

# Where the Cat Sat: A Multilingual Benchmark for Spatial Language Understanding

Anonymous ACL submission

## Abstract

Spatial language understanding is fundamental to human communication and reasoning, enabling tasks from robot navigation to document analysis and geographic information systems. Current spatial language understanding exhibits biases toward English and in particular prepositional marking. We present a novel framework for spatial language understanding and a multilingual benchmark decomposing spatial relations into surface elements (figure, ground, predicate, markers) and semantic components (dynamicity, stasis). Evaluating frontier Large Language Models (LLMs) on Spanish, Basque, and Chinese, we find high accuracy on surface element identification but persistent gaps in semantic classification. Basque case affixes remain most challenging—small models achieve as low as 15.3% on spatial markers—suggesting morphological complexity poses difficulties even for large models. These results suggest that surface parsing does not entail spatial understanding, and that evaluation must include languages with diverse spatial marking strategies beyond prepositions.

## 1 Introduction

Spatial language understanding remains a fundamental challenge in NLP, as evidenced by ongoing workshop series dedicated to the topic (Kordjamshidi et al., 2018; Bhatia et al., 2019; Kordjamshidi et al., 2020; Alikhani et al., 2021; Kordjamshidi et al., 2024). Current approaches exhibit systematic biases toward English and well-resourced Indo-European languages (Ulinski et al., 2019; Olek and Piasecki, 2024), limiting our understanding of how well computational models genuinely comprehend spatial relations across typologically diverse languages.

Existing work exhibits four limitations: spatial markers are understood predominantly as prepositions, overlooking case marking and spatial nouns;

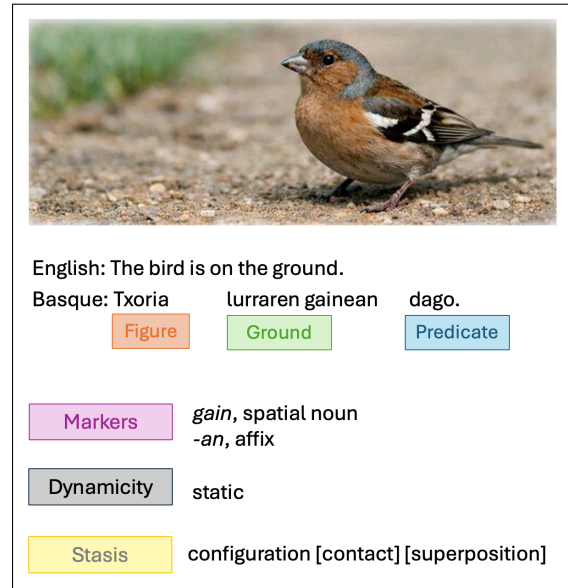


Figure 1: Example annotation for a Basque sentence describing a bird on the ground. Surface elements: *Txoria* (‘the bird’) is the figure; *lurraren gainean* (‘the ground’s surface’) is the ground; *dago* (‘is’) is the spatial\_predicate, classified as copula; *gain* (‘surface’) and *-an* are spatial\_markers, a spatial noun and affix respectively. Semantic components: dynamicity is static; configuration is [+contact][superposition].

spatial relations are labeled using English prepositions as categories rather than decomposed into semantic primitives (Liu et al., 2025; Beekhuizen, 2025), hampering analysis and comparison across languages; relations requiring frames of reference (Levinson and Wilkins, 2006) are ignored or conflated with proximity-based ones; and finally spatial relations that involve motion have been rarely studied systematically.

This paper presents a multilingual benchmark for spatial language understanding that addresses these limitations, covering three typologically diverse languages: Spanish (Indo-European), Basque (Isolated), and Chinese (Sino-Tibetan).

056	Our evaluation framework identifies six compo-	multiple marker types beyond prepositions (i.e. ad-	105
057	nents: four surface elements ( <code>figure</code> , <code>ground</code> ,	positions, affixes, spatial nouns) as well as predi-	106
058	<code>spatial_predicate</code> , <code>spatial_markers</code> ) and two	cates, while decomposing semantics into topolog-	107
059	semantic components ( <code>dynamicity</code> , <code>stasis</code> ), with	ical versus projective relations using primitives	108
060	<code>stasis</code> divided into <code>configuration</code> (topological	rather than English prepositions, with systematic	109
061	primitives) and <code>projective</code> (frame-of-reference	<code>dynamicity</code> annotation. Next, we detail each com-	110
062	distinctions), the definitions of which are de-	ponent.	111
063	scribed in Section 2.2.		
064	Our contributions include: a novel framework	<b>2.2 Surface Elements</b>	112
065	for spatial language understanding grounded in	<b>Figure.</b> The noun phrase denoting the entity	113
066	cross-linguistic typology; a multilingual bench-	whose location is being described (e.g., <i>Txoria</i> ‘the	114
067	mark covering Spanish, Basque, and Chinese	bird’). We extract the complete noun phrase includ-	115
068	with gold-standard annotations; recognition of	ing all determiners, modifiers, and complements.	116
069	<code>spatial_markers</code> beyond prepositions; primitive		
070	decomposition for topological relations and ex-	<b>Ground.</b> The reference object relative to which	117
071	PLICIT frame-of-reference annotation; and system-	the <code>figure</code> ’s location is specified (e.g., <i>lurraren</i>	118
072	atic coverage of <code>dynamicity</code> × <code>stasis</code> combina-	<i>gainean</i> ‘the ground’s surface’). We include spa-	119
073	tions.	tial nouns like <i>gain</i> when they function as part of	120
074	Our evaluation reveals systematic limitations in	the <code>ground</code> phrase, essential for languages where	121
075	current LLMs. Basque <code>spatial_markers</code> prove	spatial nouns are integral to spatial description	122
076	particularly challenging across model scales, with	(Levinson and Meira, 2003).	123
077	small models showing severe difficulties and fron-		
078	tier models demonstrating variable performance—	<b>Spatial Predicate.</b> The predicate expressing the	124
079	often the weakest compared to Spanish and Chi-	spatial relation linking <code>figure</code> and <code>ground</code> (e.g.,	125
080	nese markers. This systematic challenge appears	<i>dago</i> ‘exists/is’, type: copula). We distinguish <i>verb</i>	126
081	to stem from these models’ inability to reliably seg-	(i.e. content/lexical verbs) from <i>copula</i> (i.e. gram-	127
082	ment and classify case affixes as spatial markers.	matical elements that link <code>figure</code> to <code>ground</code> ).	128
083	Across languages, models perform better on sur-		
084	face elements than semantic components.	<b>Spatial Markers.</b> Elements that encode the spa-	129
		tial relationship between <code>figure</code> and <code>ground</code> . In	130
085	<b>2 A Multilingual Framework for Spatial</b>	the example (Figure 1), we identify two markers:	131
086	<b>Language Understanding</b>	<i>gain</i> (type: spatial noun) and <i>-an</i> (type: affix, the	132
		locative case marker). We recognize three types:	133
087	<b>2.1 Framework Overview</b>	adposition (prepositions, postpositions, or gram-	134
088	Our framework decomposes spatial relations into	maticalized multi-word constructions), affix (case	135
089	six components organized along two dimensions.	markers or morphological affixes encoding spatial	136
090	Figure 1 shows an example annotation for a Basque	meaning (Creissels, 2008)), and spatial noun (lexi-	137
091	sentence describing a spatial scene. We first de-	cal nouns with spatial meaning functioning within	138
092	fine the surface elements that must be extracted	the <code>ground</code> phrase).	139
093	from text, then specify the semantic components		
094	that must be inferred.	<b>2.3 Semantic Components</b>	140
095	<b>Surface Elements.</b> Linguistic forms present in	<b>Dynamicity.</b> The temporal aspect of the spa-	141
096	the sentence that models must identify by pars-	tial relation. Three values are possible: static	142
097	ing: <code>figure</code> , <code>ground</code> , <code>spatial_predicate</code> , and	(the <code>figure</code> ’s spatial relation to the <code>ground</code> re-	143
098	<code>spatial_markers</code> .	mains constant), source (the <code>figure</code> moves away	144
099	<b>Semantic Components.</b> Conceptual categories	from the <code>ground</code> ), and goal (the <code>figure</code> moves to-	145
100	not explicitly marked that models must infer:	ward the <code>ground</code> ). In our example, the value is	146
101	<code>dynamicity</code> and <code>stasis</code> (type of spatial relation-	static: the bird is not in motion. However, dif-	147
102	ship: topological or projective).	ferent visual stimuli can elicit sentences with dif-	148
103	This organization addresses limitations of exist-	ferent <code>dynamicity</code> values—“The bird lands on the	149
104	ing frameworks by expanding surface extraction to	branch” would be goal, while “The bird flies from	150
		the branch” would be source.	151

**Stasis.** The overall type of spatial relationship between figure and ground. Following [Levinson and Wilkins \(2006\)](#), we distinguish two spatial domains: configuration (topological) and projective relations. Table 4 (Appendix B) provides a complete mapping of all primitive values to representative English sentences for reference.

**Configuration** These relations involve contiguity or close proximity between figure and ground, without requiring directional specification. Rather than labeling these with atomic categories, we characterize them through semantic primitives. In our example, the value is [+contact][superposition], indicating “the bird” is in physical contact with the ground and positioned above it. All configuration values specify contact status ([+contact] or [-contact]). Beyond contact, we identify four configuration types: containment (figure within boundaries of ground, specified as [open] or [closed]), attachment (figure mechanically fastened to ground), superposition (figure superior to ground, optionally [top]), and subposition (figure inferior to ground). This primitive-based approach contrasts with existing work, which labels topological relations using English prepositions as categories ([Liu et al., 2025](#); [Beekhuizen, 2025](#)).

**Projective** These relations apply when figure and ground are separated in space and directional specification becomes necessary. Unlike topological relations, projective relations require external coordinate systems—frames of reference. Following [Levinson and Wilkins \(2006\)](#), we distinguish three frame types: relative (coordinates from the observer’s bodily axes; values: left, right, front, back); intrinsic (coordinates from inherent facets of the ground object; values: front, back); and absolute (fixed environmental bearings; values: north, south, east, west). Appendix A.2 provides annotated examples of projective relations.

