Comparing human and LLM politeness strategies in free production

Anonymous ACL submission

Abstract

Polite speech poses a fundamental alignment challenge for large language models (LLMs). Humans deploy a rich repertoire of linguistic strategies to balance informational and social goals - from positive approaches that build rapport (compliments, expressions of interest) to negative strategies that minimize imposition (hedging, indirectness). We investigate whether LLMs employ a similarly context-sensitive repertoire by comparing human and LLM re-011 sponses in both constrained and open-ended production tasks. We find that larger models $(\geq 70B \text{ parameters})$ successfully replicate key preferences from the computational pragmatics 016 literature, and human evaluators surprisingly prefer LLM-generated responses in open-ended 018 contexts. However, further linguistic analy-019 ses reveal that models disproportionately rely on negative politeness strategies even in positive contexts, potentially leading to misinterpretations. While modern LLMs demonstrate an impressive handle on politeness strategies, these subtle differences raise important questions about pragmatic alignment in AI systems.

1 Introduction

027

028

034

042

Speakers do not always say exactly what they mean. For example, we might say a friend's poem "wasn't terrible" rather than saying "it was bad" to avoid hurting their feelings (Yoon et al., 2020), or just compliment specific elements that we liked without mentioning other elements (Brown and Levinson, 1987; Goffman, 1967; Pinker et al., 2008). These kinds of politeness strategies allow speakers to balance competing goals, conveying accurate information while maintaining positive relationships (Hill et al., 1986; Leech, 2014). As large language models (LLMs) are increasingly deployed in openended interactions across sensitive social domains like healthcare and education, their ability to appropriately use and understand polite language remains an important alignment challenge.

Politeness theory provides a valuable framework for addressing these questions. Seminal work by Brown and Levinson (1987) distinguishes between positive politeness strategies that affirm the listener (compliments, expressions of interest) and negative politeness strategies that minimize imposition (hedging, indirectness). This distinction is crucial because different contexts call for different strategies, and mismatches can lead to communication breakdowns. If an AI system employs distancing, hedge-filled language ("I'm somewhat concerned that this approach might not be optimal") in contexts where human speakers would use rapportbuilding strategies ("I love your creativity here, and wonder if we could build on it by ... "), users may perceive the system as cold, insincere, or lacking genuine engagement-even when its literal content is appropriate.

043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

079

This pragmatic misalignment represents a critical gap in our understanding of LLMs as social agents. While considerable attention has been paid to whether models can recognize politeness or generate polite language in constrained settings, effective social interaction depends not just on understanding politeness norms in the abstract but on actively selecting and applying appropriate strategies from a diverse linguistic repertoire. Human speakers navigate this complexity intuitively, deploying hedging, elaboration, indirect speech acts, and numerous other strategies to balance competing communicative goals in context-sensitive ways. To fully understand LLMs' grasp of politeness strategies, we need to examine whether they exhibit similar patterns of strategy selection and deployment across different contexts. This requires moving beyond limited-choice evaluations to examine open-ended language generation, where models have access to the full range of linguistic choices.

In this work, we investigate polite speech generation in both humans and LLMs, making the following contributions:

- We test whether LLMs reproduce human patterns of goal sensitivity in polite feedback using constrained response sets from Yoon et al. 086 (2020).
 - We collect and analyze a new dataset of openended responses from both humans and LLMs to identical social scenarios, enabling direct comparison of politeness strategies.
 - We perform detailed linguistic analyses to identify systematic differences in how humans and LLMs deploy various categories of politeness strategies.

Our results reveal that while LLMs have acquired important aspects of human-like pragmatic competence in polite language production - enough to be preferred by human evaluators - they also show systematic differences in strategy deployment that raise intriguing questions about the mechanisms underlying their social language capabilities. In particular, we find that models disproportionately rely on negative politeness strategies (minimizing imposition) even in contexts where humans prefer positive politeness strategies (building rapport), suggesting important differences in how these systems navigate social interactions.

Related Work 2

100

101

102

103

104

105

106

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126 127

128

129

130

131

2.1 **Computational Models of Politeness**

Research on politeness in linguistics and cognitive science has evolved from descriptive frameworks to quantitative models of pragmatic language use. Foundational work by Brown and Levinson (1987) established a systematic taxonomy of strategies, which has provided conceptual scaffolding for subsequent computational approaches. Studies in computational linguistics have since documented various linguistic markers of politeness across languages and contexts, examining formal features such as hedging, indirectness, and specific syntactic constructions correlate with perceived politeness (Danescu-Niculescu-Mizil et al., 2013; Aubakirova and Bansal. 2016).

Recent models in the Rational Speech Act (RSA) framework have explained the use of polite language as emerging from tradeoffs between informational utility capturing the desire for accuracy, a social utility representing the goal of making listeners feel good, and a self-presentational term reflecting speakers' desire to be perceived as both kind and honest (Yoon et al., 2020; Lumer and Buschmeier, 2022; Carcassi and Franke, 2023; Gotzner and 133 Scontras, 2024). This body of work has established 134 a solid theoretical foundation for analyzing polite-135 ness as a pragmatic phenomenon arising from un-136 derlying tradeoffs. However, existing models have 137 primarily focused on explaining choices among a 138 small number of constrained utterance alternatives 139 rather than modeling the rich variety of strategies 140 humans employ in open-ended generation contexts. 141

132

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

2.2 Pragmatic Capabilities in LLMs

Recent research has explored various aspects of pragmatic competence in large language models. Studies have examined LLMs' ability to understand indirect speech acts (Ruis et al., 2024; Jian and Narayanaswamy, 2024), recognize conversational implicatures (Hu et al., 2022; Lipkin et al., 2023), and interpret non-literal language (Yerukola et al., 2024; Liu et al., 2024). These investigations predominantly employ multiple-choice formats, presenting models with pragmatic puzzles and evaluating their ability to select contextually appropriate interpretations. Results generally suggest that modern LLMs demonstrate sophisticated pragmatic understanding, often approaching human-like performance on benchmark tasks. However, these studies primarily assess recognition rather than production capabilities, leaving open questions about whether models can actively deploy pragmatic strategies in their own generated outputs.

2.3 Polite Language Generation in LLMs

Work on generating polite language in AI systems represents a smaller but growing research area. Early approaches focused on style transfer, with systems like those developed by Niu and Bansal (2018) demonstrating that neural models could transform neutral text into more polite versions through specific syntactic transformations. Subsequent work explored paraphrasing to increase politeness (Fu et al., 2020), politeness-focused style transfer (Madaan et al., 2020), and creating polite chatbots (Mukherjee et al., 2023). However, these systems typically focused on surface-level transformations rather than strategic deployment of politeness based on contextual factors. As noted in a recent survey (Priya et al., 2024), existing approaches to polite language generation have predominantly emphasized isolated features (hedging expressions, please markers, specific lexical choices) rather than examining the full repertoire



Figure 1: (A) Correlations between human and model response probabilities for the top 4 models with specific prompting strategies we tested. Both the base and instruct-tuned versions of Qwen2.5-72B are shown here for comparison. Error bars are 95% confidence intervals across vignettes. (B) Comparing the pattern of human and LLM responses across different communicative goals and ratings. Model results are from Llama-3.3-70B-Instruct using the multi-choice-persona prompting strategy; human responses are from Yoon et al. (2020).

of politeness strategies and how they're selected based on communicative context. This leaves a significant gap in our understanding of whether LLMs can approximate the context-sensitivity that characterizes human politeness. Our work addresses this gap by directly comparing politeness strategies in humans and LLMs across varying communicative goals, examining whether models align with human preferences for positive versus negative politeness strategies in different contexts.

182

183

184

185

189

190

191

192

193

194

195

196

199

200

204

207

3 Experiment 1: Constrained settings

To what extent are LLMs sensitive to the goals that give rise to politeness in human speech? To address this question, we first examined whether LLMs could reproduce the patterns of goal-sensitive language use reported by Yoon et al. (2020). Their study provided empirical evidence for a computational model of politeness where speakers strategically balance informational accuracy with social goals. Most notably, they found that when giving negative feedback, humans often deploy negation (e.g., "wasn't terrible" rather than "bad") (see also Gotzner and Scontras, 2024).

