

# Parser-Validated Minimal-Pair Preferences for DPO: A Case Study in Legal “Empathy”

Anonymous ACL submission

## Abstract

Aligning LLMs to produce responses perceived as empathetic typically relies on costly human preference data and offers limited insight into which linguistic cues drive those preferences. We study legal question answering and introduce an automatic preference-data pipeline based on parser-validated *syntactic minimal pairs*. Grounded in linguistic accounts of perspective-taking, we generate rule-labeled minimal pairs along five dimensions (pronouns, voice, tense, polite imperatives, evaluative adverbs) and validate the intended contrast with dependency parsing (82.7% success). From 1,785 questions, we produce 7,378 minimal pairs and fine-tune three 7–8B model families (LLaMA-3, Mistral, Gemma) with DPO. Human evaluation (3,309 judgments, 35 raters) prefers DPO over an SFT baseline (68.8% vs. 31.2%,  $p < 0.001$ ,  $h=0.40$ ), robust to length controls. Feature validation shows voice is the dominant above-chance contributor (80%,  $h=0.64$ ), while other edits are register-sensitive. Overall, parser-validated minimal pairs provide an interpretable, scalable route to preference optimization and identify which cues align with human judgments in-domain.

## 1 Introduction

The alignment of Large Language Models (LLMs) with human values represents one of the most critical challenges in AI safety and deployment (Ouyang et al., 2022; Bai et al., 2022a). As LLMs become the backbone of conversational AI systems, ensuring they can engage empathetically with users has emerged as a key alignment objective. However, empathy in AI systems can be understood through two distinct lenses: emotional empathy (feeling what others feel) and linguistic empathy—the grammatical encoding of perspective-taking and cognitive understanding (Kuno and Kaburaki, 1977). While most alignment research focuses on subjective emotional responses requiring human

evaluation, we take a different approach by leveraging linguistically grounded perspective-taking cues that can be operationalized as parser-checkable feature controls.

Reinforcement Learning from Human Feedback (RLHF) has become the dominant approach for alignment (Ziegler et al., 2020; Kaufmann et al., 2023), but suffers from fundamental limitations: expensive human annotations, reward model training instabilities, and lack of interpretability. These challenges are particularly acute for subjective concepts like emotional empathy, where human annotators show low agreement and definitions vary across contexts (Casper et al., 2023).

Direct Preference Optimization (DPO) (Rafailov et al., 2023) simplifies alignment by eliminating the reward modeling step, instead learning directly from preference data. However, this approach inherits the fundamental bottleneck of RLHF: dependence on human-annotated preferences, which poses significant scalability and consistency challenges. We propose an alternative paradigm that reduces reliance on human preference labels for training by automatically generating preference pairs through systematic manipulation of linguistically motivated feature contrasts. Specifically, we leverage Kuno’s theory of linguistic empathy (Kuno and Kaburaki, 1977; Kuno, 1987), which argues that grammatical choices systematically encode speaker perspective and alignment with discourse participants. This yields an explicit inventory of contrasts that can be instantiated and verified at the level of surface realization.

**Scope of “meaning preservation.”** Throughout the paper, when we say that a transformation “preserves meaning,” we mean that it aims to preserve *propositional (truth-conditional) content*. We do not aim to hold constant other aspects (e.g., register or discourse stance), which are instead part of the controlled contrast under study.

083	Our key insight is that certain grammatical pat-	Direct Preference Optimization (DPO) ( <a href="#">Rafailov</a>	131
084	terns, such as voice alternation, pronoun usage, and	<a href="#">et al., 2023</a> ) addresses some technical limitations	132
085	tense selection, can be implemented as controlled	by eliminating the reward model and directly opti-	133
086	<i>syntactic minimal pairs</i> in generated responses. Ac-	mizing policies from preference data. While this	134
087	cordingly, we construct <i>minimal pairs</i> that differ in	simplifies training and improves stability, DPO	135
088	one explicitly targeted feature and validate the real-	inherits RLHF’s core bottleneck: dependence on	136
089	ization of that feature with dependency parsing. We	human-annotated preferences.	137
090	then assign preference labels by the pre-specified		
091	feature target (i.e., which side instantiates the con-	<b>2.2 Automated preference generation</b>	138
092	trast), rather than by an LLM-as-judge pipeline.	Recent work has explored automated preference	139
093	As illustrated in Figure 1, we develop an auto-	generation to overcome the scalability barrier of	140
094	matic preference generation pipeline that:	human annotation. Constitutional AI ( <a href="#">Bai et al.,</a>	141
095	1. creates minimal pairs through rule-based gener-	<a href="#">2022b</a> ) generates feedback using predefined princi-	142
096	ators (e.g., active vs. passive voice),	ples, allowing models to self-critique based on con-	143
097	2. validates them using dependency parsers to	stitutional rules. Self-Instruct ( <a href="#">Wang et al., 2023</a> )	144
098	check that the intended feature contrast is re-	enables models to generate their own training data	145
099	alized, and	through iterative instruction following and refine-	146
100	3. trains via DPO using these rule-defined labels,	ment. UltraFeedback ( <a href="#">Cui et al., 2023</a> ) leverages	147
101	followed by human evaluation to test align-	GPT-4 to create multi-aspect preference judgments	148
102	ment with perceived empathy in-domain.	across various dimensions of response quality.	149
103	Our approach yields efficient learning of feature	However, these approaches share critical limita-	150
104	discrimination on the constructed minimal pairs	tions: (1) they rely on LLM-as-judge paradigms	151
105	across model families. Importantly, high discrim-	where AI systems make subjective quality assess-	152
106	ination accuracy on constructed pairs should be	ments, potentially amplifying biases; (2) they lack	153
107	interpreted as learnability of the targeted contrast,	theoretical grounding explaining <i>why</i> certain re-	154
108	not as a direct measure of empathetic generation	sponses should be preferred; and (3) they typically	155
109	quality. In our legal QA case study, we addition-	do not provide a direct, feature-level linkage be-	156
110	ally conduct human evaluations to quantify which	tween the training label and a verifiable linguis-	157
111	contrasts align with perceived empathy and to pro-	tic contrast, and often require separate validation	158
112	vide an empirical ranking of which linguistically	against human judgments. In contrast, our ap-	159
113	motivated cues transfer in-domain.	proach grounds preferences in formal linguistic	160
114		theory, uses parsers for objective verification, and	161
115	<b>2 Related Work</b>	validates alignment with human judgment through	162
116	<b>2.1 LLM alignment methods</b>	evaluation (80% acceptance for voice transforma-	163
117	Current LLM alignment methods face a fundamen-	tions).	164
118	tal tension between effectiveness and scalability.	<b>2.3 Empathy modeling in NLP</b>	165
119	Reinforcement Learning from Human Feedback	Empathy modeling in NLP has gained attention	166
120	(RLHF) ( <a href="#">Christiano et al., 2017</a> ; <a href="#">Ziegler et al.,</a>	through shared tasks and benchmarks ( <a href="#">Hasan et al.,</a>	167
121	<a href="#">2020</a> ) has become the dominant paradigm, using	<a href="#">2023, 2024b,a</a> ). However, these approaches typi-	168
122	human preferences to train reward models that	cally focus on emotional empathy, relying on sub-	169
123	guide models toward helpful, honest, and harm-	jective human annotations of empathetic responses	170
124	less behavior ( <a href="#">Bai et al., 2022a</a> ). However, RLHF	and facing inter-annotator agreement challenges.	171
125	suffers from well-documented limitations: reward	Existing work largely treats empathy as an emer-	172
126	hacking, distribution shift, prohibitive annotation	gent property to be learned from labels rather than	173
127	costs, and annotator biases ( <a href="#">Casper et al., 2023</a> ).	a phenomenon with systematic linguistic correlates	174
128	These challenges are particularly acute for subjec-	( <a href="#">Kann, 2017</a> ). This conflation of emotional and	175
129	tive qualities like emotional empathy, where def-	linguistic empathy limits both the consistency of	176
130	initions vary across contexts and inter-annotator	training signals and the interpretability of learned	177
	agreement remains low ( <a href="#">Kaufmann et al., 2023</a> ).	behaviors. Our work specifically targets linguistic	178
		empathy, objectively measurable through grammat-	179
		ical analysis.	180

