

# Family Matters: Cross-Provider LLM-Judge Committees and Iterative Prompts for Turkish PropBank Argument-Frame Prediction

Anonymous ACL submission

## Abstract

Large language models are increasingly used not only to label data but also to judge whether another model’s labels are correct. Using Turkish PropBank argument-frame prediction as a testbed, we study two questions. First, does agreement between LLM judges reflect shared capability or shared provider family? Across four judges from two provider families (two Gemini, two OpenAI), within-family agreement consistently exceeds cross-family agreement. However, a two-judge committee combining one model from each family improves exact-match precision by 11.9 points over the unfiltered baseline, more than doubling the gain of the best single-judge filter, while adding a third same-family judge yields little benefit. Second, how do targeted prompt revisions affect other error types? A prompt edit for Turkish intransitive activity verbs improves the intended ARG0-only class by 21.7 points but also increases ARG2 omissions by 30.9 points; a follow-up clarification partially recovers overall performance. These results suggest that LLM-judge committees benefit more from family diversity than larger committee size, and that prompt revisions for closed semantic taxonomies should be evaluated by error class rather than aggregate metrics alone.

## 1 Introduction

PropBank-style annotation links each predicate sense to the semantic participants licensed by that sense (Palmer et al., 2005). These participants are represented with core roles such as ARG0, ARG1, and ARG2.<sup>1</sup> The framework is closely related to Frame Semantics (Fillmore, 1976; Baker et al., 1998) and has been widely studied in semantic role labeling (Gildea and Jurafsky, 2000; Carreras and Màrquez, 2005; He et al., 2017; Shi and Lin, 2019).

<sup>1</sup>We adopt the PropBank core-role inventory (ARG0–ARG4), but predict only ARG0, ARG1, and ARG2 because they cover more than 99.5% of valid frames in the gold standard.

Turkish makes sense-level argument-frame prediction particularly challenging because evidence for semantic roles is distributed across lexical meaning, derivational morphology, and case marking. As an agglutinative language, Turkish encodes voice, valency, and case-related information through suffixation (Ofazer, 1993; Göksel and Kerslake, 2005; Marşan et al., 2021). As a result, surface syntax does not always transparently reveal the underlying argument frame. A verb may appear intransitive while still licensing both ARG0 and ARG1, and dative-marked complements may function as core arguments rather than adjuncts. These distinctions make Turkish PropBank annotation highly dependent on expert linguistic analysis.

Large language models (LLMs) offer a potential way to reduce this annotation burden, but they also introduce new methodological questions. Recent work shows that LLMs can predict frame-like semantic structures from definitions and examples (Tayyar Madabushi et al., 2025), especially when combined with retrieval-augmented prompting (Lewis et al., 2020; Brown et al., 2020). However, prompt revisions are often evaluated only through aggregate metrics, without examining which error classes improve or deteriorate. In parallel, LLMs are increasingly used as evaluators (Zheng et al., 2023; Wang et al., 2024), yet it remains unclear whether agreement among LLM judges reflects semantic reliability or shared provider-family biases.

This paper studies Turkish PropBank argument-frame prediction as a closed multi-label task. Given a predicate sense, definition, and optional example sentence, models predict which of ARG0, ARG1, and ARG2 are licensed. We evaluate five prediction models under zero-shot, static few-shot, and retrieval-based prompting on 500 held-out senses from a 17,692-sense gold standard, and add a targeted GPT-5-mini ablation to test whether a reasoning-enabled OpenAI model changes the

081	family-level picture.	
082	Our findings are threefold. First, LLM-judge	setting differs in that we predict the set of roles
083	agreement shows a clear but asymmetric provider-	licensed by a predicate sense rather than label-
084	family effect: within-family agreement exceeds	ing spans in context. <a href="#">Tayyar Madabushi et al.</a>
085	cross-family agreement, but Gemini judges are sub-	(2025) show that LLMs can generate FrameNet-
086	stantially more internally consistent than OpenAI	style frame assignments, motivating our use of gen-
087	judges. Second, two-judge committees that com-	erative models for closed multi-label PropBank
088	bine one OpenAI and one Gemini model outper-	prediction.
089	form larger single-family committees, improving	
090	exact-match precision by 11.9 points over the un-	<b>In-context learning and retrieval-augmented</b>
091	filtered baseline. Third, targeted prompt edits re-	<b>prompting.</b> In-context learning enables LLMs
092	distribute errors rather than uniformly improving	to adapt from instructions and examples without
093	performance: guidance that fixes ARG0 errors for	parameter updates ( <a href="#">Brown et al., 2020</a> ), while re-
094	Turkish intransitive activity verbs also substantially	trieval methods improve prompting by selecting
095	increases ARG2 omissions. Together, these results	relevant demonstrations at test time ( <a href="#">Lewis et al.,</a>
096	suggest that LLM-judge committees benefit more	<a href="#">2020</a> ). Prior work shows that demonstration selec-
097	from family diversity than committee size, and that	tion strongly affects performance ( <a href="#">Liu et al., 2022;</a>
098	prompt revisions for closed semantic taxonomies	<a href="#">Rubin et al., 2022; Min et al., 2022</a> ). We therefore
099	should be evaluated by error class rather than ag-	compare static and retrieved few-shot prompting us-
100	gregate metrics alone.	ing BM25 ( <a href="#">Robertson and Zaragoza, 2009</a> ), dense
101	<b>2 Related Work</b>	multilingual sentence representations ( <a href="#">Reimers and</a>
102	<b>Turkish PropBank and Turkish NLP infrastruc-</b>	<a href="#">Gurevych, 2019; Karpukhin et al., 2020; Devlin</a>
103	<b>ture.</b> Turkish PropBank extends PropBank-style	<a href="#">et al., 2019; Conneau et al., 2020</a> ), and Maximal
104	semantic role annotation to Turkish verb senses	Marginal Relevance for diversity ( <a href="#">Carbonell and</a>
105	by recording which core roles each sense licenses.	<a href="#">Goldstein, 1998</a> ).
106	The resource builds on earlier Turkish NLP infras-	
107	tructure, including morphological analysis ( <a href="#">Ofłazer,</a>	<b>LLM-as-a-Judge and evaluator bias.</b> LLM-
108	<a href="#">1993; Yıldız et al., 2019</a> ) and Turkish Universal De-	based evaluation has become increasingly common
109	pendencies treebanks ( <a href="#">Şahin and Adalı, 2018; Tsar-</a>	( <a href="#">Zheng et al., 2023; Wang et al., 2024</a> ), but prior
110	<a href="#">faty et al., 2013</a> ). The sense inventory of <a href="#">Marşan</a>	work shows that LLM judges are sensitive to fac-
111	<a href="#">et al. (2021)</a> provides ARG0–ARG4 role descrip-	tors such as output length, verbosity, and presenta-
112	tions in a synset-style organisation; our experi-	tion order ( <a href="#">Liu et al., 2023; Chiang and Lee, 2023;</a>
113	ments use the resulting inventory of 17,692 senses	<a href="#">Dubois et al., 2023</a> ). These concerns are especially
114	with at least one filled core-role field.	relevant for semantic role prediction, where outputs
115	The distinctions relevant to this task are stan-	may be nearly correct despite differing by a single
116	dard in semantic role theory. Unergative intrans-	role. We therefore examine not only agreement
117	itives typically license an agent-like ARG0, un-	with gold labels, but also whether judges from the
118	accusatives license an undergoer-like ARG1, and	same provider family behave more similarly than
119	transitive or triadic predicates may additionally li-	judges from different families.
120	cence ARG2 ( <a href="#">Levin, 1993; Göksel and Kerslake,</a>	
121	<a href="#">2005</a> ). In Turkish, these distinctions are often ex-	<b>Evaluation methodology.</b> Because each predi-
122	pressed through morphology and case marking,	cate sense may license multiple roles, the task is a
123	making sense-level frame prediction linguistically	multi-label classification problem ( <a href="#">Tsoumakas and</a>
124	challenging.	<a href="#">Katakis, 2007</a> ). We report exact match together
125	<b>Semantic role labelling and generative frame</b>	with macro- and micro-F1. For judge reliability, we
126	<b>prediction.</b> Traditional semantic role labelling	use quadratic-weighted Cohen’s $\kappa$ ( <a href="#">Cohen, 1960;</a>
127	identifies argument spans in a sentence and assigns	<a href="#">Landis and Koch, 1977</a> ). Prompt comparisons are
128	semantic roles using statistical or neural models	evaluated with paired bootstrap confidence inter-
129	( <a href="#">Gildea and Jurafsky, 2000; Carreras and Màrquez,</a>	vals ( <a href="#">Efron and Tibshirani, 1993; Dror et al., 2018</a> ),
130	<a href="#">2005; He et al., 2017; Shi and Lin, 2019</a> ). Our	which estimate the stability of differences under
		resampling of matched evaluation items.