## 2.4 Limitations of Existing Approaches

Current spatial language evaluation frameworks were designed primarily for English and exhibit systematic biases that limit their applicability to typologically diverse languages. In terms of surface form, existing approaches extract only figure + preposition + ground, overlooking the broader range of adpositions and other encoding strategies such as case marking ([Creissels, 2008](#)), spatial

nouns, or zero marking ([Haspelmath, 2019](#); [Stolz et al., 2014](#)). Spatial nouns tend to be ignored as contributing elements to spatial relations, despite some languages relying heavily or entirely on such elements ([Levinson and Meira, 2003](#); [Lakoff, 1987](#)). Spatial predicates, which contribute to the overall spatial meaning, are also often excluded from evaluation.

Beyond surface elements, semantic evaluation presents additional challenges. Configuration and projective relations are typically amalgamated, with projective relations often ignored entirely or conflated with topological ones ([Levinson and Wilkins, 2006](#)). The semantic categories themselves are problematic: topological relations are labeled using English prepositions as theoretical categories (e.g., “on”, “in”, “above”) ([Liu et al., 2025](#); [Beekhuizen, 2025](#)) rather than being decomposed into semantic primitives that could apply across languages. Finally, dynamicity—whether a spatial relation is static or involves motion—is rarely studied systematically, and the interaction between dynamicity and spatial configuration is almost never evaluated jointly.

## 3 Multilingual Data Collection

We now demonstrate the applicability of this framework through data collection across three typologically diverse languages.

### 3.1 Language Selection

Spanish (Indo-European) encodes spatial relations primarily through prepositions, representing the marking pattern typical of major European languages. Basque (isolate) employs agglutinative morphology with spatial cases combined with spatial nouns, thus allowing us to evaluate LLM’s ability to recognize affixes as spatial markers. Chinese (Sino-Tibetan) distributes spatial meaning across prepositions and localizers.<sup>1</sup>

### 3.2 Visual Stimuli Design

Following established methods for cross-linguistic spatial elicitation ([Bowerman and Pederson, 1992](#)), we created canonical images depicting spatial relations to elicit natural descriptions from native speakers.<sup>2</sup>

<sup>1</sup>Language-specific analytical decisions—regarding marker boundaries, contracted forms, and the grammatical status of elements such as Chinese localizers—are detailed in Appendix D.

<sup>2</sup>Both images and gold standard annotation will be made publicly available through a GitHub repository upon publica-

246	Our stimuli target specific primitive values for	These models represent current frontier capabilities	291
247	topological relations (e.g., contact vs. non-contact,	and serve as the primary benchmark for our	292
248	open vs. closed containment), include projective	evaluation.	293
249	relations covering all three frames of reference (rel-	<b>Small models.</b> We additionally evaluate	294
250	ative, intrinsic, absolute), and systematically vary	three smaller multimodal models: Claude-3.5-	295
251	dynamycity (static, source, goal) for each spatial	Haiku (Anthropic, 2024), GPT-4o-mini (OpenAI	296
252	relation. This design yields systematic coverage of	et al., 2024), and Qwen2.5-VL-7B-Instruct (Bai	297
253	the stasis $\times$ dynamycity space: 63 unique com-	et al., 2025). All models are presented as multi-	298
254	binations in total.	lingual, though language support details vary by	299
255	<b>3.3 Elicitation Procedure</b>	provider.	300
256	Data collection proceeded in two phases.	This dual-scale design allows us to distinguish	301
257	<b>Pilot phase.</b> For each language, we designed a	limitations that reflect genuine difficulty in spatial	302
258	pilot questionnaire with 12 images covering repre-	understanding from limitations that reflect model	303
259	sentative spatial relations: 1 configuration type	capacity. Patterns of failure that appear at both	304
260	$\times$ 3 dynamycity values, plus 9 projective re-	scales suggest fundamental challenges; patterns	305
261	lations (3 frames $\times$ 3 dynamycity values). Na-	that appear only at smaller scales suggest capacity-	306
262	ative speaker experts produced sentences describing	dependent limitations.	307
263	each image and segmented the grammatical com-	Models were evaluated with text-only input (sen-	308
264	ponents (figure, ground, spatial_predicate,	tence alone). A comparative multimodal condition	309
265	spatial_markers). Informed by our investigation	(sentence plus image) was run to verify whether	310
266	on spatial literature, these initial parses informed	visual grounding could address observed semantic	311
267	the development of language-specific parsing cri-	difficulties.	312
268	teria.	<b>4.2 Prompting Strategy</b>	313
269	<b>Full phase.</b> The procedure was expanded to	To isolate the understanding of specific spatial	314
270	cover all 63 stasis $\times$ dynamycity combinations.	components and mitigate error propagation, we	315
271	To reduce annotator burden, we simplified the task:	employ a component-wise prompting strategy.	316
272	native speakers produced only sentences, without	Rather than requesting a single complex JSON ob-	317
273	grammatical analysis. Segmentation was then per-	ject, the model is queried sequentially for each of	318
274	formed applying the criteria established in the pilot	the six components.	319
275	phase.	For each inference call, the prompt comprised	320
276	<b>Verification.</b> A second native speaker expert re-	of: (i) the operational definition for that com-	321
277	viewed all analyses to ensure consistency with the	ponent (identical to the annotation guidelines in	322
278	operational definitions.	Appendix C), and (ii) the target sentence. This	323
279	<b>3.4 Resulting Gold Standard</b>	approach ensured that evaluation measured the	324
280	The final parallel dataset comprises 189 sentences	model’s ability to apply schema definitions rather	325
281	across the three languages, with gold-standard an-	than its ability to maintain long-context coherence.	326
282	notations for all six spatial relation components.	<b>4.3 Inference Configuration</b>	327
283	<b>4 Experimental Setup</b>	Inference was performed using the OpenRouter	328
284	<b>4.1 Model Selection</b>	API <sup>3</sup> for all models, with temperature 0.2 and max-	329
285	We evaluate two groups of models to assess the ef-	imum token limit of 8,192. To account for stochas-	330
286	fect of scale on spatial language understanding.	tic variance, we conducted three independent runs	331
287	<b>Frontier models.</b> We evaluate three large-scale	per model-language pair.	332
288	multimodal models: Claude-3.5-Sonnet (An-	<b>4.4 Evaluation Metrics</b>	333
289	thropic, 2024), GPT-4o (OpenAI et al., 2024),	Our evaluation employs granular scoring that re-	334
290	and Qwen2.5-VL-72B-Instruct (Bai et al., 2025).	wards partial understanding while penalizing hal-	335
	tion.	lucinations. All scores range from 0.0 to 1.0 and	336
		are reported as percentages in tables.	337
		<sup>3</sup> <a href="https://openrouter.ai">https://openrouter.ai</a>	

**Figure and Ground.** Evaluated using normalized Levenshtein similarity between gold and predicted text spans, yielding scores from 0.0 (no match) to 1.0 (exact match).

**Spatial Predicate.** Evaluated as the average of two components: (1) normalized Levenshtein similarity between gold and predicted text, and (2) binary match for predicate type (verb vs. copula).

**Spatial Markers.** A position-based metric that evaluates each marker independently, then averages across all positions (using the maximum of gold/predicted marker counts). For each marker, we compute the average of: (1) normalized Levenshtein similarity between marker texts (with length penalty for overprediction), and (2) binary match for marker type (adposition, affix, or spatial noun). Missing markers score 0.0.

**Configuration.** Features extracted from bracket notation (e.g., [+contact], [superposition]) are evaluated using F1-score, treating each primitive as an independent feature.

**Projective.** Frame type and directional value are treated as independent features and evaluated using F1-score.

**Dynamicity.** Exact match accuracy over three categories: static, source, and goal.

## 5 Results

### 5.1 Overall Performance

Table 1 summarizes performance across all components and languages for frontier models. These models achieve strong performance on surface elements across languages, with figure and ground identification generally high (83.0%–100%). `Spatial_predicate` identification remains strong for Spanish and Basque but lower for Chinese. Semantic components show more variation: configuration and projective accuracy is lower than surface element extraction.

Small models show lower and more variable performance compared to frontier models (see Table 6 in Appendix E for small model comparisons). Figure and ground identification remain relatively stable for Claude and Qwen-7B, though GPT drops to 40.8% on Basque figure extraction. The most pronounced difference between scales appears in semantic components and Basque `spatial_markers`.

Across both model scales, surface elements outperform semantic components, with this gap particularly pronounced for Chinese in frontier models.

### 5.2 Spatial Markers

Table 1 shows `spatial_markers` accuracy for frontier models. `Spatial_markers` extraction reveals cross-linguistic asymmetry in our evaluation. For frontier models, marker accuracy varies substantially: Claude maintains comparable performance across languages, while GPT and Qwen show notably lower accuracy for Basque markers despite strong performance on other surface elements.

Small models demonstrate larger gaps in performance (see Table 6 in Appendix E for comparison). Marker accuracy decreases across all languages, with Basque markers showing the most dramatic decline. Basque marker extraction achieves the lowest accuracy among marker types across all small models (15.3%–29.6% vs. 49.1%–83.5% for Spanish), suggesting that segmenting and classifying case affixes poses particular difficulty at smaller scales.

### 5.3 Surface Elements vs Semantic Components

Table 2 and Table 3 aggregate performance into surface elements versus semantic components for frontier and small models respectively.

Surface element accuracy exceeds semantic accuracy across nearly all conditions. The magnitude of this gap varies by language and model scale.

For frontier models, Spanish exhibits a narrow gap (89.6% vs. 86.6% for Claude), with surface and semantic accuracy nearly equivalent. Chinese shows substantial gaps: GPT-4o drops from 82.3% surface accuracy to 67.5% semantic accuracy, highlighting that correct parsing does not guarantee semantic comprehension. Basque demonstrates an intermediate pattern—surface element accuracy remains high, but semantic accuracy falls below Spanish despite comparable parsing success.

Small models show the same directional pattern but with lower overall performance and higher variance. The qualitative character of errors also differs between scales: frontier model errors cluster around specific components and languages, while small model errors are more uniformly distributed.

Table 1: Overall frontier (%). Figure/Ground: Levenshtein similarity; Predicate/Markers: avg. Levenshtein/type match; Dynamicity: accuracy; Configuration/Projective: F1-score.