## 2.4 Linguistic perspectives on empathy

While NLP approaches typically treat empathy as a holistic emotional quality requiring human judgment, linguistic research reveals that empathy is systematically expressed through grammatical choices. Kuno and Kaburaki’s seminal work shows that linguistic empathy can be displayed via grammatical structure choice without altering propositional meaning (Kuno and Kaburaki, 1977). This linguistic empathy refers to the speaker’s identification with participants in the described event, as reflected through syntactic choices (Kuno, 1987), namely, a formal property distinct from emotional empathy.

For instance, voice alternation shifts perspective without changing semantic content: structuring a sentence in active voice rather than passive voice (e.g., “John hit Mary” vs. “Mary was hit by John”) changes the recipient of empathy by shifting *the camera angle* (Kuno and Kaburaki, 1977). Active voice centers the agent’s perspective (John), while passive voice centers the patient’s perspective (Mary). These are objectively verifiable grammatical properties, not subjective emotional assessments. Similarly, pronoun usage reveals psychological distance in discourse (Kuno, 1987). Choosing “he” or “she” instead of a proper noun indicates closer alignment with the referent (Kuroshima and Iwata, 2016). Recently, Kann (Kann, 2017) provided empirical validation for these linguistic theories by developing computational methods to measure perspective-marking features, finding significant correlations with psychological empathy scales.

Building on these theoretical foundations, we also incorporate evaluative adverbs based on Ernst’s (Ernst, 2002) classification of agent-oriented adverbs. These adverbs (e.g., ‘carefully,’ ‘thoughtfully,’ ‘thoroughly’) express the speaker’s assessment of how an agent performs an action, indicating cognitive effort and procedural care. While not direct empathy markers, we hypothesize that in professional contexts like legal discourse, explicit markers of deliberate consideration correlate with cognitive perspective-taking.

## 3 Method

### 3.1 Task formulation and approach

We formulate empathetic alignment as a preference learning task where responses exhibiting specific linguistic features associated with cognitive

perspective-taking are preferred over those lacking them. Crucially, we target linguistic empathy (formal grammatical properties that encode perspective) rather than emotional empathy, enabling objective verification enabling parser-based verification of the targeted feature contrasts without requiring human-labeled preference data for training.

Given a legal question  $q$ , we generate preference pairs  $(r_{accepted}, r_{rejected})$  that differ only in a single targeted feature while aiming to preserve propositional (truth-conditional) content. For each feature  $f \in \{\textit{pronoun}, \textit{voice}, \textit{tense}, \textit{politeness}, \textit{evaluative}\}$  and question  $q$ :

- $r_{accepted}^f$ : response exhibiting feature  $f$  (objectively verifiable)
- $r_{rejected}^f$ : response lacking feature  $f$  (objectively verifiable)

We then train models using Direct Preference Optimization (DPO) (Rafailov et al., 2023) to learn preferences for these perspective-encoding features. Given our generated dataset  $\mathcal{D} = \{(x^{(i)}, y_w^{(i)}, y_l^{(i)})\}_{i=1}^N$ , we optimize:

$$\mathcal{L}_{\text{DPO}} = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right] \quad (1)$$

following the standard DPO formulation. The high quality of our automatically generated contrast pairs—labeled by explicit feature targets and validated for their realization—provide clear supervision signals without requiring subjective human judgment or reward models.

### 3.2 Computational rules for linguistic empathy

We implement five computational rules that systematically transform linguistic features associated with perspective-taking and speaker stance. Each rule operationalizes a formal linguistic property that can be objectively verified through grammatical analysis:

**Rule 1: Pronoun** Replace full noun phrases with pronouns to reduce referential distance (Kuno, 1987). The rule identifies main subject positions and substitutes appropriate pronouns while preserving grammatical correctness.

- ✓ “It must consider all evidence.”
- ✗ “The tribunal must consider all evidence.”

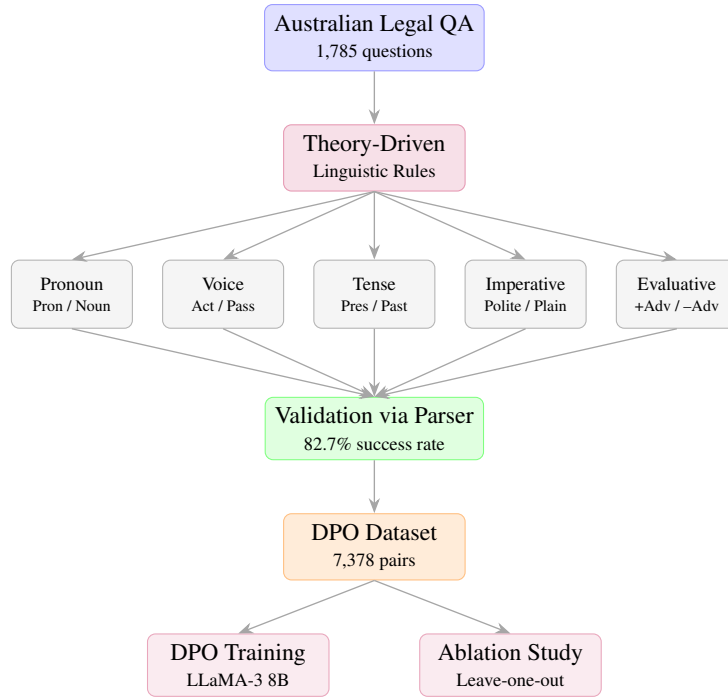


Figure 1: Automatic preference generation pipeline. Legal questions (1,785) are processed through five rule-based generators to create minimal pairs, validated by dependency parsers (82.7% success rate), producing 7,378 training pairs for DPO training.