We reimplemented this experiment with LLMs to assess their pragmatic competence in a constrained setting. In each scenario, a character gives feedback about another character's performance (e.g., a piano play or presentation), with the true quality ranging from 0 to 3 hearts. The speaker has one of four communicative goals: to be *informative*, to be *kind*, or *both*. We also added a *default* condition with no explicit goal specified to understand how LLMs behave by default. Models selected from the same set of eight responses used by humans, combining either "was" or "wasn't" with four adjectives (terrible, bad, good, amazing). 208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

226

227

228

229

230

231

232

233

234

We tested a range of open-source (8B-72B parameters) and closed-source models using two prompting strategies: an "original" strategy that presented scenarios verbatim, and a "persona" variant that systematically varied speaker characteristics (e.g., gender, occupation, background) to better approximate the diversity in the population of human participants (Murthy et al., 2025; He-Yueya et al., 2024). For each model and prompting strategy, we sampled 30 responses with temperature $\tau = 1.0$ per scenario (see Appendix A for details).

3.1 Model comparison

We report Spearman correlation and mean squared error (MSE) between LLM and human responses as an overall measure of fit (see Table 1). These results suggest that model size plays a crucial role in capturing human-like politeness strategies.

	Co	omparison	with huma	Compari	son with d	efault goal	
LLMs	Spearman Original Persona		MSE Original Persona		vs. Both	vs. Inf.	vs. Social
GPT-40	0.75	0.76	0.026	0.031	0.62	0.99	0.31
Claude-3.5-Sonnet	0.41	0.47	0.048	0.046	0.73	0.49	0.19
Llama-3.1-8B	0.11	0.15	0.052	0.052	0.77	0.86	0.71
Llama-3.1-8B-Instruct	0.17	0.17	0.061	0.063	0.87	0.75	0.78
Llama-3.1-70B	0.66	0.67	0.034	0.030	0.86	0.58	0.57
Llama-3.1-70B-Instruct	0.73	0.74	0.023	0.024	0.74	0.75	0.53
Llama-3.3-70B-Instruct	0.67	0.66	0.018	0.019	0.80	0.64	0.40
Mixtral-8x7B	0.36	0.35	0.043	0.044	0.74	0.83	0.19
Mixtral-8x7B-Instruct	0.43	0.39	0.080	0.082	0.54	0.41	0.10
Qwen2.5-72B	0.65	0.66	0.028	0.029	0.83	0.75	0.73
Qwen2.5-72B-Instruct	0.66	0.64	0.033	0.034	0.63	0.55	0.54

Table 1: Comparison of LLM response patterns with humans (Yoon et al., 2020) and between default goal and other communicative goals (using the "multi-choice-original" prompting strategy and Spearman correlation).

235 Smaller models (Llama-3.1-8B) showed essentially no correlation with human responses, often 236 failing to perform the multi-choice task at all, while 237 intermediate-sized models like Mixtral-8x7B (effective model-size is 13B (Jiang et al., 2024)) 239 showed only modest correlations. However, larger 240 models (≥70B parameters) demonstrated much 241 stronger alignment with human behavior, with 242 243 Llama-3.3-70B-Instruct achieving the highest correlations among open-source models (Spearman 244 r = 0.67). Among closed-source models, GPT-40 245 displayed particularly strong performance (Spearman r = 0.75), while Claude-3.5-Sonnet lagged 247 behind with more modest correlations (r = 0.41). These findings suggest that sophisticated pragmatic competence for politeness emerges primarily in larger models, potentially reflecting the greater contextual sensitivity needed to balance competing 253 communicative goals.

3.2 Error analysis

256

258

259

261

263

264

267

Despite strong overall correlations (see Figure 1A), even the best-performing models showed systematic differences from human responses. To better understand these patterns, we conducted a detailed comparison with human responses following the visualization approach in Yoon et al. (2020). The results in Figure 1B show that the best-fitting opensource model captures many key features of the human response patterns. Most notably, when rating a poor performance (0/3 hearts) with both informational and social goals, the model appropriately deploys negation as a politeness strategy, just as humans do: both humans and LLMs prefer to say "wasn't terrible" rather than "was bad". The model also closely tracks human preferences for positive ratings (2-3 hearts), showing appropriate sensitivity to the quality of the performance.

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

286

287

290

291

292

293

294

295

298

299

300

However, key differences emerged in the granularity of responses. Where humans show graded preferences across response options (distributing probability mass across multiple choices), LLMs tend toward more categorical binary choices, either strongly preferring or completely avoiding certain responses. They consistently choose one single option given a context, rating, and goal combination in most cases—despite our efforts to increase response diversity through temperature sampling ($\tau = 1.0$) or persona variation in prompting.

Closer analysis also revealed systematic differences in how well LLMs captured human behavior across different communicative goals. Models showed stronger alignment with human responses for the *social* goal but underperformed when the goal was to be purely *informative*. For example, when humans prioritize being informative about a poor performance, they often select direct negative feedback ("was bad"), while LLMs sometimes persist with softened language. This pattern suggests that, while LLMs have acquired some aspects of sophisticated politeness strategies, they may overapply these strategies even when directness would be more appropriate, potentially reflecting their training to be generally "helpful and harmless".

3.3 Default goal analysis

We included a *default* goal condition (no explicitly specified goal) to evaluate how LLMs respond

306

307

310

311

312

314

315

317

322

323

328

329

331

334

338

342

347

350

301

without specific communicative instructions. This condition helps reveal the implicit goals that might have been induced through various stages of model training. Although the overall fit to human data varies across models, we can ask which explicit goal produces the closest response pattern to the default goal, as measured by Spearman correlation.

Overall, we find a stronger resemblance to the *both* goal (see Figure 1B), suggesting that models generally attempt to balance informativeness and social considerations by default. However, Table 1 reveals varying correlation patterns across different LLMs. While most models show stronger correlations with the *both* goal, others correlate more strongly with the *informative* goal. For instance, L1ama-3.3-70B-Instruct appears to implicitly align with *both* (Spearman r = 0.80), whereas GPT-40 shows much stronger alignment with *informative* (Spearman r = 0.99).

These varied patterns suggest that the implicit goals guiding different LLMs' polite speech may reflect differences in their training objectives and alignment procedures. The dominant pattern of alignment with the *both* goal is consistent with the general instruction to models to be both helpful (informative) and harmless (socially appropriate). However, the variability across models indicates that while these systems have acquired sophisticated politeness capabilities, the specific ways they balance competing goals may differ from model to model.

4 Experiment 2: Open-ended generation

While our multiple-choice experiment demonstrated that larger LLMs can reproduce basic patterns of goal-sensitive politeness strategies, such as the strategic use of negation, this constrained format limits our understanding of how models deploy politeness in naturalistic settings. In real-world interactions, speakers draw from a rich repertoire of linguistic devices beyond those provided in fixed-choice scenarios. This raises a critical question: how do LLMs perform when given the freedom to generate polite language from scratch?

To address this question, we designed an openended generation experiment that uses the same scenarios as our multiple-choice study but removes the response constraints. This approach allows us to examine whether LLMs employ a similarly diverse and context-sensitive set of politeness strategies as humans when both have access to the full expressivity of language, and directly compare to results in the constrained setting.

4.1 Methods

We used the scenarios from Yoon et al. (2020), preserving the same performance ratings (0-3 hearts) and communicative goals (informative, social, both, default). We collected 3 open-ended responses per scenario from 156 human participants via Prolific (each responding to 4 distinct scenarios) and three responses from LLMs that performed well in our first experiment: GPT-40, Claude-3.5-Sonnet, and Llama-3.3-70B-Instruct. Models were instructed to "keep responses short and concise" to ensure comparable length with human responses. To assess preferences for these responses, we then conducted a two-alternative forced-choice evaluation with 156 human evaluators, each viewing four different scenarios. Evaluators made five judgments per scenario: (1) comparing human vs. LLM responses, (2-3) comparing goal-congruent vs. goal-incongruent responses for both sources, and (4-5) comparing rating-congruent vs. ratingincongruent responses for both sources. We randomized presentation order and ensured evaluators saw responses from different sources across blocks (see Appendix B for full details).

