) and challenging (more number of hops). Concretely, we create the datasets SpaRC-PS3 with the properties PS3 = {PO, RC, QU}, and SpaRC-PS4 with the properties PS4 = {EO, RC, QU}. Their composition rules, elaborated upon in Section 4, are formalized by the Algorithm 1 and Algorithm 2 respectively.

We confine our study to these four property sets as they encompass the two existing benchmarks and have covered the most interesting spatial rules

Not	$\forall(X, Y) \in Entities$	$R \in \{Dir \vee PP\}$	IF $R(X, Y)$	$\implies$ NOT( $R_{reverse}(X, Y)$ )
Inverse	$\forall(X, Y) \in Entities$	$R \in \{Dir \vee PP\}$	IF $R(Y, X)$	$\implies R_{reverse}(X, Y)$
Symmetry	$\forall(X, Y) \in Entities$	$R \in \{Dis \vee (RCC - PP)\}$	IF $R(Y, X)$	$\implies R(X, Y)$
Transitivity	$\forall(X, Y, Z) \in Entities$	$R \in \{Dir \vee PP\}$	IF $R(X, Z), R(Z, Y)$	$\implies R(X, Y)$
Combination	$\forall(X, Y, Z, H) \in Entities$	$R \in Dir, *PP \in PP$	IF $*PP(X, Z), R(Z, H), *PPi(Z, Y)$	$\implies R(X, Y)$

Table 3: Spatial Rules reproduced from SPARTUN (Mirzaee and Kordjamshidi, 2022). *Dir*: Directional relations (e.g., LEFT), *Dis*: Distance relations (e.g., FAR), *PP*: all Proper parts relations (NTPP, NTPPI, TPPI, TPP), *RCC - PP*: All RCC8 relation except proper parts relations. *\*PP*: one of TPP or NTPP. *\*PPi*: one of NTPPi or TPPi.

Dataset	$\mathcal{F}$	Properties	Textual Reason.	Split	# Context	# Ques.
SPARTUN	$\mathcal{T}_R, \mathcal{D}, \mathcal{S}_Q$	EO, RI, QU	✗	Train	6039	18400
				Dev	915	2818
				Test	925	2830
SpaRP-PS1			✓	Train	5806	16348
				Dev	877	2392
				Test	872	2301
StepGame	${}^{2D}\mathcal{D}, \mathcal{S}_U$	PO, RC, QS	✗	Train	50000	50000
				Dev	5000	5000
				Test	100000	100000
SpaRP-PS2			✓	Train	49243	49243
				Dev	4927	4927
				Test	98614	98614
SpaRP-PS3	${}^{2D}\mathcal{D}$	PO, RC, QU	✓	Train	44666	44666
				Dev	4494	4494
				Test	78092	78092
SpaRP-PS4		EO, RC, QU	✓	Train	41436	41436
				Dev	4171	4171
				Test	69474	69474

Table 4: Comparison between the extended (SpaRP) dataset and the source datasets. Descriptions of the properties are provided in Section 3.1. Relations contained in the formalisms are presented in Table 1. All the questions are of Find Relations (FR) types.

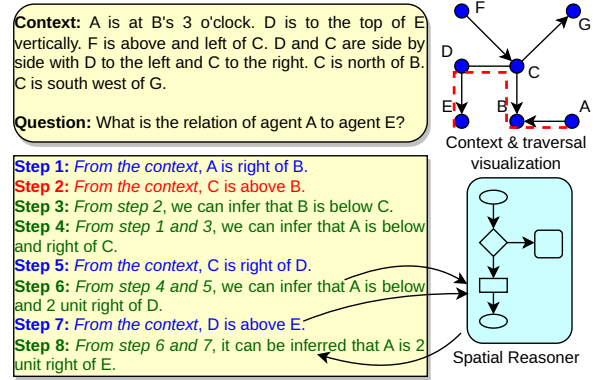


Figure 2: Our step-by-step deductive Spatial Reasoning Paths (SpaRP) generation. A context graph and node traversal from the head to the tail entity in a question is identified and verbalized. Blue indicates context relations  $r^c$ , red indicates inverse context relations  $r^{ic}$ , and green indicates deduced relations  $r^d$  between entities while traversing the reasoning path A-B-C-D-E.

and composition; the additional extension is more routine and we leave that as future work.

#### 4 The Spatial Reasoning Paths (SpaRP)

Reasoning paths are an integral part of reasoning models and critical for analyzing and enhancing such models. To the best of our knowledge, unlike other reasoning tasks such as mathematical reasoning, there exist no datasets with spatial reasoning paths. In this section, we develop deductively verified spatial reasoning paths by verbalizing the symbolic steps.

Existing spatial reasoning datasets can be considered as a collection of context-question-answer ( $\mathcal{C}, \mathcal{Q}, \mathcal{A}$ ) tuples. Formally, we denote a context  $\mathcal{C} = \{(h, r, t)\}_{i=1}^N$  defined over a set of entities  $\mathcal{E}$  and binary relations  $\mathcal{R}$  as a set of  $(h, r, t)$  tuples, where  $h \in \mathcal{E}$  is the head entity,  $t \in \mathcal{E}$  is the tail entity and  $r \in \mathcal{R}$  is the binary relation. For a given  $(\mathcal{C}, \mathcal{Q}, \mathcal{A})$  tuple, seeking relation between the head  $h_q$  and tail  $t_q$  entities, we define a symbolic reasoning path  $\mathcal{P} = (l_i)_{i=1}^L$  as a sequence of  $L$  reasoning links  $l_i = (h_i, r_i^{\cup}, t_i)$  such that  $h_1 = h_q$ ,  $t_L = t_q$ , and  $h_i = t_{i-1}$  for  $1 < i \leq L$ . We de-

fine  $r^{\cup} = r^c \cup r^{ic} \cup r^d$ , where  $r^c$  denotes the set of relations present in the context,  $r^{ic}$  denotes the inverse relations present in the context i.e. relations from  $t$  to  $h$ , and  $r^d$  denotes the set of deduced relations. Following the format of deductively verified chain-of-thought (Ling et al., 2023), we verbalize the reasoning path  $\mathcal{P}$  as a series of step-by-step reasoning sentences, where each step receives their necessary context and premises (Figure 2). The overall process is as given below:

- Entities and their relations in the contexts are either pre-annotated (SPARTUN) or extracted using regex pattern matching (StepGame) to construct the symbolic context  $\mathcal{C}$ .
- A traversal path  $\mathcal{P}$  is identified from  $h_q$  to  $t_q$  by constructing a network graph over  $\mathcal{C}$ . The deduced relations  $r^d$  are initialized to be the inverse of  $r^{ic}$ , to traverse and merge steps in a single direction from  $h_q$  to  $t_q$  (Figure 2).
- We traverse the path  $\mathcal{P}$ , progressively merging the links (as  $h_i = t_{i-1}$ ) and updating the deduced relations  $r^d$  based on the properties set PS and their spatial composition rules:
  - For SPARTUN we reuse the rules from

**Algorithm 1** Relative Direction composition for set of properties PS2 and PS3 in 2D.

**Input:** Pairs to compose  $\{pair1, pair2\}$ .  
 $quantitative \in \{true, false\}$ .  
**Output:** merged pair.

```

1: /* initialized pair starts with  $dx = dy = 0$  */
2: merged  $\leftarrow$  InitializePair
3: merged.head  $\leftarrow$  pair1.head
4: merged.tail  $\leftarrow$  pair2.tail
5: for pair  $\in \{pair1, pair2\}$  do
6:   for delta  $\in \{dx, dy\}$  do
7:     delta  $\leftarrow$  merged.delta + pair.delta
8:     /* Handle direction reversal and quantitatively
9:     unspecified */
10:    if (merged.delta  $\times$  pair.delta < 0) and not
11:    quantitative then
12:      /* Set as NaN to invalidate compositions
13:      from now on in this direction */
14:      merged.delta  $\leftarrow$  NaN
15:    else
16:      merged.delta  $\leftarrow$  delta
17:    end if
18:  end for
19: end for

```

Mirzaee and Kordjamshidi (2022), reproduced in Table 3.