**Rule 2: Voice** Convert between active and passive voice to shift perspective focus (Kuno and Kaburaki, 1977). The rule identifies transitive verbs and applies syntactic transformations—a formal grammatical operation verifiable through dependency parsing.

- ✓ “The court reviews each application carefully.”
- ✗ “Each application is reviewed carefully by the court.”

**Rule 3: Tense** Transform between present and past tense to manipulate temporal distance (Kann, 2017). The rule ensures consistent tense mapping throughout the response.

- ✓ “The department processes requests promptly.”
- ✗ “The department processed requests promptly.”

**Rule 4: Politeness imperatives** Add or remove politeness markers in imperatives (Brown and Levinson, 1987). The rule detects imperative constructions and systematically modifies politeness markers, a binary feature verifiable through syntactic analysis.

- ✓ “Please submit your documents by Friday.”
- ✗ “Submit your documents by Friday.”

**Rule 5: Evaluative adverbs** Insert evaluative adverbs from a predefined set indicating cognitive effort (Ernst, 2002). Presence/absence of specific

lexical items from our defined set—objectively verifiable.

- ✓ “The tribunal carefully examines each case.”
- ✗ “The tribunal examines each case.”

### 3.3 Automatic preference generation pipeline

As illustrated in Figure 1, our pipeline implements these computational rules through two core components that enable fully automatic preference data generation with objective verification:

**Rule-based generators** For each computational rule, we develop specialized generators that produce minimal pairs differing only in the target linguistic feature. We use Claude-3-Haiku (Anthropic, 2024) to generate base responses, then apply deterministic rule-based transformations to create *accepted* and *rejected* variants. Crucially, the preference label is determined by the transformation rule—not by LLM judgment—ensuring rule-defined labels (i.e., the preferred side is determined by the targeted feature), rather than by an LLM judge. Detailed prompts are provided in the Appendix.

**Dependency validators** We implement rule-specific parsers using spaCy (Honribal et al., 2020) to validate correct feature implementation through formal grammatical analysis. Pronoun validators

prevent ambiguity by preserving relative clause boundaries, while voice parsers detect auxiliary verbs and past participles to ensure grammatically correct active-passive conversions. Tense validation maintains consistency across finite verbs but preserves modals and infinitives where tense marking doesn't apply. For imperatives, our parser distinguishes genuine commands from declarative sentences, and for evaluative adverbs, we follow Ernst's classification (Ernst, 2009) to verify proper verb modification. This comprehensive validation achieves 82.7% success rate, with failed pairs re-generated up to three times using targeted feedback. Throughout, formatting and punctuation remain constant to prevent spurious learning signals. The validation process ensures that each pair instantiates the intended single-feature contrast under parser-based checks, with labels defined by the targeted feature.

## 4 Experiments

### 4.1 Dataset generation and processing

We apply our pipeline to 1,785 unique questions from the Open Australian Legal QA dataset (Butler, 2023). Due to varying complexity of linguistic transformations, we observe differential success rates across rules. From initial generation attempts, we successfully create 7,378 validated preference pairs: imperative (1,641 pairs), evaluative adverbs (1,557 pairs), tense (1,459 pairs), voice (1,445 pairs), and pronoun substitution (1,276 pairs). The higher success rates for imperative and evaluative rules reflect their straightforward transformations, while pronoun substitution proves most challenging due to complex co-reference resolution requirements.

To validate cross-architecture robustness, we train three model architectures: Meta-Llama-3-8B-Instruct (Dubey et al., 2024; Meta, 2024), Mistral-7B-Instruct-v0.2 (Jiang et al., 2023; Mistral AI, 2023), and Gemma-7B-IT (Gemma Team, 2024; Google, 2024). For ablation studies requiring controlled comparisons, we use a subset of 519 questions where all five rules successfully generated valid pairs, yielding 2,595 pairs ( $519 \times 5$  rules). Each leave-one-out condition contains 2,076 pairs ( $519 \times 4$  rules).

### 4.2 Training configuration

We fine-tune three model architectures (LLaMA-3 8B Instruct, Mistral-7B-Instruct, and Gemma-7B-

Model	Steps to 100%	Final Margin
Mistral-7B	50	15.26
LLaMA-3 8B	75	16.22
Gemma-7B-IT	100	13.95

Table 1: Cross-architecture convergence. All models achieve perfect preference discrimination within 100 steps.

IT) using the TRL library's DPO implementation. Our configuration employs LoRA (Hu et al., 2021) ( $r=16$ ,  $\alpha=32$ ,  $\text{dropout}=0.05$ ) for parameter-efficient training with batch size 16 ( $4 \times 4$  gradient accumulation), learning rate  $5e-5$  with cosine scheduling, and  $\beta=0.1$ . We implement early stopping with patience of 2 evaluations (every 10 steps) to prevent overfitting. All experiments run on a single NVIDIA A100 80GB GPU. Training dynamics are shown in Appendix F.

## 5 Results

### 5.1 Training and Cross-Architecture Validation

We first verify that our linguistically-grounded minimal pairs provide learnable supervision signals across different model architectures. Training three models (LLaMA-3 8B, Mistral-7B, and Gemma-7B-IT) on 7,378 pairs, all achieved 100% preference discrimination within 50–100 steps with consistent margins (Table 1). This cross-architecture consistency confirms that our linguistic rules provide clear, model-agnostic signals rather than exploiting architecture-specific biases.

The 100% accuracy requires careful interpretation. This reflects *feature discrimination*—the ability to distinguish presence from absence of specific grammatical markers—not empathetic generation quality. The question of whether these learned preferences translate to human-perceived empathy requires direct evaluation.

### 5.2 Human Evaluation: DPO vs. SFT

To answer this question, we conducted an IRB-approved study comparing DPO-aligned models against a supervised fine-tuning (SFT) baseline trained on the same accepted responses. Qualified MTurk workers ( $\geq 97\%$  approval,  $\geq 1000$  HITs, US location) judged which of two responses to 50 legal questions was "more empathetic." Each item received 16 judgments, yielding 784 valid responses after quality control.

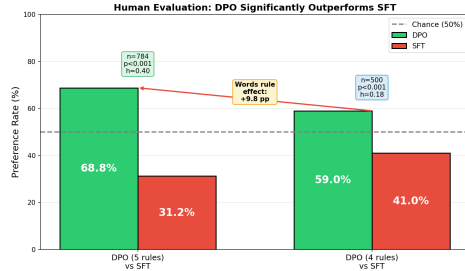


Figure 2: Human evaluation: DPO significantly outperforms SFT (68.8% vs. 31.2%,  $p < 0.001$ , Cohen’s  $h=0.40$ ).

The results were clear: DPO responses were preferred in 68.8% of judgments versus 31.2% for SFT ( $p < 0.001$ , Cohen’s  $h=0.40$ ; Figure 2). This medium-sized effect demonstrates that preference learning on linguistic features produces responses humans perceive as meaningfully more empathetic than supervised learning alone.

### 5.3 Controlling for Length Bias

One potential confound is response length: if DPO models simply generate longer responses, annotators might prefer them regardless of empathy. We tested this through correlation analysis and length-matched comparisons.