4.2 Results

Overall preferences Surprisingly, human evaluators showed a marked preference for LLMgenerated responses over human-generated ones across all goal types (66% of all trials; see Figure 2A). A mixed-effects logistic regression containing random intercepts at the evaluator and item level confirmed this preference was significantly different from chance (z = 7.63, p < 0.001). This pattern held for each of the four communicative goals, with the largest effect observed for the informative goal (22% above baseline) and the smallest effect observed for the default goal (8.3% above baseline; see Figure 2A). However, there were systematic differences in the strength of these preferences across goals; a model including a fixed effect of goal accounted for significantly more variance than the intercept-only model, according to a likelihood-ratio test $\chi^2(3) = 12.54, p = 0.006.$

Goal Sensitivity Next, we considered the extent to which human-generated and LLM-generated utterances were goal-sensitive by calculating the proportion of trials where participants preferred 354 355

351

352

353

356

357

358

359

360

361

362

363

369

370

371

372

373

374

375

376

377 378 379

380

381

382

383

386

387

389

390

391

392

393

394

395

396

397

398

399



Figure 2: Human evaluation results. The bars show the relative preference (50% is chance). Bars above the 50% line indicate the percentage to which responses are preferred as expected, and below indicate the percentage to which responses are preferred as unexpected. (A) Evaluators systematically prefer LLM generations over human generations. (B) Both humans and LLMs are sensitive to goals and (C) ratings. Error bars are bootstrapped 95% confidence intervals.

a congruent utterance (i.e., an utterance actu-400 ally produced to achieve the given goal) over 401 402 an incongruent utterance (i.e., one produced under a different goal). We found that both hu-403 mans and LLMs demonstrated sensitivity to com-404 municative goals: evaluators preferred the goal-405 congruent human response 15.9% above-baseline 406 (z = 6.89, p < 0.001), and preferred the goal-407 congruent LLM response even more strongly at 408 25.0% above baseline (z = 9.59, p < 0.001). 409 Moreover, Figure 2B suggests that LLMs main-410 tained greater or equal goal sensitivity across all 411 four goals, indicating they successfully tailored 412 413 their language to the specified communicative objective. 414

415 **Rating Sensitivity** Finally, as a sanity-check, we asked whether utterances were sensitive to 416 the actual state of the speaker (i.e., the number 417 of hearts they felt about the performance being 418 evaluated). Again, both groups showed strong 419 rating sensitivity. Human responses achieved a 420 20.8% above-baseline preference for aligned rat-421 ings (z = 8.32, p < 0.001), while LLM responses 422 demonstrated even higher sensitivity with a 26.1% 423 above-baseline preference (z = 9.59, p < 0.001). 424 As shown in Figure 2C, LLMs maintained equal 425 or improved sensitivity across all goals, indicating 426 that they are not simply producing generically po-427 428 lite utterances but are modulating their responses appropriately as a function of both the basic infor-429 mation to be conveyed (the rating) and the speci-430 fied communicative goal (e.g., being informative 431 vs. making someone feel good). 432

5 Linguistic Analysis of Politeness

While our evaluations show that LLMs successfully generate polite language that human evaluators prefer, these preferences alone don't reveal whether models use the same linguistic mechanisms as humans. To understand the specific politeness strategies employed by both humans and LLMs, we conducted a detailed linguistic analysis of the openended responses. 433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

5.1 Negation

As a first step in our analysis, we examined how frequently the strategic use of negation documented in Yoon et al. (2020) and tested in Experiment 1 is employed in open-ended responses. Among all 1,248 responses collected, 527 (42.2%) used the specific pattern of adjective evaluation studied by Yoon et al. (2020). Within this subset, 35 responses (6.6%) employed negation as a politeness strategy, and negation was most common in low-rating (0 or 1 heart) scenarios, which qualitatively replicates our findings from Experiment 1 (see Figure 5 in Appendix for details). Thus, the negation strategies studied in constrained settings do appear in openended production, but represent just one of many politeness devices available to speakers.

5.2 Word usage patterns

To better understand differences between human and LLM responses, we analyzed unigram distributions using Pointwise Mutual Information (PMI) and Jensen-Shannon Divergence (JSD). First, examining unigrams with highest PMI, we found that human responses more frequently incorporated casual language and expressions (e.g., "awesome," "great"), whereas LLM-generated responses tended

A human strategy use

Figure 3: Proportion of different politeness strategies across ratings and goals for (A) human and (B) LLMs.

toward more formal linguistic choices (e.g., "fabulous," "excellent"). Both groups effectively employed personalization as a politeness strategy, such as directly mentioning the performer's name (e.g., "Your app is pretty good, Henry!"). Additionally, both humans and LLMs adapted their lexical choices based on context, with minimal overlap in high-PMI words across different goals and ratings.

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492 493

494

495

496

497

Next, we quantified differences in empirical word frequency distributions by calculating the Jensen-Shannon Divergence (JSD). Interestingly, the JSD between the lexical distributions of preferred and non-preferred response groups was quite small (JSD = 0.013) though still significantly different than a permuted null distribution (p < 0.001), while all other group comparisons showed much larger differences (JSD > 0.13, p < 0.001; see Table 8). This suggests that simple lexical choice may not be the primary driver of human preferences in polite language.

Finally, we conducted higher-dimensional analyses using SBERT embeddings (Reimers and Gurevych, 2019) to distinguish between response categories. These analyses (described in Appendix C.2) revealed that while human vs. LLM responses were readily distinguishable in embedding space (83% accuracy), preferred vs. non-preferred responses were much harder to classify (54% accuracy). LLM responses were more distinguishable across different communicative goals than human responses, suggesting more stereotyped strategies.

5.3 Annotated politeness strategies

498

499

500

501

502

503

504

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

To obtain a comprehensive picture of the politeness strategies employed in human and LLM responses, we conducted a detailed annotation using the politeness framework from Brown and Levinson (1987), supplemented by markers from Danescu-Niculescu-Mizil et al. (2013). This framework distinguishes four broad categories of politeness strategies with many subtypes (see Appendix Table 10): positive politeness (e.g. compliments and expressions of interest), negative politeness (e.g. hedging and indirectness), off-record strategies (indirect hints that maintain plausible deniability), and bald-on-record strategies (direct statements without politeness). We used LLMs (GPT-4.1 and Claude-3.7-Sonnet) as annotators, following best practices (Tan et al., 2024). Annotations were manually verified and corrected where necessary. This approach allowed us to identify specific politeness markers and their associated strategies across the dataset. Appendix C.3 provides full details on the annotation process and framework.

Strategy Distribution. As shown in Figure 3, both humans and LLMs rely primarily on positive politeness (rapport-building) and negative politeness (minimizing imposition) strategies, with relatively low use of off-record and bald-on-record approaches. However, a key difference emerged in strategy selection patterns: while both humans and LLMs increased their use of positive polite-

ness strategies as ratings increased, LLMs showed 528 systematically higher use of negative politeness strategies even in positive contexts (higher ratings), where humans tended to reduce such strategies. These strategy distributions were significantly 532 different under a permutation test (JSD = 0.023, 533 p < 0.001). The overreliance on negative polite-534 ness strategies in positive contexts represents a fundamental misalignment with human communication patterns. It suggests LLMs may have been 537 trained to prioritize non-imposition over positive affirmation, potentially reflecting alignment pro-539 cedures that emphasize harm reduction over more 540 natural social interactions. 541

Goal and Rating Sensitivity. We observed that both humans and LLMs appropriately varied their strategy distribution by communicative goal, with the informative goal showing the most distinct pattern. For the informative goal with high ratings (2-3 hearts), LLMs showed unexpectedly higher use of negative politeness strategies compared to humans, who shifted toward positive strategies in these contexts. This pattern suggests that LLMs may overuse hedging, conventional indirectness, and other distancing strategies even when giving positive feedback, potentially explaining some of the stylistic differences observed in the evaluation. However, as with word usage patterns, the differences between strategy distributions for preferred and non-preferred responses were not significant (JSD = 0.008, p = 0.087), suggesting that preference judgments may be driven by higher-order social factors beyond the mere presence or absence of specific words or politeness strategies.