## 4 Experiments

We evaluate six LLMs from three provider families (Table 1). The Google Gemini models (Team et al., 2025) represent a recent proprietary API family. The OpenAI models include GPT-4o-mini and GPT-5-mini (reasoning-enabled), following instruction-tuned model lines (Ouyang et al., 2022). We also include Qwen2.5-7B-Instruct as an open-weights baseline (Qwen et al., 2025), similar in spirit to other open releases such as LLaMA (Touvron et al., 2023).

Provider	Model	Access
Google	gemini-2.5-flash	API
Google	gemini-2.5-pro	API
Google	gemini-3-pro-preview	API
OpenAI	gpt-4o-mini	API
OpenAI	gpt-5-mini (reasoning)	API
Alibaba	qwen2.5:7b-instruct	local

Table 1: Models evaluated.

**Decoding.** We use deterministic decoding when supported (temperature 0.0, seed 42). For models with controllable reasoning, we use the lowest available reasoning budget unless stated otherwise. Outputs that fail JSON parsing are treated as empty predictions, which penalizes recall rather than being discarded. We analyze GPT-5-mini separately in Section 5.3 because its default reasoning settings occasionally produce structured-output failures.

**Sample size.** We run a 50-item pilot for the first prompt tour and full evaluations on 500 held-out senses for the second and third tours. The pilot includes nine Gemini runs to enable matched comparisons between prompt versions. The final tour uses the retrieved few-shot setup for all main models, with an additional ablation over prompting setups for GPT-5-mini.

## 5 Results

### 5.1 Frame Prediction across Models and Setups

Table 2 reports exact match, macro-F1, and micro-F1 for the second prompt tour, where exact match requires the full predicted role set to match the gold set. Across all models, retrieved few-shot prompting is the strongest setup, indicating that semantically similar annotated senses provide useful grounding for both lexical choice and argument structure. The best overall result is obtained by

Gemini 2.5 Flash with retrieval (EM = 0.628). Figure 1 confirms this ranking across all metrics. The smaller Qwen model is consistently behind the API models, especially on cases involving ARG2, which require more fine-grained sense distinctions.

Averaging over non-reasoning models, retrieval improves exact match (0.591) over static few-shot (0.552) and zero-shot (0.542). Gains are more stable at the full evaluation size than in the pilot, where retrieval and zero-shot were nearly tied. Error analysis shows that remaining mistakes are structurally plausible rather than random, as shown in Figure 2: most errors collapse three-role frames into  $\{ARG0, ARG1\}$  or confuse unergative ARG0-only cases with two-role frames.

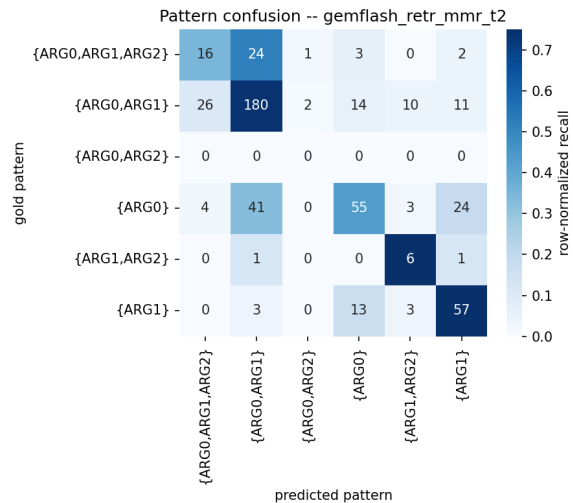


Figure 2: Pattern confusion for the best t2 prediction run (gemini-2.5-flash with retrieved few-shot prompting).

### 5.2 Prompt Iteration across Three Tours

We next analyze how prompt revisions change error behavior. For tours 1 and 2, we compare nine matched Gemini configurations (same model and setup before/after revision). While tour 1 uses a smaller pilot sample, the direction of changes is highly consistent across runs.

**Effect of targeted edits.** The t2 revision successfully fixes its intended issue: ARG0-only accuracy increases by 21.7 points across all runs, with consistent directionality and a confidence interval excluding zero.