The correlation between length difference and DPO preference was weak and non-significant ( $r=0.178$ ,  $p=0.222$ ; Figure 3, left). More directly, in pairs where responses differed by only 0–5 words, DPO was still preferred in 58.2% of cases (Figure 3, right). Effect decomposition confirmed that quality contributes +10.5 percentage points while length contributes only  $-0.7$ pp. The DPO advantage reflects response quality, not verbosity.

### 5.4 Feature-Level Validation

Having established that DPO outperforms SFT overall, we next ask: which of our five linguistic rules actually drive perceived empathy? Using the Communicative Empathy framework (Buechel et al., 2018), nine raters evaluated 75 accepted/rejected training pairs across three dimensions (Understanding, Supportiveness, Helpfulness), yielding 2,025 judgments.

The results revealed striking asymmetry (Figure 4). Voice transformations achieved 80.0% acceptance ( $h=0.64$ ,  $p < 0.001$ ), consistent across all three empathy dimensions. The remaining rules fell at or below chance: Imperative (36.8%), Evaluative adverbs (19.8%), Tense (17.8%), Pronoun

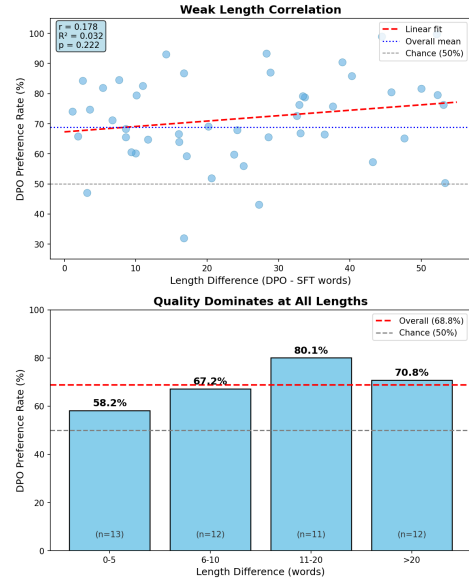


Figure 3: Length bias analysis. Top: Weak correlation ( $r=0.178$ ,  $p=0.222$ ). Bottom: DPO preferred even in length-matched pairs (58.2% at 0–5 word gap).

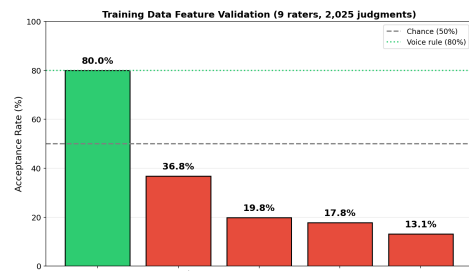


Figure 4: Feature validation ( $n=2,025$ ). Only voice achieved above-chance acceptance (80%); other features conflict with legal register norms.

(13.1%). Voice emerges as the dominant feature for perceived empathy in legal domain.

### 5.5 Feature Ablation

The low individual acceptance rates raise a puzzle: if four of five features fail validation, why does the full model outperform SFT so substantially? We hypothesized that features might contribute synergistically even when individually weak.

To test this, we trained a model without the evaluative adverbs rule and repeated our human evaluation ( $n=500$ ). Removing this rule reduced DPO preference from 68.8% to 59.0% ( $-9.8$ pp,  $p=0.033$ , Cohen’s  $h=0.18$ ; see Appendix E for details). Features that fail isolated validation can still contribute meaningfully in combination.

465	<b>5.6 Training Dynamics Ablation</b>	
466	Finally, we examined how individual rules affect	
467	training dynamics through leave-one-out ablation	
468	on a balanced 519-question subset.	
469	Removing voice produced a surprising result:	
470	margins <i>increased</i> $2.4\times$ ( $21.32 \rightarrow 51.43$ ) while	
471	maintaining perfect accuracy. This contrasts	
472	sharply with voice’s strong human validation	
473	(80%). Removing tense had the opposite effect: ac-	
474	curacy degraded to 99.5% with the lowest margins	
475	(18.57). Other rules (pronoun, politeness, evalua-	
476	tive) showed minimal impact on training dynamics,	
477	with margins ranging from 19.47 to 20.49.	
478	This apparent paradox (i.e., voice is the strongest	
479	human-validated signal yet yields lower training	
480	margins when included) requires interpretation.	
481	This contrast suggests that optimization difficulty	
482	(as reflected by training margins) and human-	
483	perceived desirability need not align, especially	
484	when a feature is pragmatically appropriate only	
485	in a subset of legal contexts. In other words, mar-	
486	gins primarily reflect how separable the constructed	
487	preference signal is under the model and dataset,	
488	whereas human judgments reflect context-sensitive	
489	pragmatic appropriateness. We revisit this discrep-	
490	ancy in the Discussion.	
491	<b>6 Discussion</b>	
492	<b>6.1 From Black-Box Preferences to</b>	
493	<b>Interpretable Alignment</b>	
494	Traditional RLHF operates as a black box: annota-	
495	tors judge which response is “better,” but the model	
496	cannot identify <i>which features</i> drive this preference.	
497	When thousands of subjective judgments coalesce	
498	into model behavior, the learned patterns resist in-	
499	terpretation.	
500	Our approach inverts this paradigm. By decom-	
501	posing “empathy” into five linguistically grounded	
502	computational rules, we create preferences with	
503	rule-defined labels tied to an explicit, parser-	
504	checkable contrast. When our model prefers re-	
505	sponse A over B, we can trace this to a specific	
506	grammatical transformation: A uses active voice	
507	while B uses passive, or A includes evaluative ad-	
508	verbs while B omits them.	
509	This transparency enabled discoveries invisible	
510	to conventional methods: the dominance of voice	
511	(80% human acceptance), the domain-specificity	
512	of other features, and the synergistic effects that	
513	emerge only in combination.	
	<b>6.2 Comparison with Automated Preference</b>	<b>514</b>
	<b>Methods</b>	<b>515</b>
	Recent work has explored automated preference	516
	generation to overcome annotation costs. Constitu-	517
	tional AI (Bai et al., 2022b), Self-Instruct (Wang	518
	et al., 2023), and UltraFeedback (Cui et al., 2023)	519
	use LLM-as-judge paradigms where models evalu-	520
	ate response quality. While scalable, these ap-	521
	proaches share critical limitations our method ad-	522
	dresses.	523
	First, existing methods lack theoretical ground-	524
	ing explaining <i>why</i> certain responses should be	525
	preferred. Our preferences derive from formal lin-	526
	guistic theories of perspective-taking, where each	527
	grammatical choice encodes documented effects	528
	on speaker-referent alignment.	529
	Second, LLM-generated preferences cannot be	530
	objectively verified. Our dependency parsers con-	531
	firm the realization of the targeted contrast, provid-	532
	ing verification independent of any LLM-as-judge	533
	preference labeling.	534
	Third, prior work assumes LLM preferences ap-	535
	proximate human perception without validation.	536
	We validate directly: voice achieves 80% human	537
	acceptance ( $h=0.64$ ), confirming theoretical pre-	538
	dictions, while also revealing that other features	539
	require domain adaptation.	540
	<b>6.3 Voice as the Dominant Signal</b>	<b>541</b>
	Human evaluation revealed striking asymmetry:	542
	voice achieved 80% acceptance while other fea-	543
	tures fell below chance in legal domain. This val-	544
	idates Kuno and Kaburaki’s (1977) central claim	545
	that voice alternation encodes speaker perspective.	546
	Active voice positions the speaker closer to the	547
	agent, while passive voice distances the speaker.	548
	This “camera angle” shift is perceived as empathy.	549
	Yet removing voice during training <i>increased</i>	550
	margins $2.4\times$ . This apparent paradox reflects	551
	domain-specific pragmatics: in legal writing, pas-	552
	sive voice serves objectivity (“The evidence was	553
	examined”), agent backgrounding (“Damages were	554
	awarded”), and precedent citation (“This principle	555
	was established in Smith v. Jones”). Our training	556
	data thus contained inconsistent signals. The 80%	557
	acceptance reflects the general preference for ac-	558
	tive voice; the 20% represents legal contexts where	559
	passive is appropriate.	560
	This finding has implications for alignment: uni-	561
	versal linguistic preferences require domain adap-	562
	tation. Future systems could implement context-	563