6 Conclusion

542

543

544

545

547

551

553

555

557

559

562

563

564

565

566

567

571

573

574

577

While LLMs demonstrate impressive pragmatic competence in politeness, we find that they also systematically differ from humans in how they deploy different strategies. Most notably, LLMs over-rely on negative politeness strategies even in positive contexts where humans prefer rapport-building approaches. This mismatch persists despite LLM responses being consistently preferred by human evaluators in open-ended settings, suggesting that surface-level preferences may not fully capture important pragmatic differences. The contrast between constrained and open-ended tasks further highlights how LLMs leverage a richer repertoire of strategies when given more expressivity, such as the "while COMPLIMENT, SUGGESTION" constructions we observed in low-rating scenarios.

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

The distributional differences in politeness strategy deployment may have significant implications for human-AI communication. When AI systems consistently employ distancing strategies even when delivering positive feedback, this contextual mismatch could lead to pragmatic misinterpretations, where humans interpret hedged positive feedback as more negative than intended. Such patterns may also reduce social presence by making AI communication feel unnecessarily formal and detached, ultimately diminishing information transfer if humans expend additional cognitive resources parsing unnecessarily indirect language. These consequences extend beyond mere stylistic differences to potentially impact the fundamental goals of communication: accurate information transfer and relationship maintenance.

Our results underscore the need to investigate not just whether AI systems can produce polite language, but whether their specific patterns of politeness strategies facilitate or hinder effective communication. For example, future work should directly test how differences in strategy deployment affect human listeners' ability to recover underlying information and speakers' intentions, particularly in extended multi-turn conversations and across diverse cultural contexts. The systematic divergence we observed raises a concern that even if LLM responses are preferred in isolation, their departure from human pragmatic patterns may lead to miscommunications in real-world interactions where precise understanding of intentions and attitudes is crucial. These findings could inform future LLM training approaches to better align with human pragmatic patterns, particularly in balancing positive and negative politeness strategies.

Finally, our results have implications for rational models of politeness in computational linguistics: given the richer space of strategies exposed in openended production, it will be essential to develop richer computational models that can probe the implicit tradeoffs guiding LLM behavior. For instance, we observed that LLMs commonly use constructions like "COMPLIMENT but SUGGESTION" in low-rating scenarios, a strategy that balances kindness with informativeness in ways that cannot be accomplished via simple negation. But, if we could be so bold as to make one suggestion, we'd greatly appreciate it if future work could help systems learn when to use positive strategies rather than minimize imposition.

630 Limitations

While our work gave a comprehensive picture of comparing the polite language use in humans and 632 LLMs, there are still limitations that could be ad-633 dressed in future work. First, throughout our analyses, we still cannot answer the question of what makes human evaluators prefer the responses they 636 prefer, as all our analyses showed very minimal differences between preferred and non-preferred responses. One guess is that even ratings and goals are made very clear in the provided scenarios, human evaluators still may not pay enough attention 641 to, and optionally omit this information, instead, 642 they tend to pick whichever one in the given pair that sounds nicer. Future research, for example, testing LLMs as evaluators and comparing LLM-645 as-evaluator preference results with humans, could give us more insight into this question. Additionally, as our results show that LLMs are still not quite human-like in picking the right politeness strategies in a context-sensitive way, future research on how to develop computational methods and algorithms to make LLMs better at polite language use and as social agents will be necessary. 653

References

654

655

657

667

670

671

672

673

674

675

678

- Malika Aubakirova and Mohit Bansal. 2016. Interpreting neural networks to improve politeness comprehension. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2035–2041.
 - Penelope Brown and Stephen C Levinson. 1987. Politeness: Some universals in language usage. Cambridge university press.
 - Fausto Carcassi and Michael Franke. 2023. How to handle the truth: A model of politeness as strategic truth-stretching. In *Proceedings of the Annual Meeting of the Cognitive Science Society*.
 - Cristian Danescu-Niculescu-Mizil, Moritz Sudhof, Dan Jurafsky, Jure Leskovec, and Christopher Potts. 2013. A computational approach to politeness with application to social factors. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 250–259.
- Liye Fu, Susan Fussell, and Cristian Danescu-Niculescu-Mizil. 2020. Facilitating the communication of politeness through fine-grained paraphrasing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 5127–5140.
 - Erving Goffman. 1967. Interaction ritual: Essays in face-to-face behavior. Routledge.

Nicole Gotzner and Gregory Scontras. 2024. On the role of loopholes in polite communication: Linking subjectivity and pragmatic inference. *Open Mind*, 8:500–510.

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

- Joy He-Yueya, Wanjing Anya Ma, Kanishk Gandhi, Benjamin W Domingue, Emma Brunskill, and Noah D Goodman. 2024. Psychometric alignment: Capturing human knowledge distributions via language models. *arXiv preprint arXiv:2407.15645*.
- Beverly Hill, Sachiko Ide, Shoko Ikuta, Akiko Kawasaki, and Tsunao Ogino. 1986. Universals of linguistic politeness: Quantitative evidence from japanese and american english. *Journal of pragmatics*, 10(3):347–371.
- Jennifer Hu, Sammy Floyd, Olessia Jouravlev, Evelina Fedorenko, and Edward Gibson. 2022. A finegrained comparison of pragmatic language understanding in humans and language models. *arXiv preprint arXiv:2212.06801*.
- Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*, 43(2):1–55.
- Mingyue Jian and Siddharth Narayanaswamy. 2024. Are LLMs good pragmatic speakers? In *NeurIPS Workshop on Behavioral Machine Learning.*
- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, Lélio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, and 7 others. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.

Geoffrey Leech. 2014. The pragmatics of politeness.

- Benjamin Lipkin, Lionel Wong, Gabriel Grand, and Josh Tenenbaum. 2023. Evaluating statistical language models as pragmatic reasoners. In *Proceedings of the Annual Meeting of the Cognitive Science Society.*
- Ryan Liu, Theodore Sumers, Ishita Dasgupta, and Thomas L Griffiths. 2024. How do large language models navigate conflicts between honesty and helpfulness? In *Forty-first International Conference on Machine Learning*.
- Eleonore Lumer and Hendrik Buschmeier. 2022. Modeling social influences on indirectness in a rational speech act approach to politeness. In *Proceedings of the Annual Meeting of the Cognitive Science Society*, volume 44.

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

791

735

736

- 742 743
- 744 745
- 746
- 747 748 749 750 751
- 752 753 754 755 756
- 757 758 759 760 761
- 763 764 765 766
- 768 769 770 771 772 773
- 774 775