- For StepGame and SpaRC-PS3, we represent the relative positions as signed integers on the  $x$  and  $y$  axis, and numerically compose them (Algorithm 1). Without the quantitative knowledge of backtracking along a given axis, e.g.  $x$ -axis for {LEFT, RIGHT}, no subsequent inferences can be made for those directions.
- For SpaRC-PS4, the relations in context can be expressed as logical conjunction  $\wedge$  of inequalities, refer Section 3, Table 2, and Figure 1. For composition of relations to merge reasoning steps, consistency of inequalities for relations  $r \in \mathcal{D}$  is checked and the deduced relations set  $r^d$  is updated (Algorithm 2).

4. We finally *verbalize* the reasoning path  $\mathcal{P}$  link-by-link (Figure 2) following the format of deductively verified chain-of-thought (Ling et al., 2023). However, instead of generating and self-verifying LLM outputs, we use spatial reasoners for ground truth generation.

We denote the extended dataset as **Spatial Reasoning Paths** (SpaRP). Specifically, we extended SPARTUN, StepGame, SpaRC-PS3, and SpaRC-PS4, to be SpaRP-PS1, SpaRP-PS2, SpaRP-PS3 and SpaRP-PS4, respectively, by enriching the former with the reasoning paths. A comparison of the derived datasets with the original datasets is summarized in Table 4.

**Algorithm 2** Relative Direction composition for set of properties PS4 in 2D.

**Input:** Pairs to compose  $\{pair1, pair2\}$ .  
current set of constraint inequalities  $ineq$   
**Output:** merged pair and updated inequalities  $ineq$ .

```

1: /* initialize an empty pair */
2: merged  $\leftarrow$  InitializePair
3: merged.head  $\leftarrow$  pair1.head
4: merged.tail  $\leftarrow$  pair2.tail
5: for rel  $\in \{LEFT, RIGHT, ABOVE, BELOW\}$  do
6:   candidate_ineq  $\leftarrow$  substitute_entities(
7:   rel.ineq, merged.head, merged.tail)
8:   consistent  $\leftarrow$  check_consistency(
9:   candidate_ineq, ineq)
10:  if consistent then
11:    insert(candidate_ineq, ineq)
12:    insert(rel, merged.relations)
13:  end if
14: end for

```

## 5 Experimental Setup

**Dataset.** Due to the expense and resource limitations for running LLMs, for each of the four subsets of SpaRP, we randomly sample 2000, 500, and 1000 datapoints as our training, validation, and test set, respectively. We call them small SpaRP, or **SpaRP-S**. We also randomly sample equal number of instances for each *number of hops* in the reasoning path. Additionally, we collect *five* diverse sets of *human-generated* natural language descriptions of the properties relevant to spatial compositions, and construct a *system prompt* template with a unified task instruction using these descriptions.

**Implementation Details.** To help replicability, we include implementation details such as dataset sampling, system prompt, and training parameters in Appendix-C.

**Evaluation Metrics.** We use exact-match accuracy and macro-averaged F1-scores<sup>1</sup>.

## 6 Results and Analysis

We run experiments with three state-of-the-art LLMs — Llama-2-13B, Llama-2-70B (Touvron et al., 2023), and GPT-4, each with both single greedy decoding and self-consistency (Wang et al., 2023) with majority voting over 20 generations with sampling (SC=20). Inputs are provided with a “system prompt” containing task instructions and 5-shot CoT with randomly sampled exemplars from the *relevant dev-set*. We also finetune Llama-2 13B and 70B models, indicated by FT in Table 5, using

<sup>1</sup>We used the `scikit-learn` v1.3.2 library

Dataset	Model	Acc.	F1
SpaRP-S-PS1 (SPARTUN)	Llama-2-13B	0.2	0.49
	Llama-2-13B-FT	18.9	22.23
	Llama-2-70B	10.1	23.37
	Llama-2-70B-FT	28	36.49
	Llama-2-70B <sub>SC=20</sub>	17.1	27.95
	GPT-4	46.8	54.30
	GPT-4 <sub>SC=20</sub>	<b>54.3</b>	<b>60.32</b>
SOTA (PISTAQ)		94.52	–
SpaRP-S-PS2 (StepGame)	Llama-2-13B	0.1	0.47
	Llama-2-13B-FT	13.7	33.23
	Llama-2-70B	10.6	26.41
	Llama-2-70B-FT	16.6	34.63
	Llama-2-70B <sub>SC=20</sub>	20.30	38.96
	GPT-4	23.9	41.09
	GPT-4 <sub>SC=20</sub>	<b>28.6</b>	<b>43.01</b>
SOTA (LLM-ASP)		90.88	–
SpaRP-S-PS3	Llama-2-13B	0.2	0.92
	Llama-2-13B-FT	27.3	32.01
	Llama-2-70B	9.4	25.27
	Llama-2-70B-FT	19.5	32.97
	Llama-2-70B <sub>SC=20</sub>	15.2	32.01
	GPT-4	23.8	35.17
	GPT-4 <sub>SC=20</sub>	<b>32.5</b>	<b>42.06</b>
SpaRP-S-PS4	Llama-2-13B	0.7	1.84
	Llama-2-13B-FT	30.6	31.62
	Llama-2-70B	9.0	22.13
	Llama-2-70B-FT	20	31.74
	Llama-2-70B <sub>SC=20</sub>	18.3	29.73
	GPT-4	21.7	33.02
	GPT-4 <sub>SC=20</sub>	<b>32.9</b>	<b>40.23</b>

Table 5: Performance evaluations of Llama-2 (13B and 70B) and GPT-4 models on the individual datasets.

QLoRA (Detmers et al., 2023) on the *verbalized reasoning paths* made available by SpaRP.

**Overall Results.** As shown in Table 5, we can observe that the performance of all the state-of-the-art LLMs on the spatial reasoning datasets are low, lagging significantly behind the existing state-of-the-art symbolic-based models such as PISTAQ (Mirzaee and Kordjamshidi, 2023) and LLM-ASP (Yang et al., 2023) on SPARTUN and StepGame, respectively, suggesting that if these generalist models are to be used for any spatial-reasoning-related tasks (e.g., in LLMs-based agents), caution should be exerted.

Among these models, GPT-4 under SC=20 exhibits the best performance overall, followed closely by GPT-4 with greedy decoding. The latter outperforms even the largest open-source Llama-2-70B model with SC=20.

We also observed the emergent abilities of LLMs on spatial reasoning as model sizes scale up. The

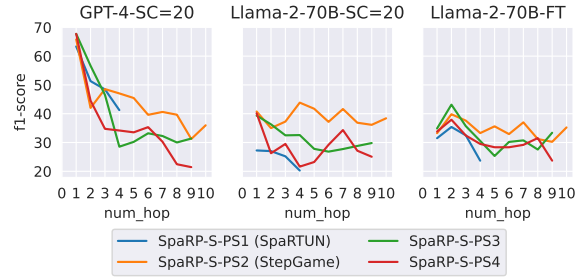


Figure 3: F1 scores vs. number of hops for spatial reasoning.

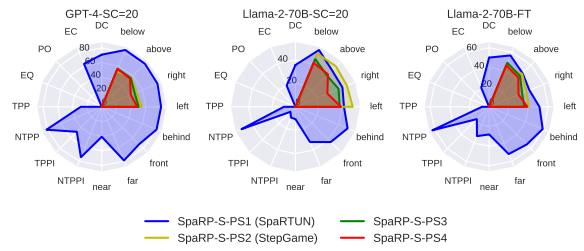


Figure 4: F1 scores of individual labels.

F1-score of Llama-2 13B model on SpaRP-S-PS1 (SPARTUN) is only 0.49 (no spatial reasoning ability), significantly lower than the 23.37 F1-score of the 70B model. Similarly, for SpaRP-S-PS3, the F1-scores are 0.92 and 25.27 for the 13B and 70B models, respectively. This is also observed on the other datasets.