**Side effects.** However, the same change redistributes errors elsewhere in the label space. Performance drops for multi-argument patterns, and

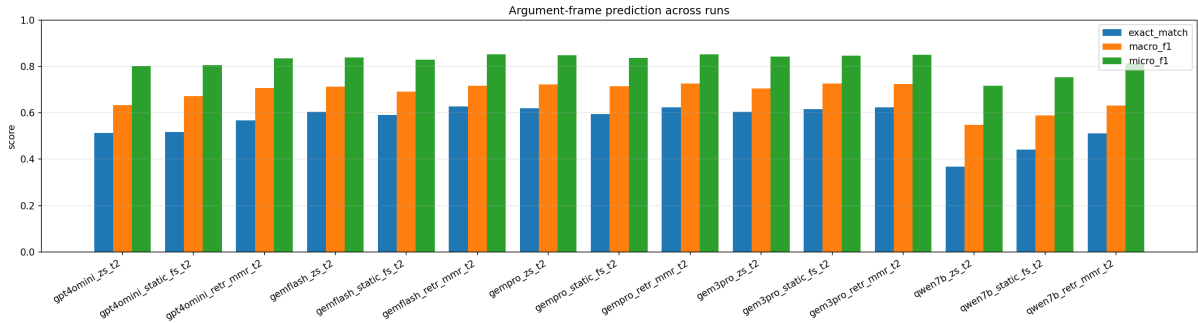


Figure 1: Main t2 benchmark across models and prompting setups. Retrieved few-shot prompting gives the strongest run within every model family, while the Qwen-7B runs lag most clearly on the strict exact-match metric.

Model	Setup	EM	macro-F1	micro-F1	ARG0 leak	ARG2 miss
gemini-2.5-flash	retr_mmr	<b>0.628</b>	0.717	0.853	0.202	0.574
gemini-2.5-pro	retr_mmr	0.624	0.726	0.853	0.083	0.481
gemini-3-pro-preview	retr_mmr	0.624	0.724	0.851	0.119	0.537
gemini-2.5-pro	zs	0.620	0.722	0.849	0.131	0.500
gemini-3-pro-preview	static_fs	0.616	0.726	0.847	0.107	0.444
gemini-2.5-flash	zs	0.604	0.713	0.838	0.214	0.500
gemini-3-pro-preview	zs	0.604	0.705	0.843	0.143	0.574
gemini-2.5-pro	static_fs	0.594	0.715	0.838	0.119	0.444
gemini-2.5-flash	static_fs	0.590	0.692	0.829	0.190	0.556
gpt-4o-mini	retr_mmr	0.568	0.707	0.835	0.119	0.556
gpt-4o-mini	static_fs	0.518	0.673	0.805	0.107	0.574
gpt-4o-mini	zs	0.514	0.633	0.802	0.119	0.778
qwen2.5-7b-instruct	retr_mmr	0.512	0.632	0.814	0.179	0.833
qwen2.5-7b-instruct	static_fs	0.442	0.589	0.754	0.119	0.796
qwen2.5-7b-instruct	zs	0.368	0.547	0.717	0.060	0.833

Table 2: Argument-frame prediction at t2,  $n=500$ . Sorted by EM.

ARG2 misses increase by 30.9 points. This indicates that correcting one linguistic distinction can shift decision boundaries for others in a closed label space.

**Tour 3 correction.** Tour 3 addresses this regression by clarifying that morphological cues are *sufficient but not necessary* for ARG2 licensing. As shown in Table 4, all Gemini models improve (up to +2.4 EM), while GPT-4o-mini degrades slightly and Qwen remains mostly unchanged. This suggests that the correction helps primarily within the Gemini family.

Overall, prompt improvements are not monotonic: they tend to shift errors between related classes rather than uniformly improving all metrics.

### 5.3 Reasoning-Model Ablation: gpt-5-mini

To extend OpenAI-family coverage, we evaluate gpt-5-mini under the t3 prompt. It improves over GPT-4o-mini in all three prompting setups, with exact-match gains of 1.2–3.8 points and substantially higher ARG0-only accuracy. However, the

ablation also exposed a structured-output issue: with the default token budget and reasoning mode, 15–20% of outputs fail JSON parsing because hidden reasoning can consume the completion budget before the required JSON is emitted. Setting `reasoning_effort` to minimal and increasing `max_completion_tokens` eliminates parse failures and makes decoding roughly three times faster, so reasoning-model results should be reported only after explicit output-budget calibration. The full ablation table is provided in Appendix A.

### 5.4 Judge Reliability and Provider-Family Clustering

We apply four judges to the 500 predictions from the best second-tour run (Gemini 2.5 Flash with retrieved few-shot prompting). Table 5 and Figure 3 report pairwise quadratic-weighted Cohen’s  $\kappa$  on the 496 cases with valid scores from all judges.

**Asymmetric cohesion.** The effect is asymmetric across families. Gemini judges are much more internally consistent (Gemini–Gemini  $\kappa = 0.701$  vs OpenAI–OpenAI  $\kappa = 0.376$ , about  $1.9\times$  higher).

Pattern	t1 mean	t2 mean	$\Delta_{t_2-t_1}$	95% CI (bootstrap)	sign consistency
{ARG0} unergative (TARGETED)	0.157	0.375	+0.217	[+0.166, +0.261]	9/9 $\uparrow$
{ARG1} unaccusative	0.833	0.775	-0.058	[-0.104, -0.009]	6/9 $\downarrow$
{ARG0,ARG1} agent-target	0.815	0.716	-0.099	[-0.128, -0.078]	9/9 $\downarrow$
{ARG1,ARG2} passive/recip.	1.000	0.694	-0.306	[-0.375, -0.236]	9/9 $\downarrow$
{ARG0,ARG1,ARG2} triadic	0.644	0.430	-0.215	[-0.318, -0.104]	8/9 $\downarrow$
<b>ARG2 miss rate</b>	0.204	0.512	+0.309	[+0.239, +0.377]	9/9 worse
ARG0 leakage rate	0.124	0.146	+0.022	[-0.027, +0.076]	5/9 worse

Table 3: Paired per-pattern and diagnostic-rate deltas,  $t_1 \rightarrow t_2$ , across the nine Gemini runs (3 models  $\times$  3 setups). Bootstrap 5000-resample 95% CIs. Bold rows are statistically significant regressions or gains (CI excludes 0).

Model	EM t2	EM t3	$\Delta$
gemini-2.5-flash	0.628	<b>0.652</b>	<b>+0.024</b>
gemini-3-pro-preview	0.624	0.634	+0.010
gemini-2.5-pro	0.624	0.628	+0.004
qwen2.5-7b-instruct	0.512	0.514	+0.002
gpt-4o-mini	0.568	0.554	-0.014

Table 4: Paired  $t_2 \rightarrow t_3$  EM deltas on the RETR\_MMR setup ( $n=500$  each). All three Gemini variants gain; gpt-4o-mini regresses; Qwen-7B is flat on aggregate EM but loses  $-0.087$  on the {ARG0,ARG1,ARG2} triadic pattern. The Gemini-vs-OpenAI signature parallels the inter-judge  $\kappa$  stratification of Table 5.