564	aware weighting based on register and communicative goals.	612
565		613
566	<b>6.4 Synergistic Feature Effects</b>	614
567	Evaluative adverbs achieved only 19.8% individual acceptance, yet removing this rule reduced overall preference from 68.8% to 59.0%. How can a weak feature contribute to success?	615
568		616
569		617
570		618
571	We hypothesize that linguistic empathy is perceived holistically rather than through individual markers. A response combining active voice, present tense, and evaluative adverbs creates a cumulative impression of speaker engagement, even if each feature alone is insufficient. This has methodological implications: isolated feature validation may underestimate contributions that emerge only in combination. Our ablation design captures these interaction effects.	619
572		620
573		621
574		622
575		623
576		624
577		625
578		626
579		627
580		
581	<b>7 Limitations and Future Work</b>	628
582	Several limitations warrant consideration. First, four of five rules conflict with legal register norms, and cross-domain evaluation is ongoing. Second, our single-sentence minimal pairs cannot capture discourse-level empathy patterns. Third, we operationalized only five computational rules grounded in linguistic empathy theory; additional phenomena ( deixis, evidentiality, modal verbs) remain unexplored. Fourth, our English-only implementation requires cross-linguistic validation. Finally, we measured perceived empathy rather than downstream outcomes such as user satisfaction or therapeutic alliance.	629
583		630
584		631
585		632
586		633
587		634
588		635
589		636
590		637
591		638
592		639
593		
594		
595	<b>8 Conclusion</b>	640
596	We introduced a parser-validated pipeline for synthesizing preference data as syntactic minimal pairs and applying DPO with substantially reduced reliance on human-labeled preference data for training. In a case study on legal question answering, we operationalized five linguistically motivated feature contrasts, validated transformations with dependency parsing, and produced thousands of preference pairs for training.	641
597		642
598		643
599		644
600		645
601		646
602		647
603		648
604		649
605		650
606		651
607		652
608		653
609		654
610		655
611		656
		657
		658
		659
		660
		661
		662
		663
		664
		665
		666
		667
		668
		669
		670
		671
		672
		673
		674
		675
		676
		677
		678
		679
		680
		681
		682
		683
		684
		685
		686
		687
		688
		689
		690
		691
		692
		693
		694
		695
		696
		697
		698
		699
		700
		701
		702
		703
		704
		705
		706
		707
		708
		709
		710
		711
		712
		713
		714
		715
		716
		717
		718
		719
		720
		721
		722
		723
		724
		725
		726
		727
		728
		729
		730
		731
		732
		733
		734
		735
		736
		737
		738
		739
		740
		741
		742
		743
		744
		745
		746
		747
		748
		749
		750
		751
		752
		753
		754
		755
		756
		757
		758
		759
		760
		761
		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779
		780
		781
		782
		783
		784
		785
		786
		787
		788
		789
		790
		791
		792
		793
		794
		795
		796
		797
		798
		799
		800
		801
		802
		803
		804
		805
		806
		807
		808
		809
		810
		811
		812
		813
		814
		815
		816
		817
		818
		819
		820
		821
		822
		823
		824
		825
		826
		827
		828
		829
		830
		831
		832
		833
		834
		835
		836
		837
		838
		839
		840
		841
		842
		843
		844
		845
		846
		847
		848
		849
		850
		851
		852
		853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
		895
		896
		897
		898
		899
		900

659	Penelope Brown and Stephen C Levinson. 1987. <i>Politeness: Some universals in language usage</i> . Cambridge University Press, Cambridge.	pages 536–541, Toronto, Canada. Association for Computational Linguistics.	713 714
662	Sven Buechel, Anneke Buffone, Barry Slaff, Lyle Ungar, and João Sedoc. 2018. Modeling empathy and distress in reaction to news stories. <i>arXiv preprint arXiv:1808.10399</i> .		715
666	Umar Butler. 2023. <a href="#">Open australian legal qa</a> .		716
667	Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, and 1 others. 2023. Open problems and fundamental limitations of reinforcement learning from human feedback. <i>arXiv preprint arXiv:2307.15217</i> .		717
674	Paul F Christiano, Jan Leike, Tom Brown, Miljan Martić, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In <i>Advances in Neural Information Processing Systems</i> , volume 30.		718
679	Ganqu Cui, Lifan Yuan, Ning Ding, Guanming Yao, Wei Zhu, Yuan Ni, Guotong Xie, Zhiyuan Liu, and Maosong Sun. 2023. Ultrafeedback: Boosting language models with high-quality feedback. <i>arXiv preprint arXiv:2310.01377</i> .		719
684	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. <i>arXiv e-prints</i> , pages arXiv–2407.		720
689	Thomas Ernst. 2002. <i>The Syntax of Adjuncts</i> . Cambridge University Press, Cambridge.		721
691	Thomas Ernst. 2009. Speaker-oriented adverbs. <i>Natural Language &amp; Linguistic Theory</i> , 27(3):497–544.		722
693	Gemma Team. 2024. <a href="#">Gemma: Open models based on gemini research and technology</a> . <i>arXiv preprint arXiv:2403.08295</i> .		723
696	Google. 2024. google/gemma-7b-it: Model card. <a href="https://huggingface.co/google/gemma-7b-it">https://huggingface.co/google/gemma-7b-it</a> . Accessed 2026-01-03.		724
699	Md Rakibul Hasan, Md Zakir Hossain, Tom Gedeon, and Shafin Rahman. 2024a. <a href="#">LLM-GEM: Large Language Model-Guided Prediction of People’s Empathy Levels towards Newspaper Article</a> . In <i>Findings of the Association for Computational Linguistics: EACL 2024</i> , pages 2215–2231, St. Julian’s, Malta. Association for Computational Linguistics.		725
706	Md Rakibul Hasan, Md Zakir Hossain, Tom Gedeon, Susannah Soon, and Shafin Rahman. 2023. <a href="#">Curtin OCAI at WASSA 2023 Empathy, Emotion and Personality Shared Task: Demographic-Aware Prediction Using Multiple Transformers</a> . In <i>Proceedings of the 13th Workshop on Computational Approaches to Subjectivity, Sentiment, &amp; Social Media Analysis</i> ,		726
707			727
708			728
709			729
710			730
711			731
712			732
			733
			734
			735
			736
			737
			738
			739
			740
			741
			742
			743
			744
			745
			746
			747
			748
			749
			750
			751
			752
			753
			754
			755
			756
			757
			758
			759
			760
			761
			762
			763
			764
			765