	Spanish			Basque			Chinese		
	Claude	GPT-4o	Qwen	Claude	GPT-4o	Qwen	Claude	GPT-4o	Qwen
Figure	100.0	100.0	100.0	99.8	96.5	83.0	99.0	96.3	97.8
Ground	85.0	83.5	86.3	99.3	97.3	83.7	99.3	100.0	99.1
Predicate	99.2	99.0	94.3	87.1	90.3	82.0	83.4	61.9	68.2
Markers	74.2	76.4	70.5	72.1	49.4	17.5	76.4	70.8	61.3
Dynamicity	98.9	97.4	92.1	97.9	96.3	89.4	95.8	97.9	88.9
Configuration	73.7	66.9	49.4	64.6	68.9	42.8	57.0	50.1	50.2
Projective	87.2	92.2	60.6	79.7	82.8	40.6	74.0	54.6	65.0

Model	Spanish		Basque		Chinese	
	Surface	Sem	Surface	Sem	Surface	Sem
Claude-3.5-Sonnet	89.6	86.6	89.6	80.7	89.5	75.6
GPT-4o	89.7	85.5	83.4	82.7	82.3	67.5
Qwen2.5-VL-72B-Instruct	87.8	67.3	66.5	57.6	81.6	68.0

Table 2: Surface vs Semantic accuracy (%) for frontier models.

Model	Spanish		Basque		Chinese	
	Surface	Sem	Surface	Sem	Surface	Sem
Claude-3.5-Haiku	90.0	74.7	73.2	73.4	82.7	65.9
GPT-4o-Mini	77.4	53.2	56.0	45.4	74.1	47.6
Qwen-2.5-VL-7B-Instruct	80.6	46.4	64.8	35.4	55.8	31.8

Table 3: Surface vs Semantic accuracy (%) for small models.

## 5.4 Stasis Breakdown

### 5.4.1 Configuration

Frontier models reveal systematic patterns across semantic primitive combinations. Closed containment, superposition, and subposition achieve moderate-to-high accuracy, while open containment proves more challenging. The contact distinction and open versus closed geometry remain difficult across models. The [top] specification shows variable performance. (See Tables 7 and 8 in Appendix E for detailed performance on topological primitive combinations for frontier and small models, respectively).

Small models show moderate performance with systematic patterns emerging: [+contact] primitives consistently outperform [-contact] primitives, particularly for containment and superposition relations. However, overall accuracy remains lower than frontier models, with greater variance across primitive combinations and languages.

### 5.4.2 Projective

Frontier models achieve strong performance on absolute frame relations across Spanish and Basque,

with Claude and GPT approaching perfect accuracy. Intrinsic frame relations also perform well for Claude and GPT, though Qwen shows notable difficulty. Relative frame relations, however, show mixed results. The [relative][left] and [relative][right] distinctions achieve high accuracy across models, while [relative][back] and [relative][front] prove more challenging. (See Tables 9 and 10 in Appendix E, which show performance on projective relations by frame of reference for frontier and small models, respectively).

Small models show strong differentiation by frame type: relative frame relations achieve high accuracy for Claude and GPT, while intrinsic frame relations prove most challenging across all models. Absolute frame performance varies by model, with Claude and GPT maintaining moderate-to-high accuracy while Qwen shows weakness.

## 6 Discussion

**Parsing Does Not Entail Understanding.** The most consistent finding across our evaluation is the dissociation between surface element identification and semantic classification. Across both frontier and small models, accuracy on surface elements—figure, ground, spatial\_predicate, spatial\_markers—exceeds accuracy on semantic components—configuration and projective relations. This gap persists regardless of model scale, though its magnitude and character differ.

Chinese provides a clear demonstration of this pattern for frontier models. Claude and GPT achieve strong ground extraction (99.3%–100%), correctly identifying spatial nouns as part of the ground phrase. Yet configuration accuracy remains moderate (50.1%–57.0%), and projective accuracy varies substantially. Models identify *what* encodes spatial meaning without necessarily understanding *what spatial relation* is encoded.

Chinese `spatial_markers` illustrate a specific incoherence in frontier model behavior. Models include localizers as part of the ground phrase when prompted for ground extraction, yet classify these same elements as adposition rather than spatial noun when prompted for marker identification. This reveals that models apply task-specific heuristics rather than maintaining stable representations of grammatical status. An experiment with joint extraction of all components (Table 11, Appendix E) confirms this instability: localizers still appear in ground phrases but are not extracted as spatial nouns in markers, demonstrating the same contradiction. Overall performance degrades substantially—marker accuracy drops by 20-30 percentage points depending on model.

The qualitative character of errors also differs between scales. Frontier model errors cluster around specific components and languages, which points to systematic challenges in particular aspects of spatial understanding. In contrast, small model errors spread more uniformly across different areas—a pattern consistent with general capacity constraints. This distinction is informative: where frontier models succeed but small models fail may reflect capacity-dependent learning; where both struggle may indicate more fundamental challenges in spatial reasoning that persist despite increased scale.

**Dynamicity as Exception.** Not all semantic components prove equally difficult. Dynamicity classification remains robust across languages and most model scales, with frontier models achieving 88.9%–98.9% accuracy. Frontier models achieve high accuracy, and most small models maintain reasonable performance on this component despite struggling with configuration and projective relations. This resilience likely reflects the nature of the distinction: dynamicity involves a three-way categorical classification with strong correlations to verb semantics, making it more accessible than spatial primitives that require geometric reasoning.

**Morphological Marking.** Basque `spatial_markers` present a diagnostic case for morphological processing in spatial language understanding. Spatial relations in Basque are encoded through case suffixes that attach directly to nouns, requiring models to segment bound morphemes rather than identify free-standing tokens.

The contrast between model scales is pronounced. Small models show severe difficulties with Basque markers, with all three models achieving their lowest marker-type scores on this component. Frontier models improve substantially, though performance varies: Claude achieves comparable accuracy across languages (72.1%–76.4%), while GPT drops to 49.4% and Qwen to 17.5% for Basque markers despite strong performance on other surface elements.

This pattern suggests that segmenting and identifying case affixes as spatial encoding poses particular challenges, especially at smaller scales. This pattern could extend to other languages with spatial case marking, such as Finnish or Hungarian. Evaluation limited to prepositional languages would miss this systematic limitation entirely.

**Configuration.** Primitives reveal asymmetric performance patterns. Closed containment, superposition, and subposition achieve moderate-to-high accuracy in frontier models, while open containment proves more challenging. The contact distinction remains difficult: [-contact] primitives underperform [+contact] by 15–30 percentage points across configuration types. The [open] versus [closed] distinction shows instability, with [open] containment achieving lower accuracy. The [top] specification shows variable performance, with accuracy dependent on lexical cues in the input.

To investigate whether visual information could address the semantic understanding gaps observed in text-only evaluation, we conducted a supplementary experiment providing models with both sentences and their corresponding images. Visual grounding produces mixed effects on configuration accuracy, with overall performance changes varying by model scale and language (Table 13, Appendix E) and no consistent improvement pattern. However, examining specific primitive combinations reveals selective gains: frontier models show systematic improvements for [+contact][containment][open] relations across all languages, with increases ranging from +5.6 to +44.4 percentage points (Table 14, Appendix E). This pattern does not extend to small models, where improvements are inconsistent (Table 15, Appendix E). Conversely, many other primitive combinations—particularly involving [-contact] and superposition—show substantial decreases when images are added (e.g., Claude-3.5-Haiku drops 18.7 percentage points on Spanish

597 configuration).

598 Small models exhibit systematic rather than ran- 648  
599 dom errors: [+contact] primitives consistently out- 649  
600 perform [-contact] primitives, particularly for con- 650  
601 tainment and superposition relations. However, 651  
602 overall accuracy remains lower than frontier mod- 652  
603 els. This suggests that even smaller architectures 653  
604 can learn basic contact distinctions, though they 654  
605 struggle with the full complexity of primitive com- 655  
606 binations. 656

607 **Projective.** Projective relations present distinct 657  
608 evaluation challenges depending on frame type. 658  
609 Absolute frame relations—cardinal directions— 659  
610 rely on lexically explicit nouns (north, south, east, 660  
611 west), making them more tractable for text-only 661  
612 evaluation. Frontier models achieve consistently 662  
613 high accuracy on absolute frames in Spanish and 663  
614 Basque for Claude and GPT, though Qwen shows 664  
615 notable weakness even on these lexically explicit 665  
616 relations. 666

617 Intrinsic and relative frame relations pose a 667  
618 more fundamental challenge: without visual con- 668  
619 text, these distinctions can be genuinely ambigu- 669  
620 ous in text alone. A sentence like “the fox is in front 670  
621 of the house” could describe either an intrinsic re- 671  
622 lation (relative to the house’s inherent front) or a re- 672  
623 lative relation (from the viewer’s perspective). Given 673  
624 this inherent ambiguity, we avoid strong claims 674  
625 about the relative difficulty of specific directional 675  
626 values or frame types, as observed patterns may re- 676  
627 flect textual cues rather than spatial reasoning per 677  
628 se. 678

629 This motivated our visual grounding experiment. 679  
630 However, adding images does not produce con- 680  
631 sistent improvements. Small models show highly 681  
632 variable changes ranging from -8.3 to +58.3 per- 682  
633 centage points across cardinal directions (Table 17, 683  
634 Appendix E). Frontier models show minimal or 684  
635 negative changes for absolute frames, with mixed 685  
636 results across other frame types (Table 16, Ap- 686  
637 pendix E). Recent work demonstrates similar chal- 687  
638 lenges in visual grounding for spatial distinctions 688  
639 (Tong et al., 2024), indicating that projective re- 689  
640 lations merit dedicated future investigation. 690

## 641 7 Conclusion

642 We introduced a novel framework for spatial lan- 691  
643 guage understanding and a multilingual bench- 692  
644 mark that addresses systematic gaps in existing 693  
645 evaluation approaches. By decomposing spatial 694  
646 relations into surface elements and semantic com-

647 ponents, and by recognizing spatial markers be- 648  
649 yond prepositions, our framework enables evalua- 650  
651 tion across typologically diverse languages. 652