Judge 1	Judge 2	$\kappa_{\text{quad}}$
<i>Within-family pairs</i>		
<b>gemini-2.5-pro</b>	<b>gemini-3-flash-preview</b>	<b>0.701</b>
<b>gpt-4o-mini</b>	<b>gpt-5-mini</b>	<b>0.376</b>
<i>Cross-family pairs</i>		
gemini-3-flash-preview	gpt-5-mini	0.351
gemini-2.5-pro	gpt-5-mini	0.303
gemini-3-flash-preview	gpt-4o-mini	0.266
gemini-2.5-pro	gpt-4o-mini	0.239

Table 5: Pairwise quadratic-weighted Cohen’s  $\kappa$  across four judges ( $n=496$ ). The Gemini-Gemini pair reaches “substantial” agreement (Landis and Koch, 1977); both within-family OpenAI and all cross-family pairs sit in the “fair” band (0.21–0.40), but within-family OpenAI (0.376) is consistently above the cross-family range (0.24–0.35). Families therefore differ in internal cohesion: Gemini is about  $1.9\times$  as internally consistent as OpenAI under this metric.

381 Notably, the OpenAI pair is closer to the strongest  
382 cross-family agreement (0.35) than to the Gemini  
383 pair. This is partly explained by score distributions:  
384 gpt-4o-mini is the strictest judge ( $\bar{s} = 3.53$ ),  
385 Gemini judges are more lenient ( $\bar{s} = 4.47, 4.49$ ),  
386 and gpt-5-mini is intermediate ( $\bar{s} = 4.03$ ), yield-  
387 ing correspondingly intermediate agreement pat-  
388 terns.

**Implication.** Judge agreement is not a single scalar property of “provider family.” High within-family agreement can reflect strong calibration consistency, while cross-family disagreement can reflect scale differences rather than semantic disagreement. For evaluation, adding a second judge from the same provider may add little new information, whereas cross-family pairing provides a more independent signal at similar cost.

## 398 5.5 Committee-Filtered Precision

399 We next evaluate whether LLM judges can identify a high-precision subset of predictions. Table 6 reports the main coverage-precision frontier; Appendix Figure 5 visualizes the same trade-off. The committee analysis uses the original three-judge panel (gpt-4o-mini, gemini-2.5-pro, gemini-3-flash-preview); the fourth judge (gpt-5-mini) is reserved for Section 5.4 to isolate

407 within-family OpenAI agreement and is excluded  
408 here to keep the comparison consistent with earlier  
409 results.

**Committee filtering.** The strongest filter requires 410 both GPT-4o-mini and Gemini 2.5 Pro to assign the 411 maximum score. This retains 41.3% of items while 412 improving exact match by 11.9 points over the full 413 set, more than twice the gain of the best single- 414 judge filter. Adding the third same-family Gemini 415 judge yields only a 0.3-point change, indicating 416 that cross-family diversity is more important than 417 committee size. 418

**Score spread as a diagnostic.** Judge disagree- 419 ment is also informative, but not monotonically so. 420 Figure 4 reports exact match by three-judge score 421 spread. Flagging cases with a spread of at least two 422 captures roughly one third of the data but selects 423 a lower-accuracy subset overall. Large disagree- 424 ments often correspond to genuinely ambiguous 425 or polysemous predicates rather than clear errors, 426 while moderate disagreement is a better indicator 427

Filter strategy	$N$	Coverage	EM on subset	$\Delta\text{EM vs. full}$
Full set (no filter)	496	100.0%	0.627	—
<i>Single-judge filters</i>				
gpt-4o-mini $\geq 4$	277	55.8%	0.682	+0.055
gemini-2.5-pro $\geq 4$	394	79.4%	0.673	+0.046
<i>Two-judge committees (cross-family)</i>				
Either gpt or pro $\geq 4$ (union)	419	84.5%	0.659	+0.032
Both gpt and pro $\geq 4$ (intersection)	252	50.8%	0.706	+0.079
Both gpt and pro = 5 (unanimous high)	205	41.3%	<b>0.746</b>	<b>+0.119</b>
<i>Three-judge committees</i>				
$\geq 2$ of 3 say $\geq 4$ (majority)	389	78.4%	0.663	+0.036
Unanimous $\geq 4$	244	49.2%	0.709	+0.082
Unanimous = 5	199	40.1%	<b>0.749</b>	<b>+0.122</b>
<i>Spread-based diagnostic</i>				
Spread $\geq 2$ (route to human)	161	32.5%	0.596	-0.031

Table 6: Coverage–precision frontier for filter strategies. Bootstrap 95% CI for the (unanimous = 5 EM – full EM) lift: [+0.046, +0.193].

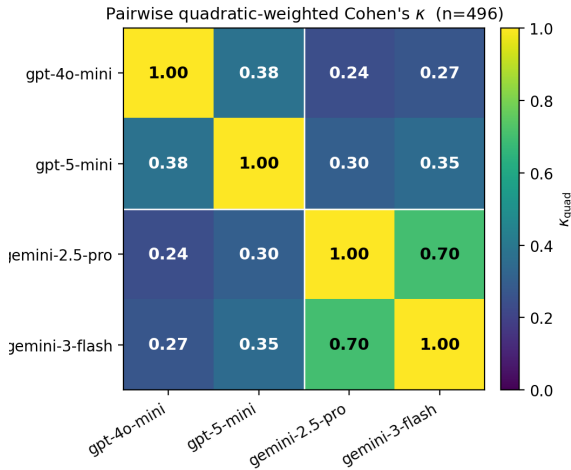


Figure 3: Four-judge agreement matrix. The two within-family blocks (OpenAI top-left, Gemini bottom-right) are both denser than the off-diagonal cross-family block, but the two within-family blocks differ in magnitude: Gemini–Gemini is 0.70, OpenAI–OpenAI is 0.38. The asymmetry visualises “provider-family clustering” as a real but non-symmetric effect rather than a uniform property.

of likely mistakes.

## 5.6 Qualitative Case Analysis

Representative prediction and judge outputs are provided in Appendix Table 9. The examples illustrate the main mechanisms behind the aggregate results: t2 repairs ARG0-only errors, ARG2 is under-generated after the morphology-focused revision, t3 recovers some lexically licensed ARG2 cases, and same-family judges can accept the same wrong frame.