**Impact of Spatial Properties and Composition Rules.** StepGame and SpaRP-PS3 consider entities as point objects (PO), however, SpaRP-PS3 does not quantify directions rendering them not composable while backtracking, e.g. RIGHT followed by LEFT is not composable. SpaRP-PS4 considers entities as real objects with extended sizes, thereby introducing added complexity to spatial relation composition (Section 4 and Algorithm 2). The F1-scores (Table 5) of both GPT-4 and Llama-2 underscore these challenges.

Furthermore, Figure 3 demonstrates that the F1-scores for both SpaRP-S-PS3 and SpaRP-S-PS4 consistently trail those of SpaRP-S-PS2 (StepGame) across varying numbers of hops. This highlights the utility of our SpaRC framework in identifying additional challenges that are not addressed by the existing benchmarks.

**Relation-wise Performance.** The performance of GPT-4 is significantly better compared to Llama-2 models on SpaRP-S-PS1 (SPARTUN), which

Errors, examples and explanations
<p><b>Error:</b> Incorrect relation extraction  <b>Context:</b> Box EEE has a tiny white rectangle and covers a midsize white diamond.  <b>Extracted:</b> Step 7: It is given that the tiny white rectangle is inside and touching the box EEE.  <b>Explanation:</b> Has only means inside.</p>
<p><b>Error:</b> Reverse answer  <b>Question:</b> What is the relation of the agent W to the agent X?  <b>Answer Step:</b> Step 8: From step 5 and 7, we can infer that X is above and left of W. Hence, the answer is above, and left.  <b>Explanation:</b> Directional relations are non-symmetric. Question is from W to X, while answer is from X to W.</p>
<p><b>Error:</b> Copied, not composed  <b>Reasoning Steps:</b> Step 6: From step 4 and 5, we can infer that A is below and right of S. Step 7: From the context, S is left of M. Step 8: From step 6 and 7, we can infer that A is below and right of M.  <b>Explanation:</b> Relation from S to M not used in composition. Instead, relation from A to S is copied in step 8.</p>
<p><b>Error:</b> Composed without connection  <b>Reasoning Steps:</b> Step 5: From step 3 and 4, we can infer that Y is right of L. Step 14: From step 12 and 13, we can infer that C is below and right of K. Step 15: From step 5 and 14, we can infer that C is below and right of L.  <b>Explanation:</b> No common entity between merged steps 5 and 14 which are 9 steps apart.</p>

Table 6: Errors, their examples (only relevant steps) and explanations in the model generated reasoning paths.

has a larger candidate set comprising of 16 relations, including 8 topological relations. In contrast, SpaRP-S-PS2 (StepGame) has a smaller candidate set consisting of only directional relations. This highlights a notable deficiency in Llama-2 regarding the understanding and composition of topological relations. More importantly, even the finetuned Llama-2 model falls short of GPT-4’s performance. The top proprietary LLMs still significantly outperform their open-source counterparts in topological spatial reasoning.

Additionally, Figure 3 demonstrates that even when controlling for the same number of hops, the F1-scores of Llama-2 on SpaRP-S-PS1 (SPARTUN) rank lowest across all hops. An examination of F1-scores on a per-relation basis (Figure 4) further confirms this difficulty of topological relations for Llama-2 models compared to GPT-4.

**Finetuning with Reasoning Paths.** We observe that finetuning the 13B and 70B models with the reasoning paths made available in SpaRP consistently improves the spatial reasoning capabilities. Finetuning consistently boosts the F1-score by 21–32 and 7–13 points for 13B and 70B models respectively, across the datasets. In specific instances, the accuracy of a finetuned 13B model surpasses that of 5-10 times larger models such as Llama-2-70B

with SC=20, and GPT-4. We hope the proposed reasoning-path generation can be further used for improving LLMs’ explainability and robustness on spatial reasoning.

**Error Analysis of Reasoning Paths.** We sampled and manually analyzed a total of 80 model generated reasoning paths across all datasets for both the GPT-4 and Llama-2 70B models. The deductive step-by-step reasoning path made available by SpaRP proves to be useful in identifying errors in the generated outputs (Table 6). Commonly observed errors include incorrect parsing or retrieval of relations from the contexts, especially for topological relations. Additionally, we observe instances of reverse answering, where relations between tail to head entities are returned instead of head to tail entities in a question. More complex reasoning failures involve copying relations from one of the reasoning steps instead of composing them. Similarly, composing relations between reasoning steps without a common entity is observed frequently over distant steps. Additional errors with examples are provided in Appendix D. These errors are more prevalent in Llama-2 models, resulting in poorer performance compared to GPT-4.

## 7 Conclusion

Spatial reasoning is one of the basic components of intelligence. We perform a study on the spatial reasoning abilities of the latest LLMs under comprehensive setups. To support the study, we introduce (SpaRC), a systematic framework to characterize spatial reasoning scenarios by identifying and defining six spatial properties of objects, spatial relations, and contexts, and their impact on the spatial composition rules. Based on that, we create the (SpaRP) reasoning paths for the datasets. We found that the state-of-the-art LLMs do not perform well on the datasets — their performances are consistently low across different setups. The spatial reasoning is an emergent capability as model sizes scale up. Finetuning both large language models (e.g., Llama-2-70B) and smaller ones (e.g., Llama-2-13B) can significantly improve their performance by 7–32 points on F1-scores. We also found top proprietary LLMs still significantly outperform their open-source counterparts in topological spatial understanding and reasoning. We provide detailed analyses and insights in our experiments.



## 568 Limitations

569 We aimed to characterize various properties of the  
570 objects, relations, contexts and the associated spa-  
571 tial composition rules. We, however, note that  
572 the spatial scenarios, relations and interactions be-  
573 tween objects can still be incomplete. Further, the  
574 existing datasets and our extensions of them still  
575 pertain to a limited combination of the character-  
576 izations *in isolation* in a context. Even with our  
577 proposed characterizations, a combination of these  
578 within a single context is common in the real world,  
579 including multi-modality with visual perception,  
580 which we haven't considered in our current study.  
581 The base datasets, although textual, are synthetic in  
582 nature. Combined with the use of symbolic reason-  
583 ers for our reasoning path generation, our dataset  
584 inherit all the associated limitations such as relative  
585 lack of linguistic diversity, types of objects, rela-  
586 tions etc. Finally we note that due to the cost and  
587 resource constraints of using LLMs, we worked  
588 with a smaller set of about 1000 test instances per  
589 dataset, which is a common data size to work with  
590 LLMs.

## 591 References

592 P. Anderson, Q. Wu, D. Teney, J. Bruce, M. Johnson,  
593 N. Sunderhauf, I. Reid, S. Gould, and A. van den  
594 Hengel. 2018. [Vision-and-language navigation: In-  
595 terpreting visually-grounded navigation instructions  
596 in real environments](#). In *2018 IEEE/CVF Confer-  
597 ence on Computer Vision and Pattern Recognition  
598 (CVPR)*, pages 3674–3683, Los Alamitos, CA, USA.  
599 IEEE Computer Society.

600 Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wen-  
601 liang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei  
602 Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu,  
603 and Pascale Fung. 2023. [A multitask, multilingual,  
604 multimodal evaluation of ChatGPT on reasoning, hal-  
605 lucination, and interactivity](#). In *Proceedings of the  
606 13th International Joint Conference on Natural Lan-  
607 guage Processing and the 3rd Conference of the Asia-  
608 Pacific Chapter of the Association for Computational  
609 Linguistics (Volume 1: Long Papers)*, pages 675–718,  
610 Nusa Dua, Bali. Association for Computational Lin-  
611 guistics.

612 Yonatan Bisk, Deniz Yuret, and Daniel Marcu. 2016.  
613 [Natural language communication with robots](#). In  
614 *Proceedings of the 2016 Conference of the North  
615 American Chapter of the Association for Computa-  
616 tional Linguistics: Human Language Technologies*,  
617 pages 751–761, San Diego, California. Association  
618 for Computational Linguistics.