777 778

7

779

781 782 783

78

- 7
- 788 789 790

- Aman Madaan, Amrith Setlur, Tanmay Parekh, Barnabas Poczos, Graham Neubig, Yiming Yang, Ruslan Salakhutdinov, Alan W Black, and Shrimai Prabhumoye. 2020. Politeness transfer: A tag and generate approach. *arXiv preprint arXiv:2004.14257*.
- Sourabrata Mukherjee, Vojtěch Hudeček, and Ondřej Dušek. 2023. Polite chatbot: A text style transfer application. In *Proceedings of the 17th Conference* of the European Chapter of the Association for Computational Linguistics: Student Research Workshop, pages 87–93.
- Sonia Krishna Murthy, Tomer Ullman, and Jennifer Hu. 2025. One fish, two fish, but not the whole sea: Alignment reduces language models' conceptual diversity. In Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 11241– 11258, Albuquerque, New Mexico. Association for Computational Linguistics.
- Tong Niu and Mohit Bansal. 2018. Polite dialogue generation without parallel data. *Transactions of the Association for Computational Linguistics*, 6:373–389.
- Steven Pinker, Martin A Nowak, and James J Lee. 2008. The logic of indirect speech. *Proceedings of the National Academy of sciences*, 105(3):833–838.
- Priyanshu Priya, Mauajama Firdaus, and Asif Ekbal. 2024. Computational politeness in natural language processing: A survey. *ACM Computing Surveys*, 56(9):1–42.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERTnetworks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Laura Ruis, Akbir Khan, Stella Biderman, Sara Hooker, Tim Rocktäschel, and Edward Grefenstette. 2024. The goldilocks of pragmatic understanding: Finetuning strategy matters for implicature resolution by Ilms. *Advances in Neural Information Processing Systems*, 36.
- Zhen Tan, Alimohammad Beigi, Song Wang, Ruocheng Guo, Amrita Bhattacharjee, Bohan Jiang, Mansooreh Karami, Jundong Li, Lu Cheng, and Huan Liu. 2024.
 Large language models for data annotation: A survey. *arXiv e-prints*, pages arXiv–2402.
- Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. 2025. Hallucination is inevitable: An innate limitation of large language models. *Preprint*, arXiv:2401.11817.
- Akhila Yerukola, Saujas Vaduguru, Daniel Fried, and Maarten Sap. 2024. Is the pope catholic? yes, the pope is catholic. generative evaluation of non-literal

intent resolution in LLMs. In *Proceedings of the* 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 265–275, Bangkok, Thailand. Association for Computational Linguistics.

Erica J. Yoon, Michael Henry Tessler, Noah D. Goodman, and Michael C. Frank. 2020. Polite speech emerges from competing social goals. *Open Mind*, 4:71–87.

A Experiment 1 Methods

We closely followed the experimental paradigm of Yoon et al. (2020). In this study, participants read short scenarios about someone seeking feedback on a performance or creative work. Each scenario specified (1) the true quality of the work on a scale from 0 to 3 hearts and (2) the speaker's communicative goal – either to be informative, to make the person feel good, or to do both. Participants then chose what they would say from a restricted set of options, combining either was or wasn't with one of four adjectives: terrible, bad, good, or amazing. Scenarios were constructed from 13 different contexts (e.g., filmmaking, songwriting, concert performance), yielding 156 unique scenarios (13 contexts \times 4 ratings \times 3 goals). We also added a "default" condition with no explicitly specified goal, bringing our total to 208 scenarios.

To test LLMs on this task, we developed two prompting strategies. In our basic approach, which we called "multi-choice-original", we simply presented each scenario verbatim and asked the model to choose from the eight possible responses (all combinations of "was"/"wasn't" with the four adjectives). To better approximate the diversity of human participants and with the hope to see that diversifying the personas of LLMs would improve their performance, we also considered a "persona" variant where we systematically varied speaker characteristics like gender, occupation, and background, where we call "multi-choicepersona". We tested these approaches across a range of current LLMs, including both closedsource (GPT-4o, Claude-3.5-Sonnet) and opensource models (Llama-3, Mixtral, Qwen2.5) of varying sizes (8B to 70B parameters). For opensource models, we compared both base and instructtuned versions where available to see the influence of the post-training stage on this task. To approximate the multiple participants in human studies, we collected 30 responses per scenario from each model using a temperature of $\tau = 1.0$.

845

847

851

852

854

858

860

867

872

873

874

876

878

884

887

A.1 Prompting strategies

The system prompt remained consistent between the original and persona prompting strategies.

Multi-choice-original/persona system prompt: You will see a scenario below. In the scenario, person A is asking for person opinion on their performance. $\$ B's Person B's true feelings in the scenario are shown on a scale of 0 to 3 hearts.\n 0 heart, the lowest rating, means the person does not like the performance at all, and 3 hearts, the highest rating, the person likes it a lot.\n means Please read the scenario carefully and answer the question ONLY with one of the eight options provided. \n Please provide your response in the following format:\n Answer: < one of the eight possible answer options in the scenario>

To construct persona prompts, we varied the following details:

- Race: {white, Black, Asian, Hispanic, American-Indian}
- Gender: {woman, man, non-binary person}
 - City: {New York, Chicago, San Francisco, Boston, Houston}
 - Years of experience: {17, 18, 19, 20, 21, 22, 23}
 - Occupation: {a critic, an expert, a teacher, a friend, a colleague, an acquaintance}

For example, here is an example of what the scenarios look like with and without a persona.

Scenario without a persona: Imagine that John just gave a presentation, but John didn't know how good it was. John approached Chris, who knows a lot about giving presentations, and asked "How was my presentation?"

Scenario with a persona: Imagine that John just gave a presentation, but John didn't know how good it was. John approached Chris, an Asian man from Boston, who has 19 years of experience as a teacher in the field and knows a lot about giving presentations. John asked "How was my presentation?" A complete instance of our scenario as fed into the LLM user prompt looks like this — the example is a multi-choice-original, 2-hearts rating, informative goal scenario.

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

Context: Imagine that Bob just gave a presentation, but Bob didn't know how good it was. Bob approached John, who knows a lot about giving presentations, and asked "How was my presentation?"

Rating: Here's how John actually felt about Bob's presentation: 2 out of 3 hearts

Question: If John wanted to give as accurate and informative feedback as possible, but not necessarily make Bob feel good, What would John be most likely to say?

Options:

1. It was terrible.	906
2. It was bad.	907
3. It was good.	908
4. It was amazing.	909
5. It wasn't terrible.	910
6. It wasn't bad.	911
7. It wasn't good.	912
8. It wasn't amazing.	913
Additional results	014

A.2 Additional results

See Table 2 for a complete comparison between LLM and human response patterns across Spearman, Pearson, and MSE metrics using both multichoice-original and multi-choice-persona prompting strategies, as a complement to Table 1 in the main text, where Pearson correlation scores were not reported.

See Table 3 and 4 for comprehensive comparison results between human and LLM responses across different goals. Since the "default" goal case was not studied in Yoon et al. (2020), we focused on "both", "social", and "informative" goals and report both Pearson and Spearman correlation scores. For the two tables, we observed that different LLMs have medium to strong correlations with human responses. Both base and instruct-tuned versions of Llama-3.1-8B and Mixtral-8x7B showed very low correlation scores and their incompetence in generating polite language.

See Table 5 and 6 for a complete comparison report between "default" and other goals as a complement to Table 1 in the main text. We reported

937	both Pearson and Spearman correlation scores for
938	"multi-choice-original" and "multi-choice-persona"
939	prompting strategies. We found that including per-
940	sonas does not have any real effect, and the findings
941	are consistent with those reported in the main text.

LIMe	Pearson		Spearman		MSE	
LLIVIS	Original	Persona	Original	Persona	Original	Persona
GPT-40	0.84	0.81	0.75	0.76	0.026	0.031
Claude-3.5-Sonnet	0.50	0.55	0.41	0.47	0.048	0.046
Llama-3.1-8B	-0.02	-0.02	0.11	0.15	0.052	0.052
Llama-3.1-8B-Instruct	-0.05	-0.01	0.17	0.17	0.061	0.063
Llama-3.1-70B	0.59	0.63	0.66	0.67	0.034	0.030
Llama-3.1-70B-Instruct	0.76	0.74	0.73	0.74	0.023	0.024
Llama-3.3-70B-Instruct	0.83	0.82	0.67	0.66	0.018	0.019
Mixtral-8x7B	0.36	0.33	0.36	0.35	0.043	0.044
Mixtral-8x7B-Instruct	0.29	0.29	0.43	0.39	0.080	0.082
Qwen2.5-72B	0.71	0.70	0.65	0.66	0.028	0.029
Qwen2.5-72B-Instruct	0.78	0.77	0.66	0.64	0.033	0.034

Table 2: Complete version of model comparison results reported in Table 1 in the main text, including Pearson correlation scores.