653 Our evaluation of frontier and small-scale LLMs 654  
655 on Spanish, Basque, and Chinese reveals three key 656  
657 findings. First, surface element identification ac- 658  
659 curacy consistently exceeds semantic classification 660  
661 accuracy, demonstrating that successful surface 662  
663 parsing does not entail spatial understanding. Sec- 664  
665 ond, morphological spatial marking poses persis- 666  
667 tent challenges: Basque case suffixes prove particu- 668  
669 larly difficult across model scales, with small mod- 670  
671 els showing severe difficulties and frontier models 672  
673 demonstrating variable performance depending on 674  
675 architecture. Third, primitive decomposition us- 676  
677 ing F1-score metrics exposes systematic patterns— 678  
679 contact asymmetries, containment geometry insta- 680  
681 bility, frame-of-reference distinctions—that would 681  
682 be invisible under atomic labeling schemes. 682

683 These findings point to several research direc- 684  
685 tions for advancing spatial language understanding 685  
686 in LLMs. First, expanding typological coverage 686  
687 to languages with diverse morphological marking 687  
688 systems is essential—current evaluation systemat- 688  
689 ically excludes case-marking strategies, limiting 689  
690 our understanding of whether models comprehend 690  
691 spatial relations or succeed primarily on preposi- 691  
692 tional patterns. Second, component-wise evalua- 692  
693 tion frameworks provide diagnostic precision that 693  
694 atomic labels cannot, enabling separation of pars- 694  
695 ing failures from semantic failures and targeted 695  
696 analysis of systematic difficulties. Third, multi- 696  
697 modal frameworks designed specifically for spatial 697  
698 reasoning merit dedicated investigation. Finally, 698  
699 projective relations require specialized methodolo- 699  
700 gies for frame-of-reference distinctions, particu- 700  
701 larly for intrinsic and relative frames where textual 701  
702 ambiguity is inherent. 702

703 Advancing spatial language understanding in 703  
704 LLMs requires expanding evaluation beyond 704  
705 prepositional languages to determine whether mod- 705  
706 els genuinely comprehend spatial meaning or suc- 706  
707 ceed primarily on surface patterns. Our benchmark 707  
708 provides initial evidence for this necessity, and our 708  
709 component-wise framework offers a methodology 709  
710 that can be extended to additional languages and 710  
711 spatial phenomena in future research. 711

## 712 8 Limitations

713 Our language sample, while typologically diverse, 713  
714 covers only three languages. The patterns we ob- 714  
715 716

697	serve—particularly the difficulty with morphological marking—should be tested on additional languages with spatial case systems (e.g., Finnish, Hungarian) to confirm their generality.	
698		
699		
700		
701	Our multimodal evaluation provides initial evidence that visual grounding does not resolve semantic challenges, but was limited in scope. The experiment used canonical static images with text-only prompts, without systematic manipulation of visual properties, attention mechanisms, or multimodal integration strategies. A comprehensive understanding of why visual input fails to help would require controlled studies examining image properties, model attention patterns, and visual-linguistic integration mechanisms—work beyond the scope of this benchmark-focused study.	
702		
703		
704		
705		
706		
707		
708		
709		
710		
711		
712		
713	Our typology of spatial markers—adpositions, affixes, and spatial nouns—does not exhaust the crosslinguistic inventory. Less common strategies such as tonal marking could encode spatial meaning in some languages, but fall outside our current annotation scheme.	
714		
715		
716		
717		
718		
719	The primitive decomposition, while more fine-grained than existing schemes, does not exhaust spatial semantics. The configuration component could be extended to capture relations like “between” or “among” and configurations involving encirclement or multiple grounds, which are not reducible to our current feature set. Similarly, the projective component could incorporate finer angular distinctions and language-specific absolute systems beyond cardinal directions.	
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		
	<b>9 Ethical Considerations</b>	
730	This research was approved by the University of Melbourne Human Research Ethics Committee. Native speaker participants were recruited through personal networks and volunteered without financial compensation for this brief linguistic elicitation task. All participants provided informed consent and were informed that their anonymized sentence productions would be used for NLP research and made publicly available as part of a benchmark dataset. The dataset contains only anonymized linguistic annotations with no personally identifiable information.	
731		
732		
733		
734		
735		
736		
737		
738		
739		
740		
741		
742		
	<b>References</b>	
743	Malihe Alikhani, Valts Blukis, Parisa Kordjamshidi, Aishwarya Padmakumar, and Hao Tan, editors.	
744		
	2021. <i>Proceedings of Second International Combined Workshop on Spatial Language Understanding and Grounded Communication for Robotics</i> . Association for Computational Linguistics, Online.	745 746 747 748
	Anthropic. 2024. <a href="#">Model Card Addendum: Claude 3.5 Haiku and Upgraded Claude 3.5 Sonnet</a> . Technical Report.	749 750 751
	Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 others. 2025. <a href="#">Qwen2.5-VL Technical Report</a> . <i>Preprint</i> , arXiv:2502.13923.	752 753 754 755 756 757 758
	Barend Beekhuizen. 2025. <a href="#">Spatial relation marking across languages: Extraction, evaluation, analysis</a> . In <i>Proceedings of the 29th Conference on Computational Natural Language Learning</i> , pages 571–585, Vienna, Austria. Association for Computational Linguistics.	759 760 761 762 763 764
	Archana Bhatia, Yonatan Bisk, Parisa Kordjamshidi, and Jesse Thomason, editors. 2019. <i>Proceedings of the Combined Workshop on Spatial Language Understanding (SpLU) and Grounded Communication for Robotics (RoboNLP)</i> . Association for Computational Linguistics, Minneapolis, Minnesota.	765 766 767 768 769 770
	Melissa Bowerman and Erik Pederson. 1992. Topological relations picture series. In <i>Space Stimuli Kit 1.2</i> , page 51. Max Planck Institute for Psycholinguistics, Nijmegen.	771 772 773 774
	Denis Creissels. 2008. <a href="#">Spatial Cases</a> . In Andrej L. Malchukov and Andrew Spencer, editors, <i>The Oxford Handbook of Case</i> , page 0. Oxford University Press.	775 776 777 778
	Redouane Djamouri, Waltraud Paul, and John Whitman. 2013. <a href="#">Postpositions vs Prepositions in Mandarin Chinese: The Articulation of Disharmony</a> . In Theresa Biberauer and Michelle Sheehan, editors, <i>Theoretical Approaches to Disharmonic Word Order</i> , page 0. Oxford University Press.	779 780 781 782 783 784
	Ricardo Etxepare. 2013. <a href="#">Basque Primary Adpositions from a Clausal Perspective</a> . <i>Catalan journal of linguistics</i> , 12:41–82.	785 786 787
	Ricardo Etxepare and Bernard Oyharçabal. 2012. <a href="#">Datives and Adpositions in Northeastern Basque</a> . In Beatriz Fernandez and Ricardo Etxepare, editors, <i>Variation in Datives: A Microcomparative Perspective</i> , page 0. Oxford University Press.	788 789 790 791 792
	Martin Haspelmath. 2019. <a href="#">Differential place marking and differential object marking</a> . <i>STUF - Language Typology and Universals</i> , 72(3):313–334.	793 794 795
	Hualde. 2011. <i>A Grammar of Basque</i> . De Gruyter Mouton.	796 797