766 Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano  
767 Ermon, Christopher D Manning, and Chelsea Finn.  
768 2023. Direct preference optimization: Your language  
769 model is secretly a reward model. In *Advances in*  
770 *Neural Information Processing Systems*, volume 36.

771 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Al-  
772 isa Liu, Noah A Smith, Daniel Khashabi, and Han-  
773 naneh Hajishirzi. 2023. Self-instruct: Aligning lan-  
774 guage models with self-generated instructions. *arXiv*  
775 *preprint arXiv:2212.10560*.

776 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B  
777 Brown, Alec Radford, Dario Amodei, Paul Chris-  
778 tiano, and Geoffrey Irving. 2020. Fine-tuning lan-  
779 guage models from human preferences. *arXiv*  
780 *preprint arXiv:1909.08593*.

781 Appendices A–D detail our automatic preference  
782 generation pipeline; Appendix E provides human  
783 evaluation methodology; and Appendix F provides  
784 training details.

## A Complete Generator Prompts

We present the full implementation of our rule-specific generators, including all prompts, transformation logic, and fallback mechanisms.

### A.1 Pronoun Substitution Generator

#### Stage 1: Base Response Generation

**System:** You answer legal questions with single sentences using full noun phrases.

**User Prompt:**

Answer this question with ONE sentence using FULL NOUN PHRASES.

**CRITICAL:** Do NOT use any pronouns (it, they, he, she, this, that).

Always use the complete noun phrase.

**Examples:**

- ✓ “The immigration tribunal must consider all evidence that the applicant submits.”
- ✗ “It must consider...” (NO! Use “The tribunal must consider...”)

Question: {question}

#### Stage 2: Pronoun Transformation

**Transformation Rules:**

- Replace ONLY the MAIN SUBJECT (before relative clause)
- KEEP relative clauses (who/that/which...) intact
- “The tribunal” → “It”
- “The person/applicant” → “They”
- “Individuals” → “They”

**Example:**

“The tribunal that reviews immigration cases must consider all evidence.”  
→ “It that reviews immigration cases must consider all evidence.”

#### Stage 3: Fallback Patterns

**Regex patterns (if LLM fails):**

*With relative clauses:*

- The (tribunal|court|board) (who|that|which)... → It
- The (applicant|person) + (who|that|which)... → They

*Simple patterns:*

- The tribunal → It
- The court → It
- The applicant → They

## A.2 Voice Transformation Generator

### Stage 1: Active Voice Generation

**System:** You answer legal questions with single sentences in active voice.

**Requirements:**

- Subject performs action (S + V + O structure)
- Include clear object for passive transformation
- Avoid “is/was/been + past participle”

**Good examples:**

- ✓ “The tribunal must review all submitted documents.”
- ✗ “Documents are reviewed by the tribunal.” (passive)

### Stage 2: Passive Transformation

**Transformation patterns:**

- Basic: “The court decides cases” → “Cases are decided by the court”
- Modal: “The court must decide” → “It must be decided by the court”
- Continuous: “The court is reviewing” → “It is being reviewed by the court”

## A.3 Tense Mapping Generator

### Present Tense Generation

**System:** Describe what tribunals ACTUALLY DO in present tense.

**Requirements:**

- Use present ACTION verbs (reviews, processes, decides)
- AVOID modals (can, must, may) - describe actual actions

**Tense mappings:**

is→was, are→were, has→had, reviews→reviewed, decides→decided

## A.4 Politeness Insertion Generator

### Imperative Generation & Transformation

**Stage 1:** Generate plain imperative (no politeness markers)  
Example: “Submit your appeal to the tribunal within 30 days.”

**Stage 2:** Add politeness marker

- Formal contexts: Add “Kindly”
- General: Add “Please”

Example: “Submit your appeal” → “Please submit your appeal”

## A.5 Evaluative Adverb Generator

### Plain Statement & Adverb Insertion

**Stage 1:** Generate without evaluative adverbs  
Example: “The tribunal considers all submitted evidence.”

**Stage 2:** Insert appropriate adverb

**Approved adverbs (agent-oriented, following Ernst’s classification):**

- Manner: carefully, thoughtfully, thoroughly, deliberately
- Procedural: properly, appropriately, correctly, duly

- Evaluative: reasonably, sensibly, fairly, wisely
- Professional: professionally, competently, efficiently, effectively

Example: “The tribunal considers” → “The tribunal carefully considers”

## B Validation Algorithms

### B.1 Base Validation Framework

#### Algorithm 1 Base Rule Validation

**Require:** Accepted  $R_a$ , Rejected  $R_r$ , Rule  $T$   
**Ensure:** Valid (bool), Diagnostics  
1:  $docA \leftarrow \text{SpacyParse}(R_a)$   
2:  $docR \leftarrow \text{SpacyParse}(R_r)$   
3: **if**  $\text{SentCount}(docA) \neq 1$  OR  $\text{SentCount}(docR) \neq 1$  **then**  
4:     **return** False, “Not single sentence”  
5: **end if**  
6: **if**  $\text{WordCount}(docA) < 10$  OR  $\text{WordCount}(docA) > 40$  **then**  
7:     **return** False, “Length violation”  
8: **end if**  
9:  $result \leftarrow \text{ValidateRule}(docA, docR, T)$   
10: **return**  $result$

### B.2 Pronoun Validation

#### Algorithm 2 Pronoun Rule Validation

**Require:** Parsed  $docA, docR$   
**Ensure:** Validation result  
1:  $hasRelClause \leftarrow \text{CheckRelativeClauses}(docA, docR)$   
2:  $subjA \leftarrow \text{ExtractMainSubjects}(docA)$   
3:  $subjR \leftarrow \text{ExtractMainSubjects}(docR)$   
4:  $pronouns \leftarrow \text{CountPronouns}(subjA)$   
5:  $nouns \leftarrow \text{CountNouns}(subjR)$   
6: **if**  $pronouns = 0$  **then**  
7:     **return** False, “No pronouns in accepted”  
8: **end if**  
9: **if**  $nouns = 0$  **then**  
10:     **return** False, “No nouns in rejected”  
11: **end if**  
12: **if**  $hasRelClause$  AND NOT  $\text{PreservedRelClause}(docA, docR)$  **then**  
13:     **return** False, “Relative clause altered”  
14: **end if**  
15: **return** True, {pronouns:  $pronouns$ , nouns:  $nouns$ }