LLMs	Both		Social		Inf.	
	Original	Persona	Original	Persona	Original	Persona
GPT-4o	0.89	0.89	0.80	0.75	0.83	0.78
Claude-3.5-Sonnet	0.58	0.65	0.49	0.52	0.46	0.53
Llama-3.1-8B	0.05	0.01	-0.09	-0.03	-0.03	-0.04
Llama-3.1-8B-Instruct	0.10	0.09	-0.15	-0.03	-0.11	-0.11
Llama-3.1-70B	0.66	0.73	0.79	0.81	0.36	0.39
Llama-3.1-70B-Instruct	0.92	0.91	0.88	0.87	0.46	0.42
Llama-3.3-70B-Instruct	0.92	0.92	0.88	0.88	0.70	0.68
Mixtral-8x7B	0.18	0.20	0.63	0.62	0.19	0.11
Mixtral-8x7B-Instruct	0.32	0.19	0.62	0.84	-0.08	-0.10
Qwen2.5-72B	0.83	0.73	0.85	0.83	0.47	0.57
Qwen2.5-72B-Instruct	0.86	0.85	0.84	0.83	0.65	0.64

Table 3: Pearson Correlation between human and LLMs over different goals

LLMs	Both		Social		Inf.	
	Original	Persona	Original	Persona	Original	Persona
GPT-40	0.70	0.71	0.74	0.74	0.80	0.80
Claude-3.5-Sonnet	0.38	0.49	0.43	0.47	0.41	0.41
Llama-3.1-8B	0.12	0.10	0.16	0.15	0.02	0.15
Llama-3.1-8B-Instruct	0.32	0.39	0.17	0.17	0.06	0.06
Llama-3.1-70B	0.65	0.68	0.63	0.64	0.69	0.72
Llama-3.1-70B-Instruct	0.66	0.74	0.67	0.67	0.86	0.83
Llama-3.3-70B-Instruct	0.64	0.68	0.58	0.60	0.77	0.72
Mixtral-8x7B	0.33	0.35	0.23	0.21	0.46	0.44
Mixtral-8x7B-Instruct	0.34	0.33	0.33	0.17	0.60	0.67
Qwen2.5-72B	0.70	0.64	0.59	0.58	0.69	0.74
Qwen2.5-72B-Instruct	0.69	0.63	0.63	0.66	0.64	0.64

Table 4: Spearman Correlation between human and LLMs over different goals

LLMs	vs. Both		vs. Inf.		vs. Social	
	Original	Persona	Original	Persona	Original	Persona
GPT-40	0.59	0.66	0.64	0.66	0.43	0.40
Claude-3.5-Sonnet	0.64	0.80	0.50	0.41	0.05	0.04
Llama-3.1-8B	0.76	0.70	0.83	0.67	0.70	0.69
Llama-3.1-8B-Instruct	0.87	0.81	0.57	0.66	0.72	0.74
Llama-3.1-70B	0.86	0.86	0.48	0.58	0.47	0.54
Llama-3.1-70B-Instruct	0.82	0.89	0.67	0.62	0.47	0.43
Llama-3.3-70B-Instruct	0.83	0.85	0.92	0.82	0.49	0.52
Mixtral-8x7B	0.69	0.61	0.78	0.87	0.01	-0.07
Mixtral-8x7B-Instruct	0.28	0.30	0.45	0.43	0.04	0.29
Qwen2.5-72B	0.77	0.78	0.66	0.72	0.70	0.73
Qwen2.5-72B-Instruct	0.84	0.86	0.87	0.75	0.57	0.54

Table 5: Pearson Correlation between default goal and other goals in Experiment 1

LLMs	vs. Both		vs. Inf.		vs. Social	
	Original	Persona	Original	Persona	Original	Persona
GPT-40	0.62	0.46	0.99	0.99	0.31	0.33
Claude-3.5-Sonnet	0.73	0.86	0.49	0.63	0.19	0.44
Llama-3.1-8B	0.77	0.72	0.86	0.68	0.71	0.67
Llama-3.1-8B-Instruct	0.87	0.82	0.75	0.71	0.78	0.72
Llama-3.1-70B	0.86	0.79	0.58	0.73	0.57	0.58
Llama-3.1-70B-Instruct	0.74	0.79	0.75	0.79	0.53	0.38
Llama-3.3-70B-Instruct	0.80	0.76	0.64	0.71	0.40	0.41
Mixtral-8x7B	0.74	0.64	0.83	0.84	0.19	-0.03
Mixtral-8x7B-Instruct	0.54	0.58	0.41	0.37	0.10	0.08
Qwen2.5-72B	0.83	0.75	0.75	0.74	0.73	0.72
Qwen2.5-72B-Instruct	0.63	0.61	0.55	0.64	0.54	0.52

Table 6: Spearman Correlation between default goal and other goals in Experiment 1

944

946

947

950

951

956

957

962

963

964

965

968

969

970

971

972

973

974

975

976

978

979

981

983

987

988

991

B **Experiment 2 Methods**

B.1 Participants

We recruited 156 participants through Prolific in the US or UK to take part in our open-ended response generation task. Participants were compensated at a rate of \$15 / hour following the approved IRB protocol at <University Anonymized>.

B.2 Stimuli

We first needed to elicit a large set of open-ended human responses to compare against the kinds of responses generated by LLMs. To do this, we recruited N = 156 participants through Prolific, located in the US or UK (compensated at a rate of \$15/hour) and gave them an open textbox to imagine what someone would say in the given scenario. Each participant was assigned 4 distinct scenarios out of the total set of 208 (see Figure 4 middle panel). We planned our sample size to collect at least 3 different responses for each scenario.

To verify comprehension, we began with three warm-up questions featuring different ratings, requiring participants to simply match visual ratings with their textual equivalent. All participants effectively matched visuals with text, though five participants each made one error out of three questions. We still included their responses after manually reviewing them and confirming their alignment with the ratings and contexts presented. To minimize response bias and create a more naturalistic experience, we interspersed filler scenarios among the main testing scenarios. While structured identically to testing scenarios, filler scenarios focused on opinions about *objects* rather than *people* (see Table 7 for examples). Each participant thus viewed a total of 8 scenarios (4 main testing scenarios and 4 filler scenarios). We controlled the presentation to ensure that each participant was presented with a series of distinct stories, with each of the 4 goals and 4 true-state ratings appearing exactly once.

Next, we needed to collect responses from LLMs for comparison. Instead of the multiple-choice task we gave in the previous section, Each model was presented with the same 208 scenarios as the human participants and was explicitly instructed to "keep your responses as short and concise as possible" to prevent excessively long answers. Each model generated one response per scenario with a temperature setting of $\tau = 0$, resulting in a total of 624 responses collected. We collected responses

from three LLMs: GPT-4o, Claude-3.5-Sonnet,	
and Llama-3.3-70B-Instruct.	

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1007

1012

1013

1014

1015

1016

1017

1018

1019

1020

B.3 Design

In the evaluation phase, we conducted a series of pairwise two-alternative forced choice comparisons, where human evaluators indicated which of a pair of responses they preferred for a given scenario. We included three kinds of comparisons:

- 1. Human vs. LLM preferences: Evaluators selected between human and LLM responses given identical scenarios, allowing us to understand which responses were preferred.
- 2. Goal Sensitivity: We compared responses gen-1004 erated for the original scenario (aligned-goal 1005 response) against those generated for scenar-1006 ios with different goals but identical ratings and contexts (misaligned-goal response). This 1008 comparison revealed preferences between re-1009 sponses with aligned versus misaligned com-1010 municative goals. 1011
- 3. Rating Sensitivity: We presented pairs consisting of responses generated for the original scenario (aligned-rating response) and responses generated with identical story and goal parameters but different ratings (misaligned-rating response). This comparison identified preferences between responses with aligned versus misaligned ratings.