798	C.-T. James Huang, Y.-H. Audrey Li, and Yafei Li.	Thomas Stolz, Sander Lestrade, and Christel Stolz.	855
799	2009. <i>The Syntax of Chinese</i> . Cambridge Syntax	2014. <i>The Cross Linguistics of Zero-Marking of Spa-</i>	856
800	Guides. Cambridge University Press, Cambridge.	<i>tial Relations</i> , 1st ed. edition. Number 15 in <i>Studia</i>	857
801	Parisa Kordjamshidi, Archana Bhatia, Malihe Alikhani,	<i>Typologica</i> . De Gruyter Mouton, Berlin ;.	858
802	Jason Baldrige, Mohit Bansal, and Marie-Francine	Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma,	859
803	Moens, editors. 2020. <i>Proceedings of the Third Inter-</i>	Yann LeCun, and Saining Xie. 2024. <i>Eyes Wide</i>	860
804	<i>national Workshop on Spatial Language Understand-</i>	<i>Shut? Exploring the Visual Shortcomings of Mul-</i>	861
805	<i>ing</i> . Association for Computational Linguistics, On-	<i>timodal LLMs</i> . <i>Preprint</i> , arXiv:2401.06209.	862
806	line.	Morgan Ulinski, Bob Coyne, and Julia Hirschberg.	863
807	Parisa Kordjamshidi, Archana Bhatia, James Puste-	2019. <i>SpatialNet: A Declarative Resource for</i>	864
808	jovskiy, and Marie-Francine Moens, editors. 2018.	<i>Spatial Relations</i> . In <i>Proceedings of the Com-</i>	865
809	<i>Proceedings of the First International Workshop on</i>	<i>binated Workshop on Spatial Language Understanding</i>	866
810	<i>Spatial Language Understanding</i> . Association for	<i>(SpLU) and Grounded Communication for Robotics</i>	867
811	Computational Linguistics, New Orleans.	<i>(RoboNLP)</i> , pages 61–70, Minneapolis, Minnesota.	868
812	Parisa Kordjamshidi, Xin Eric Wang, Yue Zhang,	Association for Computational Linguistics.	869
813	Ziqiao Ma, and Mert Inan, editors. 2024. <i>Proceed-</i>	I-hao Woo. 2021. <i>On preverbal zai in Mandarin Chi-</i>	870
814	<i>ings of the 4th Workshop on Spatial Language Under-</i>	<i>nese: Its progressive and prepositional functions.</i>	871
815	<i>standing and Grounded Communication for Robotics</i>	<i>Linguistics</i> , 59(3):513–539.	872
816	<i>(SpLU-RoboNLP 2024)</i> . Association for Computa-	<b>A Full Annotation Examples</b>	873
817	tional Linguistics, Bangkok, Thailand.	<b>A.1 Configuration Example: Bird on Ground</b>	874
818	George Lakoff. 1987. <i>Women, Fire, and Dangerous</i>	This section provides the complete annotations for	875
819	<i>Things: What Categories Reveal about the Mind</i> .	Spanish and Chinese descriptions of the spatial	876
820	Women, Fire, and Dangerous Things: What Cate-	scene shown in Figure 1 (the bird on the ground).	877
821	gories Reveal about the Mind. University of Chicago	The Basque annotation appears in the main text	878
822	Press, Chicago, IL, US.	(Figure 1).	879
823	Stephen Levinson and Sergio Meira. 2003. ‘Natural	{“sentence”: “El pajarito esta en el suelo.”,	
824	concepts’ in the spatial topological domain - adpo-	“language”: “es”,	
825	sitional meanings in crosslinguistic perspective: An	“figure”: “El pajarito”,	
826	exercise in semantic typology. <i>Language</i> .	“ground”: “el suelo”,	
827	Stephen C. Levinson and David P. Wilkins, editors.	“spatial_predicate”:	
828	2006. <i>Grammars of Space: Explorations in Cog-</i>	{“text”: “esta”, “type”: “copula”},	
829	<i>nitive Diversity</i> . Language Culture and Cognition.	“spatial_markers”: [	
830	Cambridge University Press, Cambridge.	{“text”: “en”, “type”: “adposition”}],	
831	Jingxia Lin. 2011. <i>A figure’s final location must be</i>	“dynamicity”: “static”,	
832	<i>identifiable: Localizer distribution in Chinese mo-</i>	“stasis”:	
833	<i>tion expressions</i> . <i>Annual Meeting of the Berkeley</i>	{“configuration”:	
834	<i>Linguistics Society</i> , pages 242–256.	“[+contact][superposition]”,	
835	Jingping Liu, Ziyang Liu, Zhedong Cen, Yan Zhou, Yi-	“projective”: {“type”: “”, “value”:	
836	nan Zou, Weiyang Zhang, Haiyun Jiang, and Tong	“”}}}	
837	Ruan. 2025. <i>Can Multimodal Large Language Mod-</i>		
838	<i>els Understand Spatial Relations?</i> In <i>Proceed-</i>		
839	<i>ings of the 63rd Annual Meeting of the Association</i>		
840	<i>for Computational Linguistics (Volume 1: Long Pa-</i>		
841	<i>pers)</i> , pages 620–632, Vienna, Austria. Association		
842	for Computational Linguistics.		
843	Michał Olek and Maciej Piasecki. 2024. <i>Three-Stage</i>	Figure 2: Spanish annotation for “El pájaro está en el	
844	<i>Extraction of Spatial Relationships Using Mark-</i>	suelo” (The bird is on the ground). Spanish encodes	
845	<i>ers</i> . In <i>Advances in Computational Collective In-</i>	the spatial relation using the copula <i>está</i> and a single	
846	<i>telligence</i> , pages 159–172, Cham. Springer Nature	adposition <i>en</i> .	
847	Switzerland.	<b>A.1.1 Chinese</b>	880
848	OpenAI, Aaron Hurst, Adam Lerer, Adam P. Goucher,	<b>Cross-linguistic observations.</b> The three lan-	881
849	Adam Perelman, Aditya Ramesh, Aidan Clark, A. J.	guages encode the same topological relation	882
850	Ostrow, Akila Welihinda, Alan Hayes, Alec Radford,	([+contact][superposition]) using different surface	883
851	Aleksander Mądry, Alex Baker-Whitcomb, Alex	strategies:	884
852	Beutel, Alex Borzunov, Alex Carney, Alex Chow,		
853	Alex Kirillov, Alex Nichol, and 400 others. 2024.		
854	<i>GPT-4o System Card</i> . <i>Preprint</i> , arXiv:2410.21276.		

```
{ "sentence": "Niao zai di shang.",
  "translation": "The bird is on the ground",
  "language": "zh",
  "figure": "niao",
  "ground": "di shang",
  "spatial_predicate":
    { "text": "", "type": "" },
  "spatial_markers": [
    { "text": "zai", "type": "adposition" },
    { "text": "shang", "type":
      "spatial_noun" } ],
  "dynamicity": "static",
  "stasis":
    { "configuration":
      "[+contact][superposition]",
      "projective": { "type": "", "value":
        "" } } }
```

Figure 3: Chinese annotation for “鸟在地上” (Niǎo zài dì shàng / The bird is on the ground). Chinese uses no explicit `spatial_predicate`, but employs two `spatial_markers`: the preposition *zài* and the spatial noun *shàng* (‘top/on’).

- **Spanish** relies on a copula (*está*) and a single general-purpose preposition (*en*) that does not specify the precise topological relationship.
- **Basque** uses a lexical verb (*dago*) and two markers: a spatial noun (*gain* ‘top/surface’) that specifies the topological relationship, plus a locative case affix (*-an*) marking the ground.
- **Chinese** has no explicit `spatial_predicate` but employs two markers: a preposition (*zài*) marking general spatial location, and a spatial noun (*shàng* ‘top/on’) specifying the topological relationship.

This variation demonstrates that evaluating only prepositional markers would miss crucial spatial encoding strategies. The semantic annotation (configuration and dynamicity), however, remains consistent across languages, showing that our primitive-based approach captures meaning independently of surface form.

## A.2 Projective Examples: Frame-of-Reference Relations

This section provides examples of projective relations requiring different frames of reference. Un-

like the bird example, these relations involve directional specification.

### A.2.1 Intrinsic Frame: Fox in Front of House

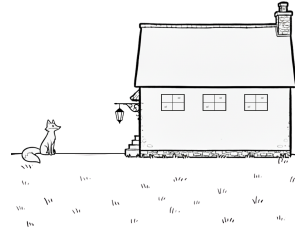


Figure 4: Visual stimulus for intrinsic frame relation: the fox’s position is described relative to the house’s inherent front facet.

```
{ "sentence": "Azeria etxearen aurrean dago.",
  "translation": "The fox is in front of the
    house",
  "language": "eu",
  "figure": "Azeria",
  "ground": "etxearen aurrean",
  "spatial_predicate":
    { "text": "dago", "type": "copula" },
  "spatial_markers": [
    { "text": "aurre", "type":
      "spatial_noun" },
    { "text": "-an", "type": "affix" } ],
  "dynamicity": "static",
  "stasis":
    { "configuration": "",
      "projective": { "type": "intrinsic",
        "value": "front" } } }
```

Figure 5: Basque annotation for “Azeria etxearen aurrean dago” (The fox is in front of the house). This is an intrinsic frame relation: the spatial description relies on identifying the house’s inherent front facet. The spatial noun *aurre* (‘front’) combined with the locative case *-an* encodes this directional relationship.

**Why this is projective.** Unlike the bird-on-ground example, which involves physical contact and can be described purely through topological primitives, the fox-house relation requires specifying a direction. The house has an inherent orientation—a canonical front side (typically where the door is located)—and the fox’s position is described relative to this intrinsic property of the ground object. The observer’s position is irrele-

921	vant; the relation holds regardless of where one	<b>B Primitive Values and Examples</b>	967
922	views the scene from.		
923	<b>A.2.2 Relative Frame Example</b>	<b>C Annotation Guidelines</b>	968
924	Relative frame relations depend on the observer’s	<b>Figure</b>	969
925	viewpoint. For example, “The ball is to the left	The noun phrase denoting the entity whose location	970
926	of the tree” describes the ball’s position using co-	or motion is being described. Extract the complete	971
927	ordinates derived from the observer’s bodily axes.	NP in its base form: include all determiners,	972
928	If the observer moves to the opposite side of the	modifiers, and complements that belong to the	973
929	tree, the same physical configuration might be de-	phrase itself, but exclude independent adpositions.	974
930	scribed as “to the right of the tree.”	For contractions (e.g., ‘al’ = ‘a’ + ‘el’), extract only	975
931	Our dataset includes relative frame examples	the NP component (‘el’).	976
932	with all four directional values: left, right, front,	<b>Ground</b>	977
933	and back.	The reference object, site, or region relative to	978
934	<b>A.2.3 Absolute Frame Example</b>	which the Figure’s location or motion is specified.	979
935	Absolute frame relations use fixed environmental	Extract the complete noun phrase that serves as	980
936	bearings. For example, “The village is south of the	the landmark, including all determiners, modifiers,	981
937	mountain” employs cardinal directions that remain	and complements. This includes any lexical el-	982
938	constant regardless of observer position or object	ements with spatial meaning (like ‘top’, ‘front’,	983
939	orientation.	‘right’, ‘side’, ‘cima’, ‘lado’, ‘上面’)	984
940	Our dataset includes absolute frame examples	when they function as the head noun of the phrase. The	985
941	with all four cardinal values: north, south, east, and	entire noun phrase built around such spatial nouns	986
942	west.	(e.g., ‘la cima del edificio’, ‘the top of the moun-	987
943	<b>A.3 Dynamicity Variation</b>	tain’, ‘大楼的上面’) must be extracted as the com-	988
944	The same visual configuration can be described	plete ground.	989
945	with different dynamicity values depending on	<b>Spatial Predicate</b>	990
946	whether motion is involved and, if so, in which di-	The predicate that expresses the static relation or	991
947	rection.	motion event linking Figure and Ground. Extract	992
948	<b>Static:</b> “The bird is on the ground” describes a	the complete predicate as it appears in the sen-	993
949	stable spatial relation without motion.	tence, including all its grammatical elements and	994
950	<b>Goal:</b> “The bird lands on the ground” or	any complements that are part of the verbal unit it-	995
951	“The bird flies onto the ground” describes mo-	self. Do NOT include spatial markers that are con-	996
952	tion toward the ground as destination. The	nected to the Ground. Includes both ‘text’ (the sur-	997
953	configuration value ([+contact][superposition])	face form) and ‘type’ (grammatical category). If	998
954	describes the <b>final</b> spatial state after the motion is	the construction does not have a spatial predicate,	999
955	complete.	the entire field should be an empty string “”.	1000
956	<b>Source:</b> “The bird takes off from the ground” or	• verb: Lexical verbs that denote actions, states,	1001
957	“The bird flies from the ground” describes motion	or processes (e.g., lies, sits, runs, goes, re-	1002
958	away from the ground as departure point. The	mains).	1003
959	configuration value describes the <b>initial</b> spatial	• copula: Grammatical elements that link the	1004
960	state before the motion begins.	figure to the ground, typically expressing at-	1005
961	This systematic variation allows us to evaluate	tribution rather than action (e.g., is, was).	1006
962	whether models understand that dynamicity and	<b>Spatial Markers</b>	1007
963	configuration are independent but interacting di-	A list of elements that encode the spatial relation-	1008
964	mensions: the same topological relationship can	ship between a figure (the entity whose location or	1009
965	occur in static contexts or as the source/goal of mo-	motion is being described) and a ground (the refer-	1010
966	tion events.	ence object). Only include elements that are gram-	1011
		matically linked to a ground NP. Each marker is	1012