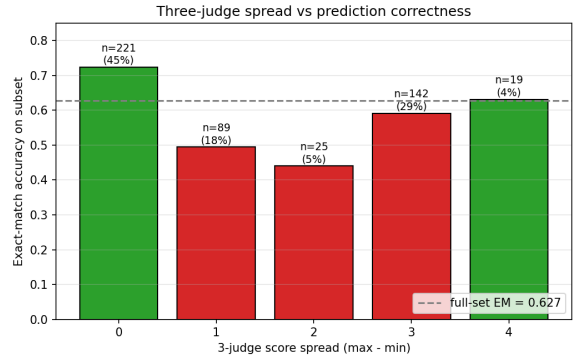


Figure 4: Exact-match accuracy by three-judge score spread. Moderate spread is the clearest warning signal; maximum spread is rare and not itself the lowest-accuracy bucket.

## 6 Discussion

### Why does the unergative guidance help?

ARG0-only senses are challenging because a single participant must be recognized as agent-like rather than defaulting to ARG1. In Turkish, additional material is often expressed as an oblique phrase (e.g., constructions similar to “laugh at X” or “cry for X”), which models may incorrectly treat as a core ARG1 argument. The revised prompt addresses this by introducing a simple diagnostic: check whether the verb sense can take a true accusative-marked theme. If not, the sole animate participant is more likely to be ARG0. This concrete test is easier to apply reliably than an abstract definition of unergativity.

**Why does the morphology guidance hurt ARG2?** The morphology block is linguistically

455	correct, but it changes how the model weights evi-	— should expect prompt revisions to redistribute	505
456	dence. Listing suffixes that often affect valency ap-	rather than purely add accuracy. The unergative-	506
457	pears to encourage the model to treat morphology as	versus-unaccusative split that hurts us in Turk-	507
458	a near-exhaustive signal for argument structure.	ish has direct analogues in, for example, Hungar-	508
459	This is harmful for ARG2, since some predicates li-	ian, Finnish, or Basque, where the same single-	509
460	cence a third role lexically rather than through overt	participant verb can be agent-like or undergoer-like	510
461	derivational morphology. For example, verbs like	depending on the sense. Likewise, any pipeline	511
462	<i>açıklamak</i> ‘explain’ can express a three-participant	that uses LLM judges as a quality gate should ex-	512
463	relation (someone explains something to someone),	pect intra-family judge pairs to share calibration	513
464	and constructions such as <i>farklı olmak</i> ‘be different	biases that an inter-family pair would expose; this	514
465	from’ also introduce a third participant without the	is a property of how the judges were trained, not	515
466	listed morphological markers. The third-tour clari-	of which language they are judging. The practi-	516
467	fication corrects this by stating that morphology is	cal recommendation — pair across families before	517
468	evidence for ARG2 but not a requirement.	adding within a family — therefore travels with the	518
469	<b>Provider family vs. scale.</b> The judge results	methodology, not with the dataset.	519
470	show that large API models are not interchange-		
471	able, but the structure is more nuanced than a single	<b>7 Conclusion</b>	520
472	“family clustering” effect. Within-family agree-		
473	ment is always higher than cross-family agreement,	We presented an iterative prompt-engineering and	521
474	yet internal cohesion differs substantially across	LLM-judge committee pipeline for Turkish Prop-	522
475	families: Gemini judges agree strongly ( $\kappa = 0.70$ ),	Bank argument-frame prediction. Three findings	523
476	while OpenAI judges are much less consistent	extend prior work: (i) within-family inter-judge	524
477	( $\kappa = 0.38$ ). Within OpenAI, gpt-5-mini behaves	agreement consistently exceeds cross-family agree-	525
478	as an intermediate scorer, with mean score and	ment, but families differ markedly in internal cohe-	526
479	pairwise $\kappa$ values between gpt-4o-mini and the	sion ( $\kappa = 0.70$ for two Gemini judges vs. $\kappa = 0.38$	527
480	Gemini models, suggesting that reasoning-tuned	for two OpenAI judges across four judges and	528
481	successors can shift scoring scales within the same	496 paired predictions), so provider-family effects	529
482	provider ecosystem. The practical implication re-	are real but not symmetric; (ii) a two-judge cross-	530
483	mains that cross-family judge pairs are preferable,	family committee matches three-judge unanimous-	531
484	but the benefit of adding a second same-family	5 precision at higher coverage, supporting commit-	532
485	judge depends on how internally calibrated that	tee design by family diversity over committee size;	533
486	family is.	and (iii) targeted prompt blocks redistribute rather	534
487	<b>When to consult humans.</b> The most useful trig-	than purely add accuracy in closed-set taxonomies,	535
488	ger for human review is not low confidence or any	motivating systematic paired evaluation of prompt	536
489	form of judge disagreement. Large disagreements	revisions.	537
490	often reflect stable differences between judge fami-		
491	lies on genuinely ambiguous senses. Instead, mod-	For practitioners, the most actionable result is	538
492	erate disagreement is more informative, as it tends	the committee design. If an LLM judge is used	539
493	to identify cases near the decision boundary. For	as a quality gate, pair two judges from different	540
494	annotation support systems, this suggests that hu-	providers and treat a third same-family judge as	541
495	man effort should be focused on near-miss cases	redundant. The compute saved on that third judge	542
496	rather than all instances where judges differ due to	is better spent either on routing uncertain cases to	543
497	scale or calibration effects.	a human annotator or on adding a fourth model	544
498	<b>Why this matters beyond Turkish.</b> Our dataset	from a yet different ecosystem; in our experiments,	545
499	is Turkish, but the two main findings are not	the marginal precision gained by a same-family	546
500	Turkish specific. Any closed semantic taxonomy	addition was within rounding error of zero.	547
501	that depends on linguistically subtle distinctions		
502	— predicate–argument structure in other morpho-	Future work will integrate LoRA fine-tuning,	548
503	logically rich languages, fine-grained relation ex-	extend judge benchmarking to local models, and	549
504	traction, sense disambiguation, frame assignment	apply the spread-based diagnostic in a deployed	550
		Turkish PropBank annotation interface.	551