B.4 Procedure

156 human evaluators are recruited from Prolific 1021 in the US or UK to take part in our evaluation 1022 task. Participants were compensated at a rate of 1023 \$15/ hour following the approved IRB protocol 1024 at <University Anonymized>. Each participant 1025 evaluated four different scenarios with five prefer-1026 ence questions per scenario: one trial comparing 1027 human vs. LLM responses, two trials assessing 1028 goal sensitivity within human and LLM sources, 1029 and two trials evaluating rating sensitivity for both 1030 sources. We ensured that each participant was pre-1031 sented with responses from distinct human sources 1032 and distinct LLM sources in each block, and each 1033 participant completed a total of 4 blocks consist-1034 ing of distinct scenarios with unique rating-goal 1035 combinations. To minimize potential confounds, 1036 we implemented several additional controls. First, 1037 we randomized both question order and response 1038 option order within scenarios to control for order 1039

Figure 4: Pipeline for comparing open-ended polite speech generation in humans and LLMs. Our study consists of two stages: an initial stage where we elicit responses for a variety of scenarios and a second stage where we ask a naive group to evaluation which of these responses they prefer.

Scenario	Rating: 0/3 hearts, Goal: Both	Rating: 2/3 hearts, Goal: Informative
Imagine that Jenny wrote a poem, but she didn't know how good it	Human : You are talented. Put in more effort and it will be superb.	Human: I think your poem has merit and it's pretty good.
was. Jenny approached Karen, who knows a lot about poems, and asked " <i>How was my poem</i> ?"	LLM: I loved the effort you put into your poem and I think there's a lot of potential, but the rhythm and flow could use some improvement.	LLM: I liked most of it but there's definitely room for improvement in a few places.
Imagine that John wanted to get Josh's opinion about a video game they just played. After Josh finished the game, John asked, " <i>What did you think?</i> "	Human: I didn't really care for it, but I had fun hanging out with you.	Human: It was a really fun video game.

Table 7: Examples of open-ended human and LLM responses in Experiment 2.

1040effects. We also inserted a transition page between1041blocks to reduce carryover effects. For the goal sen-1042sitivity and rating sensitivity comparisons, LLM1043comparisons were constrained to pairs of responses1044from the same model to control for model-specific1045variations in generation style.

B.5 Prompting strategy

In the open-ended response case, we keep the whole scenario the same as in the multi-choice version, just omitting the answer options.

Open-ended response generation system1050prompt: In the scenario, a person gave1051some performance and asked for another1052person's opinion on the performance. \n1053The person's feelings in the scenario are1054

1046

1047

1048

1049

Figure 5: (A) Distributions of how often the "was/wasn't terrible/bad/good/amazing" template studied by Yoon et al. (2020) was spontaneously produced by participants under each goal and rating. (B) How often responses use negation as a strategy among the responses that apply the Yoon et al. (2020) format under each goal and rating.

Groups	Observed JSD	Null Means
preferred vs. non-preferred	0.013	0.009
human vs. LLM	0.175	0.058
both vs. informative	0.134	0.081
both vs. social	0.195	0.096
both vs. default	0.136	0.088
informative vs. social	0.231	0.100
informative vs. default	0.134	0.093
social vs. default	0.166	0.106
0 hearts vs. 1 heart	0.119	0.087
0 hearts vs. 2 hearts	0.167	0.089
0 hearts vs. 3 hearts	0.227	0.094
1 heart vs. 2 hearts	0.129	0.088
1 heart vs. 3 hearts	0.225	0.093
2 hearts vs. 3 hearts	0.191	0.095

Table 8: JSD with word frequency counting distribution, all the p-values are < .001

shown on a scale of 0 to 3 hearts. \n 0 heart, the lowest rating, means the person does not like the performance at all, and 3 hearts, the highest rating, means the person likes it a lot.\n Please read the scenario carefully and answer the question in a complete sentence.\n Please keep your responses as short and concise as possible!\n Please only give the sentence-response without any other words!

C Additional details of linguistic analysis

1067 C.1 JSD tables

1055

1056

1057

1058

1059

1062

1063

1064

1066

1069

See the complete JSD scores of word-frequency distribution and politeness-strategy distribution at

Groups	Observed JSD	Null Means
gpt4.1 vs. claude3.7 labels	0.045	0.004
gpt4.1 vs. golden labels	0.073	0.004
claude3.7 vs. golden labels	0.026	0.003
preferred vs. non-preferred golden labels	0.008 (p = 0.087)	0.006
human vs. LLM response golden labels	0.023	0.006
both vs. informative	0.060	0.011
both vs. social	0.107	0.011
both vs. default	0.033	0.011
informative vs. social	0.180	0.013
informative vs. default	0.0197 (p = 0.006)	0.0125
social vs. default	0.131	0.013
0 hearts vs. 1 heart	0.049	0.011
0 hearts vs. 2 hearts	0.115	0.012
0 hearts vs. 3 hearts	0.286	0.013
1 heart vs. 2 hearts	0.079	0.011
1 heart vs. 3 hearts	0.263	0.012
2 hearts vs. 3 hearts	0.115	0.012

Table 9: JSD with politeness strategy frequency counting distribution, all the p-values are < .001 unless specified. We checked the agreement between gpt-4.1, claude-3.7-sonnet, and golden labels and found out the differences are quite significant (p < .001 - nontrivial). We compared the golden-labeled politeness strategies between human and LLM responses; we all use the golden-labeled politeness strategies for comparisons between goals and ratings.

1070

1071

1072

C.2 Text classification with SBERT embeddings

To analyze if there are high-dimensional features1073that differentiate each group beyond simple sta-
tistical analysis, we trained several simple classi-
fiers using SBERT (Reimers and Gurevych, 2019),
applying the pretrained sentence transformer "all-
MiniLM-L6-v2" to generate the response embed-
dings. Four different classifiers were implemented:1073

1112 1113 1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125 1126

1127

1128

1129

1130

logistic regression, random forest, SVM, and a simple MLP with a hidden layer size of 100 using scikit-learn. We then analyzed the results using the best-performing model among these four approaches.

Our analysis revealed that predicting between preferred and non-preferred responses is challenging, with performance only slightly above chance at approximately 54% across F1 score, recall, and precision metrics. In contrast, identifying human versus LLM responses was significantly more predictable, with the classifier reaching about 83% accuracy across the same metrics.

When examining the four goals comparison, we observed varying levels of predictability. Considering all responses collectively (both human and LLM), the classifier achieved approximately 60% performance in distinguishing between the four goals. Interestingly, when analyzing LLM responses in isolation, predictability increased to approximately 70%. whereas focusing solely on human responses, predictability decreased to around 40%. This suggests that LLM responses contain more distinctive patterns associated with different communicative goals compared to human responses, which exhibit greater variability in their approach to achieving the same goals.

C.3 Politeness annotation process

To handle the cases where a single politeness marker can be categorized under different politeness strategies, we allow the LLMs to assign up to three strategies per marker. Given our observation that both humans and LLMs often mix different politeness strategies in a single response, we also instruct the LLMs to identify all markers they consider reasonable. In the system prompt, we provided a comprehensive list of politeness strategies mainly from Brown and Levinson (1987) framework with additional ones from Danescu-Niculescu-Mizil et al. (2013). (see Appendix C.4). By providing a predefined, finite list of politeness strategies, we hope to unify the distribution of politeness strategies that these two LLMs can choose from and make their annotation results comparable. Since LLMs can make mistakes and often hallucinate (Xu et al., 2025; Huang et al., 2025) We manually inspected every single one of the labels annotated by the two LLMs, observing differences between the annotation results, and re-annotated any that we thought were not correctly labeled by the LLMs based on the definitions and examples

provided in the two frameworks as the golden labels.

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

There are still many cases where different researchers can label differently; our golden labels are just an instantiation of how we think what politeness strategies should be attached to politeness markers. We also note that "negation" is not considered as a specific politeness strategy in these two works, and thus we do not include it as a politeness strategy in the provided comprehensive list; we consider it as a "positive politeness - avoid disagreement" or "off-record - understate" when it appears based on the contexts. Throughout manual inspection, we indeed found out LLMs sometimes are not consistent with their labels - the same words can be labeled differently in different responses under the same goal and rating, and they sometimes hallucinate and give politeness strategies that are not in the provided list.