Primitives	Example
<i>Configuration (topological)</i>	
[+contact][containment][open]	The cat is in the box
[-contact][containment][open]	The fish is in the fish tank
[+contact][containment][closed]	The cat is in the house
[-contact][containment][closed]	The fly is flying inside the house
[+contact][attachment]	The magnet is stuck to the refrigerator
[+contact][superposition]	The bird is on the ground
[-contact][superposition]	The bird is flying above the house
[+contact][superposition][top]	The cat is at the top of the mountain
[-contact][superposition][top]	The helicopter is hovering above the roof of the building
[+contact][subposition]	The cat is under the blanket
[-contact][subposition]	The cat is under the table
<i>Projective (angular)</i>	
<i>Absolute</i>	
[absolute][east]	The car is to the east of the church
[absolute][north]	The car is to the north of the church
[absolute][south]	The car is to the south of the church
[absolute][west]	The car is to the west of the church
<i>Intrinsic</i>	
[intrinsic][back]	The fox is behind the house
[intrinsic][front]	The fox is in front of the house
<i>Relative</i>	
[relative][back]	The cat is behind the box
[relative][front]	The cat is in front of the box
[relative][left]	The cat is to the left of the box
[relative][right]	The cat is to the right of the box

Table 4: Primitive components for topological (configuration) and angular (projective) spatial relations with representative examples.

recorded separately with ‘text’ and ‘type’. Markers must be listed in their order of appearance in the sentence. If the construction has no spatial markers, the entire field should be an empty list [].

- **adposition:** Prepositions, postpositions, or fixed multi-word constructions that have grammaticalized as inseparable units (e.g., ‘on’, ‘in’, ‘at’, ‘on top of’, ‘encima de’, ‘在’). When contractions or fused forms incorporate elements that belong to the ground NP (e.g., Spanish ‘del’ = ‘de’ + ‘el’), report only the spatial marker portion (e.g., ‘de’), while the ground NP element remains part of the ground. For multi-word adpositions, determine if the expression is grammaticalized as a fixed unit: if it does not accept internal modification (e.g., ‘on top of’, ‘encima de’), mark it as a single adposition; if the spatial noun retains nominal behavior and accepts modification (e.g., ‘on the top of’, ‘en la cima de’), separate the adposition from the spatial noun.
- **affix:** Case markers or other morphological affixes that encode spatial meaning (e.g., Latin ablative ‘-o’, accusative ‘-um’, Basque inessive ‘-n’).

- **spatial\_noun:** Lexical nouns with spatial meaning that function as the head of the ground NP. Extract only the stem or head noun itself (e.g., ‘top’, ‘cima’, ‘上面’, ‘gain’), excluding determiners, modifiers, and inflectional morphology. Examples: ‘cima’ (from ‘la cima del edificio’), ‘gain’ (from ‘mahaiaren gainean’).

### Dynamicity

The temporal aspect of the spatial relation, indicating whether the relation is static or involves motion.

- **static:** The Figure’s spatial relation to the Ground is stable; no translational motion occurs.
- **source:** The Figure moves away from the Ground, which serves as the point of departure of the motion.
- **goal:** The Figure moves toward the Ground, which serves as the final location of the motion.

### Stasis

The spatial relationship between Figure and Ground. ALWAYS include both ‘configuration’

and ‘projective’ keys.

`dynamic_rule`: For motion: if Figure moves AWAY FROM Ground, describe INITIAL state; if TOWARD Ground, describe FINAL state.

`configuration`: Topological relationship as concatenated primitives in brackets (e.g., ‘[+contact][containment][open]’). Contact ([+contact] or [-contact]) must ALWAYS be first. Use empty string if projective is used.

- `contact`: [+contact] or [-contact] - physical contact between Figure and Ground. ALWAYS first.
- `containment`: [containment][open] (partial enclosure) or [containment][closed] (full enclosure)
- `attachment`: [attachment] - mechanically fastened
- `superposition`: [superposition] - Figure above Ground. Add [top] if on uppermost surface.
- `subposition`: [subposition] - Figure below Ground

`projective`: Angular relationships via frames of reference. If configuration is used, include this key with empty strings.

`type`: ‘absolute’ (fixed bearings), ‘relative’ (viewer-dependent), or ‘intrinsic’ (object’s inherent parts)

`value`: north | south | east | west | left | right | front | back

## D Language-Specific Considerations

Each language presented distinct analytical challenges.

**Spanish.** Contracted forms combining prepositions with articles (e.g., *del* from *de* + *el*) required splitting to isolate the `spatial_marker`. Multiword adpositions such as *encima de* (‘on top of’) were treated as single markers when grammaticalized as fixed units.

**Basque.** Basque employs an agglutinative system where spatial relations are encoded through case suffixes—inessive *-n*, allative *-ra*, ablative *-tik*—combined with spatial nouns (Hualde, 2011). Spatial nouns like *gain* (‘top’) frequently combine with these case markers and are treated as nouns following Etxepare (2013). Directional postpositions such as *behera* (‘down’) were analyzed as

adpositions following Etxepare and Oyharçabal (2012). Genitive cases functioning as linkers between reference objects and spatial nouns were not annotated as `spatial_markers`, as they establish possessive relationships rather than contributing spatial information per se.

**Chinese.** Chinese presents particular challenges due to the grammaticalization status of elements like *lǐ* (‘inside’), *shàng* (‘top/on’), and *xià* (‘under’). These forms—termed ‘localizers’ in the literature—derive historically from nouns, but their synchronic category remains debated, having been analyzed variously as NP enclitics, locative particles, postpositions, or as a subclass of nouns (Huang et al., 2009; Lin, 2011; Djamouri et al., 2013). We annotate these elements as `spatial_noun` markers, consistent with their nominal etymology and syntactic behavior. Regarding *zài*, we follow Woo (2021) in treating it consistently as a spatiotemporal preposition that introduces the Ground. Separately, directional complements appearing in verb-directional constructions (e.g., *chūlái* ‘exit-come’ in *zuān le chūlái* ‘burrowed out’) were annotated as part of the `spatial_predicate`, as these elements modify the motion event rather than marking the ground.

Table 5: Overall frontier (%). Figure/Ground: Levenshtein similarity; Predicate/Markers: avg. Levenshtein/type match; Dynamicity: accuracy; Configuration/Projective: F1-score.

	Spanish			Basque			Chinese		
	Claude	GPT-4o	Qwen	Claude	GPT-4o	Qwen	Claude	GPT-4o	Qwen
Figure	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	99.8 ± 0.4	96.5 ± 0.4	83.0 ± 1.1	99.0 ± 0.8	96.3 ± 0.8	97.8 ± 0.6
Ground	85.0 ± 0.5	83.5 ± 0.5	86.3 ± 0.0	99.3 ± 0.4	97.3 ± 0.1	83.7 ± 0.3	99.3 ± 0.6	100.0 ± 0.0	99.1 ± 0.3
Predicate	99.2 ± 0.0	99.0 ± 0.2	94.3 ± 0.7	87.1 ± 0.8	90.3 ± 2.6	82.0 ± 0.5	83.4 ± 3.6	61.9 ± 1.5	68.2 ± 0.3
Markers	74.2 ± 0.5	76.4 ± 1.6	70.5 ± 0.6	72.1 ± 1.6	49.4 ± 0.9	17.5 ± 0.5	76.4 ± 1.3	70.8 ± 4.2	61.3 ± 0.7
Dynamicity	98.9 ± 0.9	97.4 ± 0.9	92.1 ± 0.0	97.9 ± 0.9	96.3 ± 0.9	89.4 ± 0.9	95.8 ± 0.9	97.9 ± 0.9	88.9 ± 0.0
Configuration	73.7 ± 2.2	66.9 ± 2.4	49.4 ± 1.9	64.6 ± 1.5	68.9 ± 2.0	42.8 ± 0.9	57.0 ± 1.9	50.1 ± 3.0	50.2 ± 0.3
Projective	87.2 ± 1.0	92.2 ± 1.0	60.6 ± 2.5	79.7 ± 0.0	82.8 ± 3.7	40.6 ± 2.5	74.0 ± 0.9	54.6 ± 1.6	65.0 ± 0.0

Table 6: Overall small (%). Figure/Ground: Levenshtein similarity; Predicate/Markers: avg. Levenshtein/type match; Dynamicity: accuracy; Configuration/Projective: F1-score. (with standard deviation)

	Spanish			Basque			Chinese		
	Haiku	GPT-mini	Qwen-7B	Haiku	GPT-mini	Qwen-7B	Haiku	GPT-mini	Qwen-7B
Figure	97.3 ± 0.0	96.4 ± 0.0	94.8 ± 0.5	86.0 ± 0.5	40.8 ± 2.1	85.2 ± 1.2	88.2 ± 0.0	65.8 ± 1.0	91.0 ± 1.5
Ground	86.5 ± 0.1	89.5 ± 1.2	88.3 ± 0.9	93.0 ± 0.3	78.8 ± 1.0	67.2 ± 2.6	99.2 ± 0.5	96.4 ± 0.4	43.6 ± 0.5
Predicate	92.7 ± 0.2	74.7 ± 0.4	75.8 ± 0.5	84.1 ± 0.0	89.2 ± 1.2	88.3 ± 2.4	69.8 ± 0.0	66.9 ± 0.8	59.2 ± 0.5
Markers	83.5 ± 0.7	49.1 ± 0.5	63.4 ± 1.4	29.6 ± 0.6	15.3 ± 0.7	18.5 ± 2.5	73.7 ± 1.3	67.1 ± 0.7	29.3 ± 2.1
Dynamicity	95.2 ± 0.0	89.4 ± 2.4	75.7 ± 0.9	94.7 ± 0.9	77.8 ± 2.7	58.2 ± 0.9	93.7 ± 0.0	76.7 ± 0.9	43.9 ± 0.9
Configuration	50.3 ± 1.6	30.4 ± 0.6	30.6 ± 2.3	42.6 ± 2.1	26.4 ± 1.3	21.9 ± 0.5	31.8 ± 0.3	24.9 ± 0.2	22.6 ± 3.8
Projective	78.7 ± 1.3	39.9 ± 0.9	32.8 ± 3.2	83.0 ± 1.3	31.9 ± 1.1	26.0 ± 2.7	72.3 ± 2.8	41.0 ± 1.2	28.8 ± 5.1