552	<b>8 Limitations</b>	
553	<b>Sample-size confound in paired t1→t2.</b>	
554	The pilot t1 evaluation uses $n=50$ whereas t2/t3 use	600
555	$n=500$ ; the paired delta therefore mixes prompt	601
556	change and sample change. We mitigate this by	602
557	reporting directional consistency (9/9 across pat-	603
558	terns) and bootstrap CIs but cannot fully isolate	604
559	the prompt effect without an additional $n=500$ run	605
560	under the t1 prompt.	606
561	<b>Single fine-tuned baseline.</b>	
562	We do not parameter-efficiently fine-tune (e.g., via LoRA (Hu et al.,	607
563	2021)) a local model on PropBank. Whether do-	608
564	main adaptation of a small open model would out-	609
565	perform the API baselines on multi-label PropBank	610
566	set prediction — as has been reported for single-	611
567	label frame tasks in related work — is an open	612
568	question.	613
569	<b>Judge coverage.</b>	
570	Four judges are evaluated, all API-scale, drawn from two provider families (two	614
571	Gemini, two OpenAI). Other API-scale judges	615
572	(e.g., Anthropic Claude, Mistral) and local-model	616
573	judges (Qwen, Llama, Gemma) are not tested	617
574	here for the PropBank setting, so we do not know	618
575	whether a third provider family or open-weights	619
576	judges would behave more like the internally cohe-	620
577	sive Gemini family, the more heterogeneous Ope-	621
578	nAI family, or something else entirely.	622
579	<b>Provider-family confounding.</b>	
580	The judge comparison is observational. Provider family, training	623
581	data, alignment procedure, score calibration, and	624
582	model size are not independently controlled. We	625
583	therefore interpret family clustering as a practical	626
584	committee-design result, not as causal proof that	627
585	provider identity alone determines judge behaviour.	628
586	<b>Closed three-label set.</b>	
587	Our task does not exercise ARG3/ARG4 or thematic role labels; results	629
588	should not be interpreted as transferring to full	630
589	PropBank role labelling.	631
590	<b>t3 evaluated on a single setup only.</b>	
591	The t3 disclaimer was evaluated across all five models in	632
592	the retrieved-few-shot setup (the strongest t2 con-	633
593	figuration), plus a three-setup gpt-5-mini ablation.	634
594	We did not re-evaluate t3 on zero-shot or static	635
595	few-shot for the original five models, so the model-	636
596	family effect is established for retrieval-augmented	637
597	prompting and remains to be verified for other se-	638
598	tups.	639
	<b>9 Ethical Considerations</b>	599
	This work uses a lexical-semantic resource and	600
	does not introduce new human-subject data. The	601
	main ethical risk is automation bias: high agree-	602
	ment between LLM judges could be mistaken for	603
	correctness, causing incorrect frames to enter an	604
	annotation resource. Our results argue against fully	605
	automatic acceptance. The judge committee is use-	606
	ful as a precision-oriented filter, but ambiguous or	607
	disagreement-heavy cases should remain subject	608
	to expert review, especially because the resource	609
	supports downstream Turkish NLP systems. AI	610
	assistance was used for editing, and we avoid re-	611
	porting individual annotator behaviour; all analyses	612
	are at the model-output and frame-pattern level.	613
	<b>References</b>	614
	Collin F. Baker, Charles J. Fillmore, and John B. Lowe.	615
	1998. <a href="#">The Berkeley FrameNet project</a> . In <i>36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, Volume 1</i> , pages 86–90, Montreal, Quebec, Canada. Association for Computational Linguistics.	616
	617	618
	619	620
	621	
	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	622
	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	623
	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	624
	Askell, Sandhini Agarwal, Ariel Herbert-Voss,	625
	Gretchen Krueger, Tom Henighan, Rewon Child,	626
	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	627
	Clemens Winter, and 12 others. 2020. Language	628
	models are few-shot learners. In <i>Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS '20</i> , Red Hook, NY, USA. Curran Associates Inc.	629
	630	631
	632	
	Jaime Carbonell and Jade Goldstein. 1998. <a href="#">The use of mmr, diversity-based reranking for reordering documents and producing summaries</a> . In <i>Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '98</i> , page 335–336, New York, NY, USA. Association for Computing Machinery.	633
	634	635
	636	637
	638	639
	Xavier Carreras and Lluís Màrquez. 2005. <a href="#">Introduction to the CoNLL-2005 shared task: Semantic role labeling</a> . In <i>Proceedings of the Ninth Conference on Computational Natural Language Learning (CoNLL-2005)</i> , pages 152–164, Ann Arbor, Michigan. Association for Computational Linguistics.	640
	641	642
	643	644
	645	
	Cheng-Han Chiang and Hung-yi Lee. 2023. <a href="#">Can large language models be an alternative to human evaluations?</a> In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15607–15631, Toronto, Canada. Association for Computational Linguistics.	646
	647	648
	649	650
	651	

652	Jacob Cohen. 1960. <a href="#">A coefficient of agreement for nominal scales</a> . <i>Educational and Psychological Measurement</i> , 20:37 – 46.	Vancouver, Canada. Association for Computational Linguistics.	708 709
653			
654			
655	Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. <a href="#">Unsupervised cross-lingual representation learning at scale</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8440–8451, Online. Association for Computational Linguistics.	Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. <a href="#">Lora: Low-rank adaptation of large language models</a> .	710 711 712 713
656			
657			
658			
659		Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. <a href="#">Dense passage retrieval for open-domain question answering</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 6769–6781, Online. Association for Computational Linguistics.	714 715 716 717 718 719 720
660			
661			
662			
663			
664	Gözde Gül Şahin and Eşref Adalı. 2018. <a href="#">Annotation of semantic roles for the turkish proposition bank</a> . <i>Lang. Resour. Eval.</i> , 52(3):673–706.	J. Richard Landis and Gary G. Koch. 1977. The measurement of observer agreement for categorical data. <i>Biometrics</i> , 33(1):159–174.	721 722 723
665			
666			
667	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. <a href="#">BERT: Pre-training of deep bidirectional transformers for language understanding</a> . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	Beth Levin. 1993. <i>English Verb Classes and Alternations: A Preliminary Investigation</i> . University of Chicago Press, Chicago, IL.	724 725 726
668			
669			
670			
671			
672		Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In <i>Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS ’20</i> , Red Hook, NY, USA. Curran Associates Inc.	727 728 729 730 731 732 733 734 735
673			
674			
675			
676	Rotem Dror, Gili Baumer, Segev Shlomov, and Roi Reichart. 2018. <a href="#">The hitchhiker’s guide to testing statistical significance in natural language processing</a> . In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1383–1392, Melbourne, Australia. Association for Computational Linguistics.	Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2022. <a href="#">What makes good in-context examples for GPT-3?</a> In <i>Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures</i> , pages 100–114, Dublin, Ireland and Online. Association for Computational Linguistics.	736 737 738 739 740 741 742 743
677			
678			
679			
680			
681			
682			
683	Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpaca-farm: a simulation framework for methods that learn from human feedback. In <i>Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS ’23</i> , Red Hook, NY, USA. Curran Associates Inc.	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023. <a href="#">G-eval: NLG evaluation using gpt-4 with better human alignment</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2511–2522, Singapore. Association for Computational Linguistics.	744 745 746 747 748 749 750
684			
685			
686			
687			
688			
689			
690			
691	Bradley Efron and Robert J. Tibshirani. 1993. <i>An Introduction to the Bootstrap</i> . Chapman & Hall/CRC.	Büşra Marşan, Neslihan Kara, Merve Özçelik, Bilge Nas Arıcan, Neslihan Cesur, Aslı Kuzgun, Ezgi Sanıyar, Oğuzhan Kuyrukçu, and Olcay Taner Yıldız. 2021. <a href="#">Building the Turkish FrameNet</a> . In <i>Proceedings of the 11th Global Wordnet Conference</i> , pages 118–125, University of South Africa (UNISA). Global Wordnet Association.	751 752 753 754 755 756 757
692			
693	Charles J. Fillmore. 1976. <a href="#">Frame semantics and the nature of language</a> . <i>Annals of the New York Academy of Sciences</i> , 280(1):20–32.	Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2022. <a href="#">Rethinking the role of demonstrations: What makes in-context learning work?</a> In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 11048–11064, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	758 759 760 761 762 763 764 765
694			
695			
696	Daniel Gildea and Daniel Jurafsky. 2000. <a href="#">Automatic labeling of semantic roles</a> . In <i>Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics</i> , pages 512–520, Hong Kong. Association for Computational Linguistics.		
697			
698			
699			
700			
701	Aslı Göksel and Celia Kerslake. 2005. <i>Turkish: A Comprehensive Grammar</i> . Routledge.		
702			
703	Luheng He, Kenton Lee, Mike Lewis, and Luke Zettlemoyer. 2017. <a href="#">Deep semantic role labeling: What works and what’s next</a> . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 473–483,		
704			
705			
706			
707			