619 Howard Chen, Alane Suhr, Dipendra Misra, Noah  
620 Snaveley, and Yoav Artzi. 2019. [Touchdown: Natural](#)

[language navigation and spatial reasoning in visual  
street environments](#). In *2019 IEEE/CVF Conference  
on Computer Vision and Pattern Recognition (CVPR)*,  
pages 12530–12539. 621  
622  
623  
624

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and  
625 Luke Zettlemoyer. 2023. [QLoRA: Efficient Finetun-  
626 ing of Quantized LLMs](#). In *Thirty-seventh Confer-  
627 ence on Neural Information Processing Systems*. 628

Shibo Hao, Yi Gu, Haodi Ma, Joshua Hong, Zhen  
629 Wang, Daisy Wang, and Zhiting Hu. 2023. [Rea-  
630 soning with language model is planning with world  
631 model](#). In *Proceedings of the 2023 Conference on  
632 Empirical Methods in Natural Language Processing*,  
633 pages 8154–8173, Singapore. Association for Com-  
634 putational Linguistics. 635

Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yu-  
636 taka Matsuo, and Yusuke Iwasawa. 2022. [Large lan-  
637 guage models are zero-shot reasoners](#). In *Advances in  
638 Neural Information Processing Systems*, volume 35,  
639 pages 22199–22213. Curran Associates, Inc. 640

Geert-Jan M. Kruijff, Hendrik Zender, Patric Jensfelt,  
641 and Henrik I. Christensen. 2007. [Situated dialogue  
642 and spatial organization: What, where... and why?](#)  
643 *International Journal of Advanced Robotic Systems*,  
644 4(1):16. 645

Zhan Ling, Yunhao Fang, Xuanlin Li, Zhiao Huang,  
646 Mingu Lee, Roland Memisevic, and Hao Su. 2023.  
647 [Deductive verification of chain-of-thought reasoning](#).  
648 In *Thirty-seventh Conference on Neural Information  
649 Processing Systems*. 650

Roshanak Mirzaee and Parisa Kordjamshidi. 2022.  
651 [Transfer learning with synthetic corpora for spatial  
652 role labeling and reasoning](#). In *Proceedings of the  
653 2022 Conference on Empirical Methods in Natu-  
654 ral Language Processing*, pages 6148–6165, Abu  
655 Dhabi, United Arab Emirates. Association for Com-  
656 putational Linguistics. 657

Roshanak Mirzaee and Parisa Kordjamshidi. 2023. [Dis-  
658 entangling extraction and reasoning in multi-hop spa-  
659 tial reasoning](#). In *Findings of the Association for  
660 Computational Linguistics: EMNLP 2023*, pages  
661 3379–3397, Singapore. Association for Computa-  
662 tional Linguistics. 663

Roshanak Mirzaee, Hossein Rajaby Faghihi, Qiang  
664 Ning, and Parisa Kordjamshidi. 2021. [SPARTQA:  
665 A textual question answering benchmark for spatial  
666 reasoning](#). In *Proceedings of the 2021 Conference of  
667 the North American Chapter of the Association for  
668 Computational Linguistics: Human Language Tech-  
669 nologies*, pages 4582–4598, Online. Association for  
670 Computational Linguistics. 671

Zhengxiang Shi, Qiang Zhang, and Aldo Lipani. 2022.  
672 [Stepgame: A new benchmark for robust multi-hop  
673 spatial reasoning in texts](#). In *Proceedings of the AAAI  
674 Conference on Artificial Intelligence*, volume 36,  
675 pages 11321–11329. 676

677	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	Jason Weston, Antoine Bordes, Sumit Chopra, and	737
678	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	Tomás Mikolov. 2016. <a href="#">Towards ai-complete question</a>	738
679	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	<a href="#">answering: A set of prerequisite toy tasks</a> . In <i>4th In-</i>	739
680	Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton	<i>ternational Conference on Learning Representations,</i>	740
681	Ferrer, Moya Chen, Guillem Cucurull, David Esiobu,	<i>ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016,</i>	741
682	Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller,	<i>Conference Track Proceedings</i> .	742
683	Cynthia Gao, Vedanuj Goswami, Naman Goyal, An-		
684	thony Hartshorn, Saghar Hosseini, Rui Hou, Hakan	Zhun Yang, Adam Ishay, and Joohyung Lee. 2023. <a href="#">Cou-</a>	743
685	Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa,	<a href="#">pling large language models with logic programming</a>	744
686	Isabel Kloumann, Artem Korenev, Punit Singh Koura,	<a href="#">for robust and general reasoning from text</a> . In <i>Find-</i>	745
687	Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Di-	<i>ings of the Association for Computational Linguis-</i>	746
688	ana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-	<i>tics: ACL 2023</i> , pages 5186–5219, Toronto, Canada.	747
689	tinnet, Todor Mihaylov, Pushkar Mishra, Igor Moly-	Association for Computational Linguistics.	748
690	bog, Yixin Nie, Andrew Poulton, Jeremy Reizen-		
691	stein, Rashi Rungta, Kalyan Saladi, Alan Schelten,	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom	749
692	Ruan Silva, Eric Michael Smith, Ranjan Subrama-	Griffiths, Yuan Cao, and Karthik Narasimhan. 2023.	750
693	nian, Xiaoqing Ellen Tan, Binh Tang, Ross Tay-	<a href="#">Tree of thoughts: Deliberate problem solving with</a>	751
694	lor, Adina Williams, Jian Xiang Kuan, Puxin Xu,	<a href="#">large language models</a> . In <i>Thirty-seventh Conference</i>	752
695	Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan,	<i>on Neural Information Processing Systems</i> .	753
696	Melanie Kambadur, Sharan Narang, Aurelien Rod-		
697	riguez, Robert Stojnic, Sergey Edunov, and Thomas	Yue Zhang and Parisa Kordjamshidi. 2022. <a href="#">Explicit</a>	754
698	Scialom. 2023. <a href="#">Llama 2: Open foundation and fine-</a>	<a href="#">object relation alignment for vision and language</a>	755
699	<a href="#">tuned chat models</a> . <i>CoRR</i> , abs/2307.09288.	<a href="#">navigation</a> . In <i>Proceedings of the 60th Annual Meet-</i>	756
		<i>ing of the Association for Computational Linguistics:</i>	757
		<i>Student Research Workshop</i> , pages 322–331, Dublin,	758
		Ireland. Association for Computational Linguistics.	759
700	Takuma Udagawa and Akiko Aizawa. 2019. <a href="#">A nat-</a>		
701	<a href="#">ural language corpus of common grounding under</a>	<b>A Additional details and comparison of</b>	760
702	<a href="#">continuous and partially-observable context</a> . In <i>Pro-</i>	<b>spatial properties in SpaRC</b>	761
703	<i>ceedings of the Thirty-Third AAAI Conference on</i>		
704	<i>Artificial Intelligence and Thirty-First Innovative Ap-</i>	A symbolic context $\mathcal{C} = \{(h, r, t)_i\}_{i=1}^N$ is usually	762
705	<i>lications of Artificial Intelligence Conference and</i>	verbalized as several natural language sentences.	763
706	<i>Ninth AAAI Symposium on Educational Advances in</i>	However, we note that the <i>verbalization</i> can be a	764
707	<i>Artificial Intelligence, AAAI’19/IAAI’19/EAAI’19</i> .	conjunction of multiple tuples in a single context	765
708	AAAI Press.	sentence e.g. “Objects A and B are inside the box	766
709	Sagar Gubbi Venkatesh, Anirban Biswas, Raviteja	C”, or “Entity X is below and left of entity Y”. Such	767
710	Upadrashta, Vikram Srinivasan, Partha Talukdar, and	verbalization is common in existing benchmarks,	768
711	Bharadwaj Amrutur. 2021. <a href="#">Spatial reasoning from</a>	including SPARTUN and StepGame.	769
712	<a href="#">natural language instructions for robot manipulation</a> .		
713	In <i>2021 IEEE International Conference on Robotics</i>	<b>Fixed Orientation or Point of View (FPoV).</b>	770
714	<i>and Automation (ICRA)</i> , pages 11196–11202.	The relations are considered to axis-aligned from	771
715	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V.	a globally fixed orientation or point of view, i.e.,	772
716	Le, Ed H. Chi, Sharan Narang, Aakanksha Chowd-	fixed axes in a 2D or 3D space. We note that the	773
717	hery, and Denny Zhou. 2023. <a href="#">Self-consistency</a>	cardinal ( $\mathcal{D}_C$ ) and clock-face ( $\mathcal{D}_T$ ) directions have	774
718	<a href="#">improves chain of thought reasoning in language</a>	only 4 relations in 2D. With the set of relative direc-	775
719	<a href="#">models</a> . In <i>The Eleventh International Conference</i>	tions ( $\mathcal{D}_R$ ) being larger (6 relations in 3D), $\mathcal{D}_C$ and	776
720	<i>on Learning Representations, ICLR 2023, Kigali,</i>	$\mathcal{D}_T$ are mapped and canonicalized to four of the	777
721	<i>Rwanda, May 1-5, 2023</i> . OpenReview.net.	relative directions <i>only</i> for their label representations	778
722	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,	$\mathcal{L}$ (Table 1). Their understanding in the contexts	779
723	Barret Zoph, Sebastian Borgeaud, Dani Yogatama,	and questions is still required. Additionally, the	780
724	Maarten Bosma, Denny Zhou, Donald Metzler, Ed H.	understanding of the map to a canonicalized label	781
725	Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy	is also required to return correct answers.	782
726	Liang, Jeff Dean, and William Fedus. 2022a. <a href="#">Emer-</a>		
727	<a href="#">gent abilities of large language models</a> . <i>Transactions</i>	<b>Point Objects (PO) vs Extended Objects (EO).</b>	783
728	<i>on Machine Learning Research</i> . Survey Certifica-	Point objects are entities that are either dimension-	784
729	tion.	less i.e. their boundaries on all axes coincide, or	785
730	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	can be treated as such in a given context. Since	786
731	Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le,	they are dimensionless, in relation to other point	787
732	and Denny Zhou. 2022b. <a href="#">Chain-of-thought prompt-</a>	objects, the possible topological $\mathcal{T}_R$ relation (Table	788
733	<a href="#">ing elicits reasoning in large language models</a> . In		
734	<i>Advances in Neural Information Processing Systems</i> ,		
735	volume 35, pages 24824–24837. Curran Associates,		
736	Inc.		