### B.3 Other Validation Algorithms

[Similar compact format for Voice, Tense, Politeness, and Evaluative validations]

## C Results and Analysis

### C.1 Generation Success Rates

### C.2 Failure Analysis

#### Common Failure Patterns

##### Pronoun (29.1% failure rate):

- Complex relative clauses with multiple referents (35%)
- Nested sentence structures (28%)

Table 2: Detailed generation statistics (7,378/8,925 total pairs = 82.7%)

Rule	Pairs	Success	1st try	Avg.
Politeness	1,641	91.2%	87.2%	1.2
Evaluative	1,557	86.5%	79.3%	1.4
Tense	1,459	81.1%	69.1%	1.7
Voice	1,445	80.3%	68.4%	1.8
Pronoun	1,276	70.9%	55.7%	2.1

- Ambiguous antecedents (22%)
- Coordination issues (15%)

#### Voice (19.7% failure rate):

- “Be + noun” copular structures from listing questions (40%)
- “Who + be” questions producing identity statements (17%)
- “Who represented” questions with awkward passivization (11%)
- Already passive constructions (9%)
- Other intransitive/stative verbs (23%)

#### Tense (18.9% failure rate):

- Mixed tense requirements in legal contexts (38%)
- Perfect aspect preservation issues (29%)
- Modal verb interference (24%)
- Conditional statements (9%)

#### Politeness (8.8% failure rate):

- Non-imperative structures (52%)
- Embedded commands (31%)
- Indirect speech acts (17%)

#### Evaluative (13.5% failure rate):

- Adverb placement conflicts (44%)
- Multiple verb phrases (35%)
- Auxiliary verb complications (21%)

### C.3 Detailed Voice Failure Analysis

We analyzed 35 voice transformation validation failures in detail. The primary cause is that certain question types inherently elicit responses with copular or intransitive verbs, which cannot undergo active-passive transformation.

Table 3: Voice validation failure breakdown by question pattern (n=35)

Question Pattern	Count	%
“What was/were the X” (listing/be+noun)	14	40%
“Who was/were” (identity statement)	6	17%
“Who represented” (awkward passive)	4	11%
“Can X be” (already passive)	3	9%
Other (intransitive/stative)	8	23%

## D Examples

### D.1 Successful Transformations

#### Pronoun - Natural Relative Clause

**Q:** “What about individuals who have already filed an appeal?”

**Base:** “Individuals who have already filed an appeal must wait for the tribunal’s decision before taking further action.”

**Transformed:** “They who have already filed an appeal must wait for the tribunal’s decision before taking further action.”

✓ Relative clause preserved, only main subject replaced

#### Voice - Modal Handling

**Q:** “Must the court consider new evidence?”

**Active:** “The court must carefully examine all new evidence.”

**Passive:** “All new evidence must be carefully examined by the court.”

✓ Correct passive with modal preserved

### D.2 Failure Cases

#### Pronoun - Coordination

**Q:** “Can the person who filed the complaint and the witness who supports it both appear?”

**Base:** “The person who filed the complaint and the witness who supports the complaint may both appear.”

**Failed:** “They who filed the complaint and they who supports it may both appear.”

✗ Double pronoun creates ambiguity; coordinated subjects require special handling

#### Voice - “What were the X” (Be + Noun)

**Q:** “What were the two issues in the appeal case of Millar v Commissioner of Taxation [2016] FCAFC 94?”

**Expected answer structure:** “The two issues were X and Y.” (copular verb)

**Problem:** The verb “were” links subject to noun complement—no transitive action verb exists to passivize.

✗ Copular verbs cannot undergo active-passive transformation

#### Voice - “Who were the X” (Identity Statement)

**Q:** “Who were the parties involved in the case Kirby v Centro Properties Limited (No 2) [2011] FCA 1144?”

**Expected answer:** “The parties were X and Y.” (be + noun predicate)

**Problem:** Identity statements use “be” as a linking verb, not a transitive action verb.

✗ No action verb to convert to passive voice

#### Voice - “Who represented” (Awkward Passive)

**Q:** “Who represented the applicant and the respondent in the case of CCA Beverages v Commissioner of Taxation [1995] FCA 980?”

**Active:** “Mr. Smith represented the applicant and Ms. Jones represented the respondent.”

**Attempted passive:** “The applicant was represented by Mr. Smith and the respondent was represented by Ms. Jones.”

**Problem:** While grammatically possible, passive construction obscures the answer to “who” questions.

✗ Semantic mismatch: passive form doesn’t naturally answer “who” questions

#### Voice - Intransitive

**Q:** “When does the tribunal meet?”

**Active:** “The tribunal meets every Tuesday.”

**Failed attempts:**

- “Every Tuesday is met by the tribunal” (nonsensical)
- “Meetings are held by the tribunal” (verb changed)

✗ No valid passive for intransitive “meet”

#### Voice - Already Passive Construction

**Q:** “In the context of New South Wales law, can a company be bound by a contract independently of the Corporations Act?”

**Problem:** Question already contains passive voice (“be bound”). Generating an active-voice answer and converting it back creates circular transformations.

✗ Question structure primes passive response

#### Tense - “When did” Questions

**Q:** “When did the Health Insurance Act amendment commence?”

**Expected answer:** “The amendment commenced on January 1, 2020.” (past tense)

**Problem:** The question already uses past tense (“did”), naturally eliciting a past-tense response. Present-to-past transformation becomes meaningless.

✗ Question tense determines answer tense

#### Evaluative - Definition Questions

**Q:** “What is the definition of an RSE as per the Financial Sector (Collection of Data) reporting standard?”

**Expected answer:** “An RSE is defined as...” (stative/definitional)

**Problem:** Definition responses lack action verbs that evaluative adverbs can modify. Adding “carefully” or “thoroughly” to “is defined as” is semantically odd.