766	Kemal Oflazer. 1993. <a href="#">Two-level description of Turkish morphology</a> . In <i>Sixth Conference of the European Chapter of the Association for Computational Linguistics</i> , Utrecht, The Netherlands. Association for Computational Linguistics.	823
767		824
768		825
769		826
770		827
771	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In <i>Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22</i> , Red Hook, NY, USA. Curran Associates Inc.	828
772		829
773		830
774		831
775		832
776		833
777		834
778		
779		835
780		836
781		837
782	Martha Palmer, Daniel Gildea, and Paul Kingsbury. 2005. <a href="#">The Proposition Bank: An annotated corpus of semantic roles</a> . <i>Computational Linguistics</i> , 31(1):71–106.	838
783		
784		
785		
786	Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, and 25 others. 2025. <a href="#">Qwen2.5 technical report</a> .	839
787		840
788		841
789		842
790		
791		
792	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-BERT: Sentence embeddings using Siamese BERT-networks</a> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.	843
793		844
794		845
795		846
796		847
797		848
798		849
799		850
800	Stephen Robertson and Hugo Zaragoza. 2009. <a href="#">The probabilistic relevance framework: Bm25 and beyond</a> . <i>Found. Trends Inf. Retr.</i> , 3(4):333–389.	851
801		852
802		853
803	Ohad Rubin, Jonathan Herzig, and Jonathan Berant. 2022. <a href="#">Learning to retrieve prompts for in-context learning</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2655–2671, Seattle, United States. Association for Computational Linguistics.	854
804		855
805		856
806		857
807		858
808		859
809		860
810	Peng Shi and Jimmy Lin. 2019. <a href="#">Simple bert models for relation extraction and semantic role labeling</a> . <i>Preprint</i> , arXiv:1904.05255.	861
811		862
812		863
813	Harish Tayyar Madabushi, Taylor Hudson, and Claire Bonial. 2025. <a href="#">Generative FrameNet: Scalable and adaptive frames for interpretable knowledge storage and retrieval for LLMs powered by LLMs</a> . In <i>Proceedings of Bridging Neurons and Symbols for Natural Language Processing and Knowledge Graphs Reasoning @ COLING 2025</i> , pages 107–119, Abu Dhabi, UAE. ELRA and ICCL.	864
814		865
815		
816		
817		
818		
819		
820		
821	Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan	
822	Schalkwyk, Andrew M. Dai, Anja Hauth, Katie Millican, David Silver, Melvin Johnson, Ioannis Antonoglou, Julian Schrittwieser, Amelia Glaese, Jilin Chen, Emily Pitler, Timothy Lillicrap, Angeliki Lazaridou, and 1332 others. 2025. <a href="#">Gemini: A family of highly capable multimodal models</a> .	
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. <a href="#">Llama: Open and efficient foundation language models</a> .	
	Reut Tsarfaty, Djamé Seddah, Sandra Kübler, and Joakim Nivre. 2013. <a href="#">Parsing morphologically rich languages: Introduction to the special issue</a> . <i>Computational Linguistics</i> , 39(1):15–22.	
	Grigorios Tsoumakas and Ioannis Katakis. 2007. <a href="#">Multi-label classification: An overview</a> . <i>International Journal of Data Warehousing and Mining (IJDWM)</i> , 3(3):1–13.	
	Peiyi Wang, Lei Li, Liang Chen, Zefan Cai, Dawei Zhu, Binghuai Lin, Yunbo Cao, Lingpeng Kong, Qi Liu, Tianyu Liu, and Zhifang Sui. 2024. <a href="#">Large language models are not fair evaluators</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 9440–9450, Bangkok, Thailand. Association for Computational Linguistics.	
	Olcay Taner Yıldız, Begüm Avar, and Gökhan Ercan. 2019. <a href="#">An open, extendible, and fast Turkish morphological analyzer</a> . In <i>Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2019)</i> , pages 1364–1372, Varna, Bulgaria. INCOMA Ltd.	
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. <a href="#">Judging llm-as-a-judge with mt-bench and chatbot arena</a> . In <i>Proceedings of the 37th International Conference on Neural Information Processing Systems, NIPS '23</i> , Red Hook, NY, USA. Curran Associates Inc.	

## A GPT-5-mini Ablation

Setup	EM	Ma-F1	ARG0-only	Triadic
zero-shot	0.552	0.663	0.472	0.261
static FS	0.530	0.681	0.457	0.413
retrieved FS	<b>0.584</b>	<b>0.719</b>	0.441	<b>0.457</b>

Table 7: GPT-5-mini under t3 with minimal reasoning and a 4096-token completion budget. All runs parse successfully; EM improves over GPT-4o-mini by 1.2–3.8 points, with much higher ARG0-only accuracy (mean 0.457 vs. 0.184).

Tour	Main instruction state	Intended target	Observed effect
t1	Baseline taxonomy with five frame patterns and a three-part lexical-frame test for ARG2.	Establish whether models can apply the closed PropBank role set without task-specific corrective rules.	Pilot runs exposed severe weakness on ARG0-only unergatives.
t2	Added an explicit unergative anti-pattern block and Turkish morphological cues for causative, passive, reciprocal, reflexive, and light-verb constructions.	Correct the frequent error where animate single-participant predicates were mapped to ARG1 instead of ARG0.	ARG0-only accuracy improved by 21.7 points, but ARG2 misses increased by 30.9 points.
t3	Added a sufficiency disclaimer: morphology can support ARG2 but is not required for ARG2.	Repair the t2 side effect where models treated the morphology list as an exhaustive ARG2 trigger list.	The best Gemini retrieved-few-shot run gained 2.4 EM points; effects did not transfer uniformly to OpenAI or Qwen.