1) collapses just to {DC, EQ} i.e. outside or “disconnected”, and overlapping respectively. When combined with other formalisms such as directional relations ( $\mathcal{D}$ ), even DC becomes redundant as the presence of any directional relation implies that the objects are not at the same position. Although point objects are purely mathematical constructs, real objects can often be treated as point objects in practical contexts. For example when the sizes of the objects can be ignored in relation to the distances between them. e.g. discussing spatial (directional) relations between buildings across several towns.

Extended Objects, on the other hand, are entities that are not dimensionless, i.e. their boundaries on at least one axis extends or has a spread. All real objects are extended objects. Dimensions of objects cannot be ignored when the distances between them are comparable to their sizes for spatial rule compositions and thus they must be treated as extended objects e.g. “a number of curious silver instruments” standing on Dumbledore’s “spindle-legged tables”.

**Relation Incomplete (RI) vs Relation Complete (RC).** For a set of relations  $\mathcal{R}$ , the contexts are usually relation incomplete in several real-world scenarios or when  $|\mathcal{R}|$  is relatively large. On the other hand, the contexts can be relation complete in the real-world scenarios if  $|\mathcal{R}|$  is relatively small, and one needs to emphasize and be specific.

**Quantitatively Specified (QS) vs Quantitatively Unspecified (QU).** For our current set of formalism (Table 1), some topological relations  $r \in \mathcal{T} \setminus \{EC, EQ, TPP, TPPI\}$  and all the directional relations  $r \in \mathcal{D}$  can be quantitatively specified. However, the topological relations are usually considered qualitatively, although there are metric based calculus for RCC8 and other topological relations as well. For example, context statements “Hogwarts is 200 miles to the left of the Azkaban Fortress” and “The Quidditch Stadium is inside and 1 KM away from the Hogwarts School’s northern boundary” have LEFT and NTPP (inside) as quantitatively specified relations respectively. Quantitatively specified relations that are *reverse* of each other, such as LEFT and RIGHT, can readily be composed. For example, we can infer that Harry is 2 unit right of Ron, from the context statements – Harry is 3 unit left of Hermione, and Hermione is 5 unit right of Ron. Relations are quantitatively specified when their measurements are required in

a context directly, or to infer other spatial relations indirectly.

On the other hand, for the previous examples, the context statements “Hogwarts is to the left of the Azkaban Fortress” and “The Quidditch Stadium is inside the Hogwarts School” are quantitatively unspecified for the relations LEFT and NTPP (inside) respectively. Quantitatively unspecified relations that are *reverse* of each other, such as LEFT and RIGHT, cannot be composed unless the relations are quantified. For example, directional relation between Harry and Ron cannot be determined from the context statements – Harry is left of Hermione, and Hermione is right of Ron.

### Mutual Exclusivity of Spatial Relations.

While the *reverse* relations in any formalism cannot occur simultaneously, under RCC8 calculus, multiple topological relations  $\mathcal{T}_R$  cannot occur simultaneously for the same *ordered* pair of entities even if they are not *reverse* of each others. Thus, for a given relation  $r \in \mathcal{T}_R$  and an ordered pair of entities  $(X, Y)$ :

$$r(X, Y) \implies \text{NOT}(r'(X, Y)) \forall r' \in \mathcal{T}_R \setminus r$$

For example, TPP (inside and touching) and NTPP (inside) are exclusive in RCC8. Stating a single topological relation in  $\mathcal{T}_R$  makes the context Relation Complete (RC) in (*and only in*)  $\mathcal{T}_R$ .

However, negative implications are only for *reverse* relations in directional formalism  $\mathcal{D}$ . Orthogonal relations such as LEFT and ABOVE can be *simultaneously true* for a set of *ordered* pair of entities. As directional relations are not symmetric, we will always mean an ordered pair or sequence of entities while discussing them, unless stated otherwise. Hence, Relation Incomplete (RI) contexts can be quite common in terms of directional relations.

## B Characterization of SPARTUN and StepGame

Although the existing datasets, including SPARTUN and StepGame, do not explicitly consider the spatial properties, their contexts and spatial composition rules conform to a set of these properties referenced in Section 3.1. StepGame considers entities in a completely abstract sense placed on a grid (Figure 5). They support only directional relations (including composites such as lower-left) and an overlap. Hence, objects can either be completely overlapping or completely separate. Their placement on the grid is also in terms of implicit unit of measurements. An overlap and unrestricted

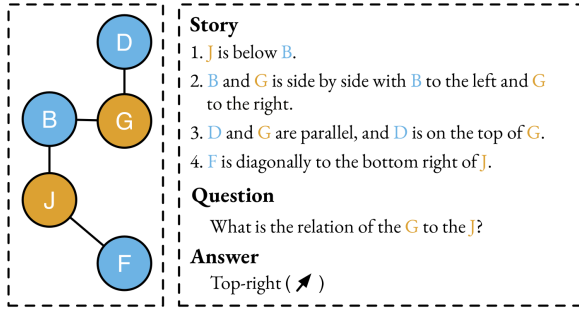


Figure 5: An example reproduced from the StepGame (Shi et al., 2022).

numerical composition of directions during their generation process coupled with the complete abstract representation of the entities essentially make them to be point objects (PO) and quantitatively specified (QS). Additionally their clear and complete expressions such as “BB is to the right of AA”, “BB is at the 3 o’clock position relative to AA”, and “AA and BB are horizontal and AA is to the right of BB” all considered equivalent means that when the relation is expressed as RIGHT, it means exactly and only RIGHT and this relation is completely known and they correspond to the relation complete (RC) context. This is why they support *strong* compositions for example presented at the beginning of Section 3 – “A is left of B and B is above C”  $\implies$  “A is to the left and above C”. Hence, the properties set of StepGame is {PO, RC, QS}.