Table 7: Configuration Breakdown by Primitive Combination frontier (with standard deviation)

	Claude-3.5-Sonnet			GPT-4o			Qwen2.5-72B		
	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH
[+contact][containment][open]	70.4 ± 12.8	55.6 ± 0.0	77.8 ± 0.0	63.0 ± 6.4	54.1 ± 4.6	77.8 ± 0.0	88.9 ± 0.0	96.3 ± 6.4	88.9 ± 0.0
[-contact][containment][open]	37.0 ± 6.4	44.4 ± 0.0	33.3 ± 0.0	33.3 ± 0.0	44.4 ± 0.0	33.3 ± 0.0	44.4 ± 0.0	55.6 ± 0.0	44.4 ± 0.0
[+contact][containment][closed]	77.8 ± 0.0	81.5 ± 6.4	80.0 ± 0.0	100.0 ± 0.0	94.4 ± 9.6	68.5 ± 17.0	77.8 ± 0.0	74.1 ± 3.2	77.8 ± 0.0
[-contact][containment][closed]	77.8 ± 0.0	51.9 ± 6.4	55.6 ± 0.0	66.7 ± 0.0	70.4 ± 6.4	66.7 ± 0.0	37.0 ± 6.4	44.4 ± 0.0	44.4 ± 0.0
[+contact][attachment]	100.0 ± 0.0	80.0 ± 0.0	60.0 ± 0.0	77.8 ± 9.6	80.0 ± 0.0	40.0 ± 0.0	83.0 ± 5.1	45.9 ± 5.1	48.9 ± 1.9
[+contact][superposition]	72.2 ± 9.6	83.3 ± 0.0	95.6 ± 3.8	74.1 ± 12.8	86.7 ± 3.8	63.0 ± 12.8	53.7 ± 3.2	74.4 ± 9.6	96.3 ± 6.4
[-contact][superposition]	80.0 ± 0.0	80.0 ± 0.0	73.3 ± 11.5	43.3 ± 0.0	44.4 ± 1.9	41.1 ± 1.9	55.6 ± 0.0	37.8 ± 7.7	55.6 ± 0.0
[+contact][superposition][top]	95.6 ± 3.8	53.3 ± 11.5	46.7 ± 0.0	88.9 ± 7.7	100.0 ± 0.0	47.8 ± 13.5	80.0 ± 0.0	56.3 ± 12.8	66.7 ± 0.0
[-contact][superposition][top]	74.1 ± 2.6	40.0 ± 7.7	44.4 ± 0.0	66.7 ± 0.0	35.6 ± 0.0	44.4 ± 0.0	26.7 ± 0.0	26.7 ± 0.0	57.8 ± 0.0
[+contact][subposition]	83.3 ± 0.0	83.3 ± 0.0	80.0 ± 0.0	48.1 ± 30.6	63.0 ± 6.4	36.7 ± 17.3	100.0 ± 0.0	63.3 ± 0.0	100.0 ± 0.0
[-contact][subposition]	66.7 ± 0.0	100.0 ± 0.0	50.0 ± 0.0	74.1 ± 6.4	85.2 ± 3.2	70.0 ± 17.3	44.4 ± 9.6	38.9 ± 9.6	50.0 ± 0.0

Table 8: Configuration Breakdown by Primitive Combination small (with standard deviation)

	Claude-3.5-Haiku			GPT-4o-mini			Qwen2.5-7B		
	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH
[+contact][containment][open]	88.9 ± 0.0	82.2 ± 11.5	88.9 ± 0.0	92.6 ± 6.4	100.0 ± 0.0	88.9 ± 0.0	49.6 ± 12.6	42.2 ± 14.6	46.7 ± 13.3
[-contact][containment][open]	40.7 ± 6.4	51.9 ± 12.8	40.7 ± 6.4	66.7 ± 0.0	59.3 ± 6.4	55.6 ± 0.0	38.5 ± 20.5	28.1 ± 6.8	22.2 ± 0.0
[+contact][containment][closed]	77.8 ± 0.0	74.1 ± 6.4	68.9 ± 0.0	63.0 ± 6.4	74.1 ± 6.4	66.7 ± 0.0	65.9 ± 11.4	43.0 ± 14.3	46.7 ± 14.7
[-contact][containment][closed]	40.7 ± 6.4	48.1 ± 6.4	43.5 ± 0.0	33.3 ± 0.0	44.4 ± 0.0	33.3 ± 0.0	34.1 ± 19.2	27.7 ± 13.4	3.7 ± 6.4
[+contact][attachment]	100.0 ± 0.0	80.0 ± 0.0	40.0 ± 0.0	60.0 ± 0.0	40.0 ± 0.0	38.7 ± 2.2	57.8 ± 4.8	11.1 ± 9.6	36.7 ± 5.8
[+contact][superposition]	66.7 ± 0.0	63.3 ± 8.8	84.4 ± 3.8	40.0 ± 0.0	53.3 ± 0.0	50.0 ± 5.8	27.8 ± 14.7	36.3 ± 15.1	48.5 ± 16.2
[-contact][superposition]	27.8 ± 1.9	51.1 ± 10.2	40.0 ± 0.0	40.0 ± 0.0	35.6 ± 7.7	26.7 ± 0.0	47.8 ± 1.9	32.2 ± 16.8	38.9 ± 9.6
[+contact][superposition][top]	100.0 ± 0.0	74.3 ± 4.7	83.3 ± 0.0	63.0 ± 12.8	63.0 ± 25.7	57.0 ± 1.3	62.2 ± 7.7	68.9 ± 16.8	71.1 ± 15.6
[-contact][superposition][top]	65.1 ± 18.0	54.8 ± 5.1	44.4 ± 0.0	44.4 ± 0.0	22.2 ± 0.0	66.7 ± 0.0	70.7 ± 8.0	55.9 ± 13.0	51.1 ± 21.4
[+contact][subposition]	63.3 ± 4.8	76.7 ± 8.8	100.0 ± 0.0	57.8 ± 0.0	40.0 ± 0.0	40.0 ± 0.0	63.0 ± 6.4	29.6 ± 12.8	48.1 ± 11.6
[-contact][subposition]	50.0 ± 16.7	66.7 ± 0.0	33.3 ± 0.0	27.4 ± 9.0	18.5 ± 3.2	0.0 ± 0.0	64.8 ± 3.2	37.0 ± 12.8	42.6 ± 22.5

Table 9: Projective Breakdown by Primitive Combination frontier (with standard deviation)

	Claude-3.5-Sonnet			GPT-4o			Qwen2.5-72B		
	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH
[absolute][east]	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	83.3 ± 16.7	50.0 ± 0.0	77.8 ± 9.6	11.1 ± 9.6	66.7 ± 0.0
[absolute][north]	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	55.6 ± 19.2	94.4 ± 9.6	66.7 ± 0.0	100.0 ± 0.0
[absolute][south]	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	77.8 ± 19.2	55.6 ± 9.6	83.3 ± 0.0	33.3 ± 0.0	100.0 ± 0.0
[absolute][west]	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	50.0 ± 0.0	66.7 ± 0.0	27.8 ± 9.6	66.7 ± 0.0
[intrinsic][back]	66.7 ± 0.0	100.0 ± 0.0	66.7 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	77.8 ± 19.2	0.0 ± 0.0	16.7 ± 0.0	16.7 ± 0.0
[intrinsic][front]	100.0 ± 0.0	66.7 ± 0.0	66.7 ± 0.0	94.4 ± 9.6	100.0 ± 0.0	100.0 ± 0.0	38.9 ± 9.6	16.7 ± 16.7	33.3 ± 0.0
[relative][back]	50.0 ± 0.0	50.0 ± 0.0	38.9 ± 9.6	77.8 ± 9.6	50.0 ± 0.0	38.9 ± 9.6	11.1 ± 19.2	0.0 ± 0.0	33.3 ± 0.0
[relative][front]	55.6 ± 9.6	33.3 ± 0.0	33.3 ± 0.0	50.0 ± 0.0	50.0 ± 0.0	44.4 ± 9.6	33.3 ± 0.0	33.3 ± 0.0	66.7 ± 0.0
[relative][left]	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	83.3 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0
[relative][right]	100.0 ± 0.0	100.0 ± 0.0	83.3 ± 0.0	100.0 ± 0.0	94.4 ± 9.6	38.9 ± 9.6	100.0 ± 0.0	100.0 ± 0.0	66.7 ± 0.0

Table 10: Projective Breakdown by Primitive Combination small (with standard deviation)

	Claude-3.5-Haiku			GPT-4o-mini			Qwen2.5-7B		
	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH
[absolute][east]	100.0 ± 0.0	100.0 ± 0.0	94.4 ± 9.6	83.3 ± 0.0	61.1 ± 9.6	100.0 ± 0.0	50.0 ± 0.0	50.0 ± 16.7	46.3 ± 17.9
[absolute][north]	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	88.9 ± 9.6	66.7 ± 0.0	94.4 ± 9.6	66.7 ± 0.0	38.9 ± 9.6	44.4 ± 9.6
[absolute][south]	100.0 ± 0.0	94.4 ± 9.6	100.0 ± 0.0	83.3 ± 0.0	33.3 ± 0.0	88.9 ± 9.6	38.9 ± 9.6	44.4 ± 9.6	55.6 ± 19.2
[absolute][west]	100.0 ± 0.0	100.0 ± 0.0	77.8 ± 19.2	77.8 ± 9.6	5.6 ± 9.6	88.9 ± 9.6	38.9 ± 9.6	38.9 ± 9.6	27.8 ± 9.6
[intrinsic][back]	50.0 ± 0.0	50.0 ± 0.0	50.0 ± 0.0	50.0 ± 0.0	50.0 ± 0.0	50.0 ± 0.0	11.1 ± 9.6	11.1 ± 9.6	0.0 ± 0.0
[intrinsic][front]	50.0 ± 0.0	50.0 ± 0.0	50.0 ± 0.0	50.0 ± 0.0	50.0 ± 0.0	50.0 ± 0.0	22.2 ± 25.5	35.2 ± 14.0	5.6 ± 9.6
[relative][back]	100.0 ± 0.0	100.0 ± 0.0	66.7 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	50.0 ± 16.7	55.6 ± 19.2	33.3 ± 16.7
[relative][front]	100.0 ± 0.0	100.0 ± 0.0	88.9 ± 19.2	100.0 ± 0.0	100.0 ± 0.0	94.4 ± 9.6	77.8 ± 19.2	50.0 ± 28.9	55.6 ± 19.2
[relative][left]	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	61.1 ± 34.7	66.7 ± 33.3	77.8 ± 38.5
[relative][right]	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	100.0 ± 0.0	94.4 ± 9.6	61.1 ± 25.5	11.1 ± 9.6	22.2 ± 9.6

Table 11: Chinese results with unified prompt. Figure/Ground: Levenshtein similarity; Predicate/Markers: avg. Levenshtein/type match; Dynamicity: accuracy; Configuration/Projective: F1-score.