Table 8: Prompt evolution across tours. The table separates the linguistic motivation from the observed empirical side effects.

## B Prompt Templates

The three prompt tours are diffs over a common base. Table 8 summarizes what changed in each tour. Below we reproduce the t3 additions. Full t1/t2/t3 prompt files and evaluation scripts will be released in an anonymized repository.

### Tour 2 — unergative anti-pattern (excerpt).

CRITICAL anti-pattern for unergatives (high model error rate observed in pilot):

- An animate participant performing a sound, behaviour, emission, or expressive act (anırmak, ağlaşmak, havlamak) is UNERGATIVE - its sole participant is ARG0, NEVER ARG1.
- Mental/emotional stative predicates with olmak / +!An- where the experiencer is the intentional/animate locus are unergative {ARG0}, NOT {ARG1}.
- Concrete test: can the verb take a true accusative-marked direct theme (X-i V-mek) within THIS sense?

### Tour 3 — ARG2 sufficiency disclaimer.

Important: the morphological cues above are SUFFICIENT but NOT NECESSARY for ARG2.

- ARG2 can be licensed by class (v) triadic causatives WITHOUT any of the morphological markers above (açıklamak, vermek, anlatmak).
- ARG2 can also be licensed in class (iv) stative two-place relations with light-verb constructions (farklı olmak, tekelinde olmak, dayak yemek).
- Do not treat the absence of morphology as evidence AGAINST ARG2.

### Judge Prompt Template.

You are a Turkish PropBank expert evaluating predicate argument-frame assignment quality.  
 ...  
 Rate the proposed frame on a 1-5 scale:

5 = perfect; 4 = mostly correct (one role over/under); 3 = partially correct; 2 = mostly wrong; 1 = completely wrong.  
 Output ONLY JSON: {"score": int 1-5, "reasoning": string}.

## C Reproducibility

All API calls and judge calls are SHA-256-cached on (model, prompt, sampling parameters). The random seed is 42 for stratified split, sampling, and API decoding where supported. Total compute cost (t1 + t2 + t3 + judges) is approximately \$25 in API spend plus  $\approx 5$  GPU-hours on an RTX 3000 Ada laptop GPU.

**Repository contents.** The release package contains the t1/t2/t3 zero-shot and few-shot prompt files, raw prediction CSVs, judge-score CSVs, aggregate metric CSVs, plotting scripts, and the figures used in this paper. The raw prediction files include the sense identifier, predicate, gold frame, predicted frame, parse status, model confidence, retrieved example identifiers when applicable, and the raw model response. The metric CSVs contain exact match, macro-F1, micro-F1, Hamming loss, per-label precision/recall/F1, per-label Cohen’s  $\kappa$ , and the diagnostic ARG0/ARG2 over- and under-generation rates.

**Dataset access.** The experiments read from the Turkish PropBank spreadsheet and filter to verb-sense rows with at least one filled core-role field. If the underlying resource cannot be redistributed under the final venue policy, the released scripts will still reproduce preprocessing, prompt construction, parsing, metric computation, judge aggregation, and figure generation once the spreadsheet is placed at the expected path.

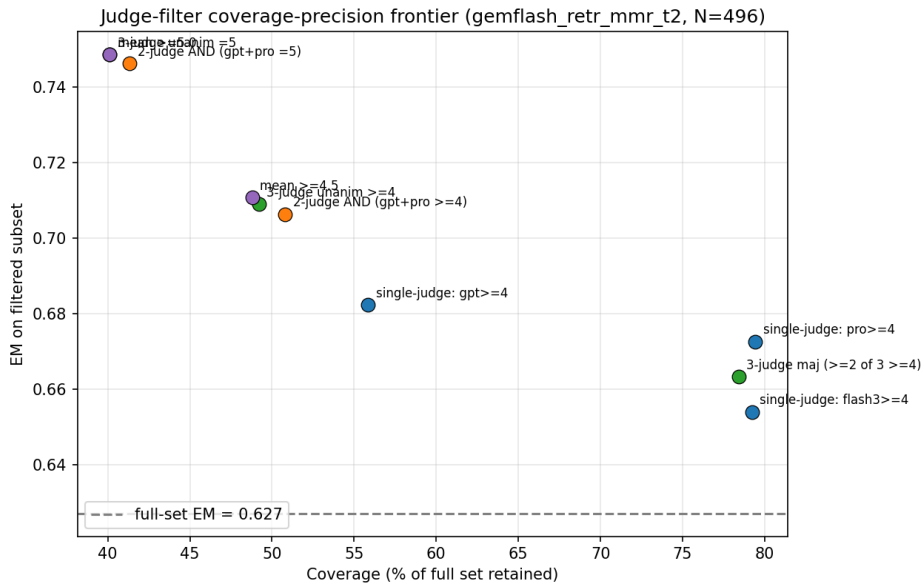


Figure 5: Coverage–precision frontier for judge filters. The best two-judge cross-family rule (GPT-4o-mini and Gemini 2.5 Pro both assign 5) is almost indistinguishable from the three-judge unanimous-5 rule, showing that the extra same-family Gemini judge adds coverage loss with negligible precision gain.

E Qualitative Examples

Phenomenon	Predicate	Gold frame	Model or judge output	Diagnostic value
Prompt repair	<i>gözü açılmak</i>	{ARG0}	t1: {ARG1}; t2: {ARG0}	The t2 unergative/stative guidance moves a single-participant sense from the dominant ARG1 error into the intended ARG0-only class.
ARG2 omission	<i>dikmek</i>	{ARG0,ARG1,ARG2} t2: {ARG0,ARG1}		The model treats the instrument-like complement as an adjunct, a typical post-t2 error where ARG2 is under-generated.
t3 recovery	<i>bel bellemek</i>	{ARG0,ARG1,ARG2} t2: {ARG0,ARG1}; t3: {ARG0,ARG1,ARG2}		The t3 reminder that lexical frames can license ARG2 without derivational morphology lets the model recover an intrinsic instrument role.
Judge disagreement	<i>açılmak</i>	{ARG0,ARG1}	GPT judge: 2; Gemini judges: 5/5 for {ARG0,ARG2}	The prediction is wrong by exact match, but same-family Gemini judges both accept it, illustrating why cross-family review is more useful than adding another same-family judge.

Table 9: Representative examples from cached prediction and judge outputs. They illustrate the main error mechanisms behind the aggregate prompt-iteration and judge-committee results.