SPARTUN on the other hand considers objects that have shapes and sizes as they built their dataset on top of NLVR images and scene graphs with different sizes of objects and blocks, and support of topological relations such as containment, inside etc. Hence, their entities are extended objects (EO). Their spatial rules (Table 3) also do not consider quantitative relations either explicitly or implicitly. Finally their spatial rules also do not make any assumption about the alignment of directional relations to be exactly parallel to an axis system. That is why a statement such as “A is to the left of B” doesn’t rule out the possibility of A additionally being above or below B i.e. the relations are not necessarily only as stated and other directional relations would still needs to be checked rather than assumed to be not present. This is in contrast with StepGame. Thus, the properties set of SPARTUN is {EO, RI, QU}. This is why applying their spatial rules (Table 3) lead to no conclusion for the previous example “A is left of B and B is above C”

presented at the beginning of Section 3. SPARTUN composition rules are thus weaker in comparison to StepGame’s composition based on these differences in their properties sets.

## C Implementation Details

### C.1 Datasets and Prompts

We created the **SpaRP-S** dataset with train, validation, and test splits of sizes 2000, 500, and 1000 respectively for each sub-dataset of SpaRP. To ensure a *fair* distribution of the difficulty level, we *randomly* sample *equal* number of instances for each *number of hops* (of entities) in the reasoning path. The final distribution is still skewed due to less number of instances for higher number of hops in SPARTUN. Additionally, we collect *five* diverse sets of *human-generated* natural language descriptions of the properties (Table 7) relevant to spatial compositions (Section 3.1). We construct a *system prompt* template with a unified task instruction and populate it with randomly sampled natural language descriptions for each instances of each sub-dataset. The system prompt template and the human-generated descriptions are presented in Table 7 through Table 13.

### C.2 Model configurations and training setup

To assess the spatial understanding and reasoning abilities of LLMs over varying model sizes, we run experiments with three model variants (all chat versions) – Llama-2-13B, Llama-2-70B, and GPT-4. The default GPT-4 i.e. GPT-4-0613 used in the experiments was accessed between Dec 1 2023 to Jan 31 2024.

We finetune a single model 13B and 70B models on all the four datasets i.e. SpaRP-S-1 (SPARTUN), SpaRP-S-PS2 (StepGame), SpaRP-S-PS3, and SpaRP-S-PS4. For finetuning, we used QLoRA (Dettmers et al., 2023) with 8-bit quantization, LoRA  $\alpha = 16$ , and LoRA config  $r = 64$ . We trained for 3 epochs with a learning rate  $lr = 1e^{-4}$ , paged AdamW optimizer, cosine  $lr$  scheduler, and an *effective* batch size of 32 using gradient accumulation.

## D Reasoning errors and their examples

We randomly sampled and manually analyzed 80 model generated reasoning paths to identify the errors and understand the discrepancy in the GPT-4 and Llama-2 70B models. A collection of several errors, their examples in terms of reasoning steps,

Terminology	Descriptions
	<p>You are an expert assistant with the knowledge of spatial relations and the rules to compose them under the assumptions that the contexts provided are of ‘{point_of_view_type}’, the objects or entities are to be treated as ‘{entity_type}’, the directions are ‘{quantitative_type}’, and ‘{relation_type}’. The description of these terminologies are as given below:</p> <p>{point_of_view_type}: {point_of_view_type_desc}{point_of_view_type_default}</p>
System Instruction Template	<p>{entity_type}: {entity_type_desc}{entity_type_default}</p> <p>{quantitative_type}: {quantitative_type_desc}{quantitative_type_default}</p> <p>{relation_type}: {relation_type_desc}{relation_type_default}</p>
	<p>You need to identify the sub-set of entities from the context that are relevant as well as combine their spatial relations with valid compositions under the above mentioned assumptions to find the spatial relations between the entities in the asked questions. The list of all possible spatial relations are: {spatial_relation_choices}. Always provide the final answer, only and only, in terms of these spatial relations. Include all the spatial relations that hold true as the answer, in case of multiple correct choices.</p>
Fixed Orientation Point of View	<p>The spatial relations are expressed from a single, consistent and unchanging perspective. This means that the observations are made from a global viewpoint that remains same and constant for all the entities in a given context. Hence, relations such as relative directions e.g. left or right always refer to the same directions and there is a one-to-one mapping between relative, cardinal and clock-face directions i.e. left is same as west or 9 o’clock position, right is same as east or 3 o’clock position, above is same as north or 12 o’clock position, and below is same as south or 6 o’clock position.</p>
Implicit Quantification	<p>Unless otherwise stated, consider the direction relations specified in the context to be of 1 unit distance. For example, the sentence, entity X is to the lower-left of entity Y means that the entity X is 1 unit to the left and 1 unit below the entity Y.</p>

Table 7: Human-generated natural language descriptions for common terminologies, defaults and system instruction. Terms inside {} are placeholders that are further replaced with their language descriptions. For current work, point\_of\_view\_type is always Fixed Orientation Point of View and the only default available is for quantitative\_type = Quantitatively Specified (QS) with quantitative\_type\_default = Implicit Quantification. All other placeholders are replaced by randomly sampled descriptions from one of their 5 diverse human-generated descriptions presented in Table 8 through Table 13. The spatial\_relation\_choices are the relevant labels  $\mathcal{L}$  from Table 1.

Diverse human-generated Descriptions for Point Objects (PO)
<p><b>Description 1:</b> Two objects can be treated as Point Objects in a given context for specifying their spatial relations if they are extremely small such that their sizes are immaterial, or if they are of similar or even varying shapes and sizes but are placed sufficiently far enough that their shapes and sizes can be ignored to state and compose spatial relations between them. This leads to a limitation on the spatial relations that can be specified between objects e.g. containment, but simpler relation compositions since shapes and sizes of the objects need not be considered. For example, a tea-cup and an apple on a table, or a school building and a warehouse that are miles away can be considered as point objects.</p> <p><b>Description 2:</b> While composing spatial relations between objects, they can be considered as Point Objects if they can be treated as dimensionless i.e. if (1) their sizes are so small that they can be neglected or (2) the size and shape of the objects are negligible compared to the great distance between the objects. Although this situation may prevent to express certain relations like containment, it provides simpler spatial relation statements and compositions over multiple objects, since the size and shape are not considered. For example, two balls on the basketball pitch or two buildings that are separated with 2 KM distance.</p> <p><b>Description 3:</b> Point Objects are small objects in a given context, whose sizes and shapes can be ignored. Thus, only their locations and orientations are considered when specifying spatial relations, leading to less number of relations and their simpler combinations over objects. A typical example of point objects can be buildings on a map or beads on a table.</p> <p><b>Description 4:</b> In this scenario, objects can be treated as Point Objects if they are extremely small or far apart to the extent that their shapes and sizes can be ignored. In such cases, certain spatial relationships, like containment, become inapplicable. Additionally, since the shapes and sizes of the objects are not important, relationship compositions can be simpler. For example, two cars that are miles apart can be considered as point objects.</p> <p><b>Description 5:</b> Entities can be treated as Point Objects when the distance between them relative to their sizes is either large or can be ignored. Therefore when providing spatial relations between them, a limited set of relations with simpler composition rules is possible. For example, when someone says, a cafe and a house that are far apart can be treated as point objects.</p>

Table 8: Five diverse human-generated natural language descriptions of Point Objects (PO).