	Chinese		
	Claude	GPT-4o	Qwen
Figure	100.0	99.4	99.4
Ground	85.3	89.9	92.9
Predicate	86.7	69.6	65.4
Markers	52.6	41.5	44.4
Dynamicity	98.4	100.0	90.5
Configuration	51.3	28.3	36.3
Projective	66.7	50.0	40.3

Table 12: Evaluation Performance by Component (multimodal) (%). Figure/Ground: Levenshtein similarity; Predicate/Markers: avg. Levenshtein/type match; Dynamicity: accuracy; Configuration/Projective: F1-score.

	Claude			Haiku			GPT-4o			GPT-mini			Qwen			Qwen-7B		
	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH
Figure	97.3	100.0	100.0	96.4	86.5	93.9	100.0	96.1	100.0	88.6	41.2	60.1	99.5	89.2	100.0	78.2	69.8	74.0
Ground	79.0	77.3	91.4	86.0	60.5	97.7	85.4	98.7	100.0	84.3	95.7	95.0	83.6	92.4	91.6	80.1	51.5	66.9
Predicate	99.2	91.0	67.7	86.2	88.7	69.3	98.4	96.6	64.2	48.0	89.2	59.9	88.2	85.7	62.2	59.9	35.2	37.5
Markers	73.0	50.7	83.2	74.3	17.8	66.0	78.9	42.0	67.8	39.8	15.0	50.3	71.6	15.6	55.0	54.7	16.3	28.0
Dynamicity	85.7	100.0	95.2	93.7	95.2	88.9	93.7	96.8	96.8	52.4	85.7	82.5	90.5	79.4	92.1	69.8	60.3	55.6
Configuration	62.9	70.1	58.9	33.1	29.9	33.7	63.9	58.1	50.4	27.8	11.8	28.0	47.1	45.6	47.8	33.0	26.0	25.5
Projective	78.6	81.8	70.8	65.0	70.0	47.7	90.0	61.8	54.4	39.6	55.0	43.3	63.3	50.0	66.7	38.4	33.9	39.1

Table 13: Performance Difference: Multimodal vs Text-only (Multimodal – Text-only, percentage points)

	Claude-3.5-Haiku			Claude-3.5-Sonnet			GPT-4o			GPT-4o-mini			Qwen2.5-72B			Qwen2.5-7B		
	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH
Configuration	-18.7	-12.7	+1.9	-8.1	+6.3	+3.7	-5.6	-9.6	+3.4	-1.9	-13.1	+3.0	-4.2	+3.0	-2.1	+4.7	+3.4	-0.7
Projective	-14.4	-12.8	-25.9	-8.1	+2.1	-2.6	-3.3	-25.3	-0.1	-0.9	+24.0	+1.3	+0.0	+6.7	+1.7	+2.0	+4.9	+12.8

Table 14: Configuration Breakdown Difference: Multimodal vs Text-only (big) (percentage points)

	Claude-3.5-Sonnet			GPT-4o			Qwen2.5-72B		
	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH
[+contact][containment][open]	+22.2	+44.4	+11.1	+38.9	+27.8	+22.2	+11.1	+5.6	+11.1
[-contact][containment][open]	+5.6	+0.0	-11.1	+22.2	+11.1	+22.2	+22.2	+11.1	+22.2
[+contact][containment][closed]	+0.0	-5.6	-11.1	-16.7	-13.9	-27.8	-11.1	+5.6	-11.1
[-contact][containment][closed]	+11.1	+5.6	-11.1	-11.1	-11.1	-22.2	-5.6	+0.0	-11.1
[+contact][attachment]	+0.0	+20.0	+0.0	-13.9	-2.2	-13.3	-23.3	+6.7	+1.7
[+contact][superposition]	-16.7	-11.1	-10.0	-8.9	-14.4	-0.0	+25.0	+17.2	-22.2
[-contact][superposition]	-46.7	-13.3	-36.7	+13.3	+10.6	+10.6	-5.6	+10.0	-12.2
[+contact][superposition][top]	-30.0	+23.3	+46.7	-23.3	-66.7	+20.0	+0.0	+4.4	+13.3
[-contact][superposition][top]	-24.4	+2.2	+15.6	+5.6	-13.3	-8.9	+22.2	+22.2	+0.0
[+contact][subposition]	-23.3	-3.3	-20.0	-16.1	-1.1	-6.7	-60.0	-23.3	-40.0
[-contact][subposition]	-33.3	-50.0	-16.7	-25.0	-33.3	-15.0	+8.3	-25.0	+0.0

Table 15: Configuration Breakdown Difference: Multimodal vs Text-only (small) (percentage points)

	Claude-3.5-Haiku			GPT-4o-mini			Qwen2.5-7B		
	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH
[+contact][containment][open]	+0.0	-12.2	+11.1	+5.6	-66.7	+11.1	+41.1	-8.6	+37.8
[-contact][containment][open]	+0.0	+0.0	+16.7	+0.0	-61.1	+11.1	+13.3	+37.8	+44.4
[+contact][containment][closed]	-11.1	-22.2	-2.2	+11.1	-38.9	+0.0	-5.6	+26.7	+6.1
[-contact][containment][closed]	-5.6	+0.0	-10.2	+5.6	-33.3	+0.0	-11.1	+14.1	+22.2
[+contact][attachment]	+0.0	-13.3	+0.0	-8.3	-40.0	+1.9	+0.6	+36.1	+26.7
[+contact][superposition]	-33.3	-40.0	-61.1	+13.3	+0.0	-8.3	+19.4	+5.0	+23.9
[-contact][superposition]	+26.7	+3.3	+33.3	-10.0	+15.6	-2.2	-3.3	+0.0	-16.7
[+contact][superposition][top]	-100.0	-59.2	+16.7	-11.1	-55.6	+42.2	-16.7	-25.6	-30.0
[-contact][superposition][top]	-53.2	-31.1	+31.1	-26.7	+0.0	+0.0	-20.0	-32.2	-43.3
[+contact][subposition]	-24.4	-15.0	-60.0	-17.8	-40.0	+0.0	+18.9	+22.2	+38.3
[-contact][subposition]	+5.0	-16.7	+11.1	-30.0	+13.9	+0.0	-2.8	-11.1	-30.6

Table 16: Projective Breakdown Difference: Multimodal vs Text-only (big) (percentage points)

	Claude-3.5-Sonnet			GPT-4o			Qwen2.5-72B		
	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH
[absolute][east]	+0.0	+0.0	+0.0	+0.0	-33.3	+0.0	+16.7	-8.3	+33.3
[absolute][north]	+0.0	+0.0	+0.0	+0.0	-16.7	-16.7	-16.7	+33.3	+0.0
[absolute][south]	+0.0	+0.0	+0.0	+0.0	+0.0	+8.3	+16.7	+66.7	+0.0
[absolute][west]	+0.0	+0.0	+0.0	+0.0	-33.3	+16.7	+33.3	+50.0	+33.3
[intrinsic][back]	+33.3	+0.0	+33.3	+0.0	+0.0	+16.7	+50.0	-16.7	+50.0
[intrinsic][front]	+0.0	+33.3	+33.3	+0.0	+0.0	+0.0	-33.3	+0.0	+0.0
[relative][back]	+0.0	+0.0	+16.7	-8.3	+0.0	+8.3	-16.7	+0.0	-33.3
[relative][front]	+8.3	+16.7	+16.7	+0.0	+0.0	+8.3	-33.3	-33.3	-66.7
[relative][left]	+0.0	+0.0	+0.0	+0.0	-50.0	-33.3	+0.0	+0.0	+0.0
[relative][right]	+0.0	+0.0	-33.3	-16.7	-33.3	+16.7	+0.0	+0.0	+0.0

Table 17: Projective Breakdown Difference: Multimodal vs Text-only (small) (percentage points)

	Claude-3.5-Haiku			GPT-4o-mini			Qwen2.5-7B		
	ES	EU	ZH	ES	EU	ZH	ES	EU	ZH
[absolute][east]	+0.0	+0.0	+0.0	+8.3	+25.0	+0.0	+16.7	+8.3	+33.3
[absolute][north]	+0.0	+0.0	+0.0	-8.3	+33.3	+8.3	+33.3	+58.3	+25.0
[absolute][south]	+0.0	+8.3	+0.0	+8.3	+16.7	+8.3	+58.3	+25.0	+0.0
[absolute][west]	+0.0	+0.0	+16.7	+16.7	+50.0	+16.7	+8.3	+16.7	+16.7
[intrinsic][back]	+0.0	+0.0	-50.0	-8.3	-16.7	+0.0	-8.3	-8.3	+0.0
[intrinsic][front]	+0.0	+0.0	-33.3	-8.3	-16.7	+0.0	+8.3	+8.3	+50.0
[relative][back]	+0.0	-33.3	-66.7	-8.3	-16.7	-16.7	+8.3	+0.0	-8.3
[relative][front]	-33.3	-50.0	-33.3	+0.0	+0.0	+8.3	+16.7	+33.3	+33.3
[relative][left]	+0.0	+0.0	+0.0	-16.7	+0.0	+0.0	+33.3	+50.0	+16.7
[relative][right]	+0.0	+0.0	+0.0	-33.3	+0.0	-16.7	+50.0	+33.3	+41.7