✗ Stative verbs incompatible with manner adverbs

835	<b>E Human Evaluation Details</b>	<b>Sample Size</b>	879
836	This appendix provides comprehensive details of	• Questions: 50	880
837	our human evaluation studies, including platform	• Judgments per question: 16	881
838	specifications, participant qualifications, quality	• Total judgments: 800 (784 after quality control)	882
839	control measures, and complete study designs.		883
840	<b>E.1 Platform and Recruitment</b>	<b>Results</b> DPO responses preferred 68.8% vs. SFT	884
841	<b>Platform</b> All studies were conducted using	31.2% (exact binomial $p < 0.001$ , Cohen's	885
842	Google Forms distributed via Amazon Mechanical	$h=0.40$ ).	886
843	Turk (MTurk). Figure 5 shows our HIT configura-	<b>E.4 Study 2: Construct Validity</b>	887
844	tion.	<b>Objective</b> Validate that our theoretical constructs	888
845	<b>Compensation</b> Participants received \$5.50 per	(minimal pair transformations) translate to perceived	889
846	HIT (estimated 20 minutes completion time), cor-	empathy differences across multiple dimen-	890
847	responding to an hourly rate of \$16.50.	sions.	891
848	<b>Participant Qualifications</b>	<b>Materials</b>	892
849	• Location: United States only	• Response pairs: 75 accepted/rejected minimal	893
850	• Language: Native English speakers	pairs	894
851	• MTurk approval rate: $\geq 97\%$	• Rules covered: pronoun, voice, tense, politeness,	895
852	• Completed HITs: $\geq 1,000$	evaluative adverbs	896
853	<b>Consent</b> All participants confirmed: "I am 18+	<b>Dimensions</b> Based on the Communicative Em-	897
854	and consent to participate" before proceeding.	pathy framework (Buechel et al., 2018):	898
855	<b>E.2 Quality Control</b>	1. <b>Understanding</b> (Perspective-taking)	899
856	<b>Attention Checks</b> Each study included manda-	2. <b>Supportiveness</b> (Affective concern)	900
857	tory attention check questions to ensure data qual-	3. <b>Helpfulness</b> (Action-oriented guidance)	901
858	ity:	<b>Procedure</b> Each of 75 pairs was evaluated on	902
859	<i>"Attention check: To show you are reading,</i>	all 3 dimensions. Questions were labeled with	903
860	<i>please select option B."</i>	dimension (e.g., "[Understanding] What was the	904
861	Participants who failed attention checks were ex-	outcome...").	905
862	cluded from analysis and their HITs were rejected.	<b>Sample Size</b>	906
863	<b>Response Order Randomization</b> For all A/B	• Pairs: 75, Dimensions: 3, Questions per rater:	907
864	comparisons, the presentation order of responses	225	908
865	was randomized to prevent position bias. The	• Raters: 9, Total judgments: 2,025	909
866	true_label field in our data files tracks which	<b>Results</b> Voice transformations achieved 80% ac-	910
867	response (A or B) corresponds to each condition.	ceptance ( $h=0.64$ , $p < 0.001$ ); other rules showed	911
868	<b>E.3 Study 1: DPO vs. SFT Comparison</b>	lower acceptance (13–37%).	912
869	<b>Objective</b> Evaluate whether DPO-aligned re-	<b>E.5 Study 3: Length Bias Analysis</b>	913
870	sponses are preferred over supervised fine-tuning	<b>Objective</b> Isolate the contribution of the evalua-	914
871	(SFT) baseline responses.	tive adverbs rule from potential length confounds.	915
872	<b>Materials</b>	<b>Materials</b>	916
873	• Response pairs: 50 held-out legal questions	• Response pairs: 50 (same questions as Study	917
874	• Conditions: DPO (5 rules) vs. SFT baseline	1)	918
875	<b>Procedure</b> Participants viewed pairs of re-	• Conditions: DPO (4 rules, evaluative adverbs	919
876	sponses (labeled A and B) to legal questions and	removed) vs. SFT	920
877	selected which response was "MORE EMPATHETIC		
878	overall."		

Settings	
Experiment 2 Aligned Unaligned 50 Pairs with No Word Rule Final	
<a href="#">View Project</a>	
Note: If you have edited the Project after publishing this Batch, you will see the latest version.	
Description:	Compare pairs of legal advice responses and select which is more empathetic. 20 minutes, \$5.50. NATIVE English speakers only.
Keywords:	empathy, legal, advice, comparison, survey, native english, evaluation, text HIT Approval Rate (%) for all Requesters' HITs greater than 97
Qualification Requirement(s):	Number of HITs Approved greater than 1000 Location is <u>US</u>
Number of Assignments per task: 10	
Reward per Assignment:	\$5.50
Batch expired on: October 13, 2025 7:44 PM PDT	
Assignment duration: 1 hour	
Auto Approval Delay: 3 days	

Figure 5: MTurk HIT settings: qualification requirements (US,  $\geq 97\%$  approval,  $\geq 1000$  HITs), compensation (\$5.50), and assignment duration (1 hour).

**Procedure** Identical to Study 1, but using a DPO model trained without the evaluative adverbs rule.

**Sample Size** 50 questions, 500 total judgments.

**Results** Removing the evaluative adverbs rule reduced DPO preference from 68.8% to 59.0% ( $-9.8\text{pp}$ ,  $p=0.033$ , Cohen's  $h=0.18$ ). Effect decomposition reveals: Quality (direct)  $+10.5\text{pp}$ , Length (mediated)  $-0.7\text{pp}$  (Figure 6).

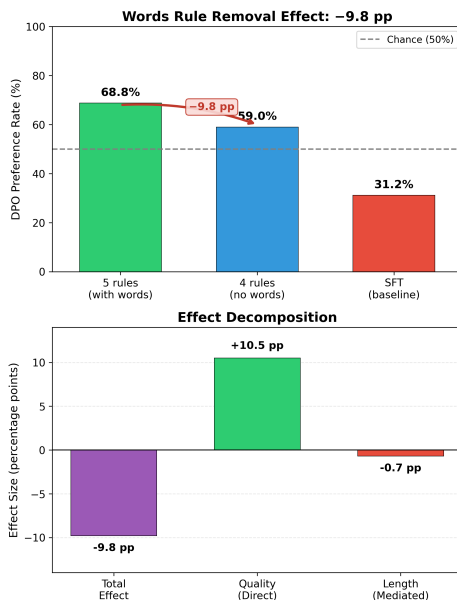


Figure 6: Words rule removal effect ( $-9.8\text{pp}$ ) and decomposition: quality ( $+10.5\text{pp}$ ) vs. length ( $-0.7\text{pp}$ ).

Study	Pairs	Judg.	Raters
1: DPO vs SFT	50	784	16/q
2: Validity	225	2,025	9
3: Length	50	500	-
<b>Total</b>	-	<b>3,309</b>	<b>35</b>

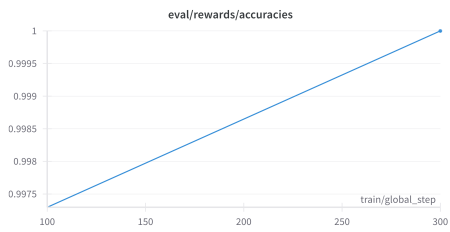
Table 4: Summary of human evaluation studies.

## F Training Details

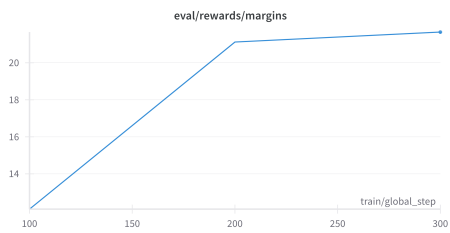
Figure 7 shows evaluation metrics for LLaMA-3 8B; similar patterns were observed for other models.

## E.6 Summary

Table 4 summarizes all human evaluation studies.



(a) Accuracy



(b) Margin



(c) Loss

Figure 7: Evaluation metrics during DPO training (LLaMA-3 8B).