---

### Diverse human-generated Descriptions for Extended Objects (EO)

---

**Description 1:** Two objects are to be treated as Extended Objects in a given context for specifying their spatial relations if their shapes and sizes in comparison to the distances between them can not be ignored to state and compose spatial relations between them. This leads to more number of possible spatial relations that can be specified between objects e.g. containment, but reduces the number while increasing the complexity of possible relation compositions, as the shapes and sizes of the objects can neither be assumed nor be discarded. For example, a tea-cup and a tube-light, or a table and a cupboard in a room are to be considered as extended objects.

**Description 2:** If the distance between objects is comparable to the shapes and sizes of the objects while specifying the spatial relations, the objects are considered as Extended Objects i.e. they can't be treated as dimensionless and they have significant length, breadth or height in comparison to the distances between the objects in the context. Although this gives an opportunity to use more specific spatial relations like touching or containment, the complexity of compositions increases. A basket and an apple in it or two entities, X and Y, in a room can be given as examples.

**Description 3:** Extended Objects refer to objects, whose shapes and sizes can affect the spatial relations that can be specified and the way they can be combined between objects. This leads to more number of relations and the combination of relations have to be minimal in the absence of the information about the shape and size of the objects. Examples of extended objects include buildings on a street or boxes in a room.

**Description 4:** In this scenario, two objects are considered to be Extended Objects if their shapes and sizes, in comparison to the distances between them, cannot be ignored. In such cases, a larger set of spatial relations between objects can be specified, although the relation composition becomes more limited when the shapes and sizes of the objects are unknown compared to when this information is known. For example, a tea-cup and a lamp or a sofa and a TV in a room can be considered as extended objects.

**Description 5:** Entities can be treated as Extended Objects if they have shapes and sizes which are not to be ignored in the context. Because of this, although a larger set of relations is possible between objects but the composition rules can become complex. For example, a cafe and a mall building can be treated as extended objects and the cafe can be a part of i.e. inside the mall building itself.

---

Table 9: Five diverse human-generated natural language descriptions of Extended Objects (EO).

---

### Diverse human-generated Descriptions for Relation Incomplete (RI) contexts

---

**Description 1:** Not all set of possible spatial relations that hold true between two objects are stated while specifying the relations between those objects. Thus, there could be multiple possible spatial configurations that conform to the stated relations between the objects. For example, the statement, object A is to the left of object B, when considered as relation incomplete could mean that A may or may not be strictly only to the left of B, i.e. it can be either only to the left, or is to the left and above, or is to the left and below B.

**Description 2:** Although some spatial relations between two objects exist, they might be overlooked while expressing the relations between those objects. Therefore, other valid configurations, which are compatible with the expression but not explicitly specified, may also exist. For instance, the relation incomplete expression, the entity X is to the left of the entity Y does not have to mean that X is to the left of Y and they are strictly aligned at the same time. The entity X can be both to the left and bottom (or above etc.) of the entity Y.

**Description 3:** An incomplete spatial relationship corresponds to the insufficient information or context to decide the exact spatial relationship between objects, leading to ambiguity. In other words, there can be multiple valid spatial arrangements or layouts that hold true to each incomplete relation. For example, given the incomplete statement that box 'one' is in front of box 'two', it holds true for multiple arrangements such as box 'one' is to the right and front of box 'two', or box 'one' is to the left and front of box 'two'.

**Description 4:** Relations are incomplete in the context statements if not all the spatial relationships that exist between two objects are stated. In such cases, multiple spatial outline or positioning of the objects are possible, without a single definitive truth. For example, consider the relationship - the fruit F is behind the object O in a 2D plane. Although O is in front of F, their relative position on the horizontal axis is incomplete, and hence, could be left, right or at the same place when considered horizontally.

**Description 5:** The provided set of spatial relations between two objects may not be enough to communicate the complete spatial position between them. Therefore, for the provided spatial information between two objects more than one arrangement is possible. For example, a metal ball is hanging below a metal beam in the workshop - can mean various spatial positions such as the metal ball is below the beam, or additionally, it can be to the right or left and away from the beam in consideration.

---

Table 10: Five diverse human-generated natural language descriptions of Relation Incomplete (RI).

---

### Diverse human-generated Descriptions for Relation Complete (RC) contexts

---

**Description 1:** All set of possible spatial relations that hold true between two objects are stated while specifying the relations between those objects. Hence, there is only one spatial configuration that conforms to the stated relations between the objects. For example, the statement, object A is to the left of object B, when considered as complete could only and only mean that A is to the left of B.

**Description 2:** All existing spatial relations between two objects are included while expressing the relations. Therefore, there is one-to-one mapping between spatial configurations and expressed spatial relations between objects. For instance, a relation complete statement, the entity X is to the right of the entity Y means that X is to the right of Y and they are aligned.

**Description 3:** Completely specified spatial relations refer to the complete sets of spatial relations that can be held as well as stated between objects. Thus, there can be only one valid spatial arrangement or layout that holds true for a relation complete statement. An example is that box 'one' is in front of box 'two' and they are in the same line that denotes front in a given fixed orientation for all.

**Description 4:** Relations are complete in a setting, if all the spatial relationships between two objects are stated. In such cases, there is a single ground truth spatial outline or positioning of the objects. For example, consider the relationship - the fruit F is behind the object O in a 2D plane. This means that O is strictly and only in front of F and are aligned on the axis i.e. can be considered to be neither left nor right but at the same position on the horizontal axis.

**Description 5:** The provided set of relations between two objects are enough to know the actual spatial position between them. Therefore, no additional information is needed to understand the actual position between two objects. For example, a metal ball is hanging below a metal beam in the workshop means that the ball is below the beam and not to its left or right.

---

Table 11: Five diverse human-generated natural language descriptions of Relation Complete (RC).

---

### Diverse human-generated Descriptions for Quantitatively Specified (QS) relations

---

**Description 1:** Spatial relations, such as directions, specified between two objects are said to be Quantitatively Specified if those relations can have a unit of measurement and are also stated, implicitly or explicitly, in the specified context. The composition of such relations is always possible when all the object parameters and the relations between any two objects in a statement are completely known. For example, with constraints such as objects A, B and C are apples lying in a line and the relation specified are of 1 unit measurement when not mentioned explicitly, the quantitatively specified statements - B is 3 units to the left of A, and C is to the right of B - can lead to the only conclusion that A is 2 units to the right of C, or its inverse equivalent i.e. C is 2 units to the left of A.

**Description 2:** Unit of measurements in spatial relations (e.g., directions) between two objects needs to be explicitly or implicitly specified for these relations to be called as Quantitatively Specified. The composition of such relations can be determined when all other object parameters and relations of two objects are given. For example, let entities X and Z be perfect round shaped balls. Let entity Y be a round basket with 10 unit radius and let centers of all objects are horizontally aligned. If X is 1 unit to the left of the center of Y and Z is 2 units to the left of X, then Z is inside the basket and 3 units to the left of the center of the basket Y.

**Description 3:** Spatial relations are Quantitatively Specified when these relations are defined with a specific unit of measurement such as meters or miles. The relation compositions over objects become deterministic if all the other object parameters and the relationships between them are provided. For example, box 'one' is 3 units above box 'two' and they are in the same line can be easily used to determine relations with respect to a third box, say box 'three', if its position is also quantitatively specified with one of them.

**Description 4:** Under this setting, spatial relations between two objects are said to be Quantitatively Specified if the relations have a unit of measurement and stated directly or indirectly in the context. In such cases, when all the object parameters and relations between any two objects in the statement are known, a deterministic relation composition is possible. For example, although there are limitations like having three apples (A1, A2, A3) arranged in a row, the statements - A2 is 2 units left of A1, and A3 is 1 unit right of A2 - provides enough information to determine the exact positions of A1 and A3 relative to each other.

**Description 5:** If the quantitative value along with the measurement unit for a spatial relation is provided then those relations are said to be Quantitatively Specified. The measurements may be a default value that is understood in the context or is explicitly provided. The composition of these relations will result in a distinctly resolved relation. For example, in the sentence, the cafe is 2 blocks north of my house and the hospital is 1 block south of the cafe, it can be easily determined that the hospital is 1 block north of my house.

---

Table 12: Five diverse human-generated natural language descriptions of Quantitatively Specified (QS) relations.

---

### Diverse human-generated Descriptions for Quantitatively Unspecified (QU) relations

---

**Description 1:** Spatial relations, such as directions, specified between two objects are said to be quantitatively unspecified if those relations can have a unit of measurement but are not stated in the specified context. The composition of such relations may not be possible even when all the object parameters and the relations between any two objects in a statement are completely known. For example, even with constraints such as objects A, B and C are apples lying in a line, the quantitatively unspecified statements - B is to the left of A, and C is to the right of B - can not lead to any conclusion regarding left, right, or overlapping relationship between A and C.

**Description 2:** In order for spatial relations between two objects to be considered as quantitatively unspecified, unit of measurement in these relations should not be specified. The exact composition or realization of such relations may not be determined even if the other object features and relations are completely known. For example, let entity X be in the basket Y of a known and stated size, and let entity Z be to the right of the entity X. It is not possible to infer whether entity Z is in the basket Y or not if its distance from X is quantitatively unspecified.

**Description 3:** Spatial relations are Quantitatively Unspecified when these relations are not defined in terms of specific units of measurement such as meters or miles. The relation compositions over objects can still not be determined even if all the other object parameters and the relationships between them are provided. For example, if box 'one' is above box 'two', it's not clear how far exactly box 'one' lies with respect to box 'two' and this will affect the conclusions to be drawn about relations with respect to other objects, say box 'three'.

**Description 4:** In this setting, spatial relations between two objects are Quantitatively Unspecified if the relations have a unit of measurement that is not specified in the context. In such cases, even when all the object parameters and relations between any two objects in the statement are known, a deterministic composition of relations may be impossible. In this scenario, although there are limitations like having three apples (A1, A2, A3) arranged in a row, the statements that lack specific quantities - A2 is on the left of A1, and A3 is on the right of A2 - do not provide enough information to determine the left, right, or overlapping positions of A1 and A3 relative to each other.

**Description 5:** If the quantitative value along with the measurement unit for a spatial relation is not provided then those relations are said to be quantitatively unspecified. The composition of these relations may not be enough to result in a distinctly resolved relation. For example, in the sentence, the cafe is to the north of my house and the hospital is to the south of the cafe, it can't be determined if the hospital is to the south or north of my house.

---

Table 13: Five diverse human-generated natural language descriptions of Quantitatively Unspecified (QU) relations.

977 the datasets to which the generated paths belong  
978 and the explanation of the errors are provided in  
979 Table 14.



---

### Errors examples and explanations

---

<p><b>Error:</b> Reverse Step Merge <b>Dataset:</b> SpaRP-S-PS2(StepGame)  <b>Reasoning Steps:</b> <i>Step 7:</i> From step 4 and 6, we can infer that <b>I</b> is 2 unit above and 2 unit right of <b>W</b>. <i>Step 9:</i> From step 8, we can infer that <b>Q</b> is right of <b>W</b>. <i>Step 10:</i> From step 7 and 9, we can infer that <b>I</b> is 2 unit above and <b>3 unit right</b> of <b>Q</b>.  <b>Explanation:</b> Requires composition from I to W to Q. Step 9 is a reverse step in this reasoning chain causing composition error.</p>
<p><b>Error:</b> Copied, not composed <b>Dataset:</b> SpaRP-S-PS2 (StepGame)  <b>Reasoning Steps:</b> <i>Step 6:</i> From step 4 and 5, we can infer that <b>A</b> is below and right of <b>S</b>. <i>Step 7:</i> From the context, <b>S</b> is <b>left</b> of <b>M</b>. <i>Step 8:</i> From step 6 and 7, we can infer that <b>A</b> is below and right of <b>M</b>.  <b>Explanation:</b> The relation from S to M is not used in composition, instead, the relation from A to S is copied for the merge step A to S to M.</p>
<p><b>Error:</b> Incorrect Composition <b>Dataset:</b> SpaRP-S-PS2 (StepGame)  <b>Reasoning Steps:</b> <i>Step 11:</i> From step 8 and 10, we can infer that <b>Z</b> is above and 2 unit right of <b>L</b>. <i>Step 13:</i> From step 12, we can infer that <b>L</b> is above and left of <b>J</b>. <i>Step 14:</i> From step 11 and 13, we can infer that <b>Z</b> is <b>3 unit above</b> and right of <b>J</b>.  <b>Explanation:</b> Quantitatively incorrect.</p>
<p><b>Error:</b> Merged without connection <b>Dataset:</b> SpaRP-S-PS3  <b>Reasoning Steps:</b> <i>Step 5:</i> From step 3 and 4, we can infer that <b>Y</b> is right of <b>L</b>. <i>Step 14:</i> From step 12 and 13, we can infer that <b>C</b> is below and right of <b>K</b>. <i>Step 15:</i> From step 5 and 14, we can infer that <b>C</b> is below and right of <b>L</b>.  <b>Explanation:</b> No common entity between merged steps 5 and 14 which are 9 steps apart.</p>
<p><b>Error:</b> Non-composable Composition <b>Dataset:</b> SpaRP-S-PS4  <b>Reasoning Steps:</b> <i>Step 16:</i> From step 13 and 15, we can infer that <b>M</b> is below <b>X</b>. <i>Step 18:</i> From step 17, we can say that <b>X</b> is above <b>A</b>. <i>Step 19:</i> From step 16 and 18, we can infer that <b>M</b> is above <b>A</b>.  <b>Explanation:</b> The reverse relations {below, above} are not composable under quantitatively unspecified (QU) criteria.</p>
<p><b>Error:</b> Reverse answer <b>Dataset:</b> SpaRP-S-PS3  <b>Question:</b> What is the relation of the agent <b>W</b> to the agent <b>X</b>?  <b>Answer Step:</b> <i>Step 8:</i> From step 5 and 7, we can infer that <b>X</b> is above and left of <b>W</b>. Hence, the answer is above, and left.  <b>Explanation:</b> Directional relations are non-symmetric. Question is from W to X, while answer is from X to W.</p>
<p><b>Error:</b> Incorrect clock extraction <b>Dataset:</b> SpaRP-S-PS3  <b>Context:</b> <b>P</b> is above <b>X</b> at 10 o'clock.  <b>Extracted:</b> <i>Step 6:</i> From the context, <b>X</b> is <b>below P</b>.  <b>Explanation:</b> Correct extraction should be below and right.</p>
<p><b>Error:</b> Incomplete relation extraction <b>Dataset:</b> SpaRP-S-PS4  <b>Context:</b> <b>I</b> is on the left side of and <b>below V</b>.  <b>Extracted:</b> <i>Step 2:</i> From the context, I is left of V.  <b>Explanation:</b> Incomplete or partial extraction, below wasn't picked.</p>
<p><b>Error:</b> Incorrect relation extraction <b>Dataset:</b> SpaRP-S-PS1 (SPARTUN)  <b>Context:</b> <b>Box EEE</b> has a <b>tiny white rectangle</b> and covers a midsize white diamond.  <b>Extracted:</b> <i>Step 7:</i> It is given that the <b>tiny white rectangle</b> is <b>inside and touching</b> the <b>box EEE</b>.  <b>Explanation:</b> <i>Has</i> only means inside.</p>
<p><b>Error:</b> Incorrect reversal <b>Dataset:</b> SpaRP-S-PS1 (SPARTUN)  <b>Reasoning Steps:</b> <i>Step 1:</i> It is given that the <b>block AAA</b> contains the <b>medium yellow square</b>. <i>Step 2:</i> From step 1, we can infer that the <b>medium yellow square</b> is <b>inside and touching</b> the <b>block AAA</b>.  <b>Explanation:</b> Reverse of <i>contains</i> should only be inside.</p>

Table 14: Observed errors and their examples in the model generated reasoning paths. Only the relevant steps are shown.