For the comprehensive list of politeness strategies provided in the annotation system prompt, there are several things to notice. First, the list is a combination of Brown and Levinson (1987) and Danescu-Niculescu-Mizil et al. (2013). We chose the list of politeness strategies from Brown and Levinson (1987) because their classic framework covers nearly all politeness strategies and is still widely used and adopted in most current work on politeness. The list of politeness strategies shown in Danescu-Niculescu-Mizil et al. (2013) is mainly adopted from Brown and Levinson (1987), with some additional strategies based on some other widely used politeness phenomena and literature. We believe that by combining the two lists, we can obtain a comprehensive set of all widely used politeness strategies in language.

A note on combining the two lists: if there is any disagreement between the two works, we follow the categorization in Brown and Levinson (1987). Specifically, we consider Deference a negative politeness strategy, in line with Brown and Levinson (1987), whereas in Danescu-Niculescu-Mizil et al. (2013), it is categorized as a positive strategy.

Our manual inspection of gpt-40 and claude-3.7-sonnet and golden-label generation follows several principles:

• We follow the definitions of politeness strate-1177 gies provided in their respective frameworks 1178 and annotate all politeness markers in a given 1179 response. A single politeness marker may cor-1180 respond to multiple politeness strategies, and 1181 1182people often mix different strategies within a1183single response. We consider all politeness1184markers and label each one with up to three1185of the most significant politeness strategies.

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1201

1202

1203

- For words that are not clearly significant enough to be considered politeness markers—cases where they could be interpreted as either common words or politeness markers—we simply accept whatever the LLMs produce.
- Could/Would are counterfactual modals, which are widely used in polite speech. They are not considered in Brown and Levinson (1987), but are included in Danescu-Niculescu-Mizil et al. (2013). In our manual labeling process, we always mark them as counterfactual modals and additionally label them with other relevant politeness strategies, such as hedging, when appropriate.

C.4 LLM annotation prompt

The following whole section is the complete system prompt:

1204You are an expert in the study of1205human polite language use, with extensive1206knowledge of the relevant literature and1207the various politeness strategies people1208employ in everyday conversation.

1209Please follow my instructions to1210help extract, annotate, and categorize1211politeness markers used in the response1212of each scenario.

1213In each scenario, Person A asks Person1214B for their opinion on A's performance.1215Person B's true feelings are represented1216on a scale of 0 to 3 hearts as the rating,1217where 0 hearts means they did not like the1218performance at all, and 3 hearts means1219they liked it very much.

1220In the question, please pay attention to1221the communicative goal mentioned (either1222to be informative, to make person A feel1223good, to do both, or to serve as the1224default with no specific goal).1225Your tasks are to:

1226 1. Read the whole scenario setup carefully. pay attention to the rating and the communicative goal 1228 in the question. Then, identify 1229 and annotate the specific word(s) or 1230

phrase(s) in Person B's response that 1231 function as politeness markers. 1232

each politeness marker 2. Categorize 1233 comprehensive list of using the 1234 politeness strategies provided 1235 below, specifying both the 1236 (positive category politeness, 1237 negative politeness, off-record, 1238 bald-on-record) and the corresponding 1239 specific politeness strategy. 1240

1241

1242

1243

1244

1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

Please present your answer in the following format for each politeness marker WITHOUT ANY additional text or explanation:

Politeness marker: [the specific 1245 word(s) or phrase(s)] 1246 Politeness strategy-1: [category 1247 + specific politeness strategy] Politeness strategy-2: [category 1249 + specific politeness strategy] 1250 (if applicable) 1251 Politeness strategy-3: [category 1252 + specific politeness strategy] 1253 (if applicable) 1254

(Repeat the above for each politeness marker found in the response)

Below is the comprehensive list of politeness strategies with examples for each strategy. The politeness markers, e.g., the specific word(s) or phrase(s) used in each strategy's example, are shown in parentheses.

Please pay attention to the usage of could/would or similar words in the following list.

I. Positive Politeness Strategies

1. Gratitude

- Example 1: "Thank you so much for 1270 your help!" (thank you) 1271
- Example 2: "I really appreciate 1272 your kindness." (I really 1273 appreciate) 1274
- 2. Greeting (social approach) 1275

1276 1277	 Example 1: "Hi there! Could you help me out?" (Hi there) 	 Example 1: "It's beautiful today, isn't it?" (isn't it?) 	1319 1320
1278	• Example 2: "Good morning! How are	 Example 2: "This solution seems ideal right?" (right?) 	1321
1275	3. Greeting (social approach)	10. Avoid disagreement	1323
		• Example 1. "Yes that might work	1394
1281	• Example 1: "H1 there! Could you	but also consider	1325
1202	• Exemple 2: "Cood manning! How and	might work)	1326
1283	• Example 2: "Good morning! How are	• Example 2: "I see your point,	1327
1204	you today: (dood morning)	though perhaps" (I see your	1328
1285	4. Positive Lexicon (positive sentiment,	point)	1329
1286	optimism)	11. Presuppose/assert common ground	1330
1287	• Example 1: "Wow, that's wonderful	• Example 1. "You know how much we	1331
1288	news!" (wonderful)	both value honesty " (we both)	1332
1289	• Example 2: "I'm thrilled about	• Example 2: "We both know how	1333
1290	your promotion." (thrilled)	difficult this can be " (We both	1334
1201	5 Notice attend to hearer's interests	know)	1335
1292	wants. needs		
1000	• Example 1. "You seem stressed-see	12. Joke	1336
1293	• Example 1. Tou seem stressed-can I assist?" (You seem stressed)	• Example 1: "If you fix this bug,	1337
1205	• Example 2: "Vou must be tired	I'll bake you cookies!" (I'll	1338
1295	nlease take a rest " (You must	bake you cookies)	1339
1297	be tired)	• Example 2: "Careful, your	1340
	·	brilliance is showing: (your	1341
1298	6. Exaggerate interest, approval,	billiance is showing)	1342
1299	sympathy	13. Assert speaker's knowledge of	1343
1300	• Example 1: "That's the best	hearer's wants	1344
1301	presentation I've ever seen!"	• Example 1: "Since I know you like	1345
1302	(the best)	chocolate, here's a cake." (Since	1346
1303	• Example 2: "Your idea	I know you like chocolate)	1347
1304	is absolutely fantastic!"	• Example 2: "Knowing you love	1348
1305	(absolutely fantastic)	adventure, I booked a trip."	1349
1306	7. Intensify interest in hearer	(Knowing you love adventure)	1350
1307	• Example 1: "I traveled across	14. Offer, promise	1351
1308	town just to see you!" (just to	• Example 1: "I'll take care of	1352
1309	see you)	that for you tomorrow." (I'll	1353
1310	• Example 2: "I've been eagerly	take care)	1354
1311	waiting to hear your story."	• Example 2: "If you're busy, I	1355
1312	(eagerly waiting)	promise to handle it myself." (I	1356
1010	0 lles in mean identity membrane	promise)	1357
1313	o. Use in-group identity markers	15. Be optimistic	1358
1314	• Example 1: "Hey mate, can you	• Evampla 1. (T'm and a set	1000
1315	give me a hand?" (mate)	• Example 1: "I'm sure you can	1359
1316			1.500
	• Example 2: "Buddy, I need your	• Example 2: "You'll definitely	1004
1317	• Example 2: "Buddy, I need your advice on something." (Buddy)	• Example 2: "You'll definitely manage to finish this in time "	1361
1317 1318	 Example 2: "Buddy, I need your advice on something." (Buddy) 9. Seek agreement 	 Example 2: "You'll definitely manage to finish this in time." (You'll definitely) 	1361 1362 1363

1364 16 1365	. Include both speaker and hearer (inclusive 'we')
1366	• Example 1: "Let's figure this out
1367	together." (Let's)
1368	• Example 2: "Why don't we start
1369	the project now?" (we)
1370 17	. Give or ask for reasons
1371	• Example 1: "Could you come with
1372	me? It'll be helpful." (It'll be
1373	helpful)
1374	• Example 2: "Why not join the
1375	group? You'd enjoy it." (You'd
1376	enjoy it)
1377 18	. Assume reciprocity
1378	• Example 1: "You helped me last
1379	time, now it's my turn." (now
1380	it's my turn)
1381	• Example 2: "I lent you my
1382	notes earlier-can I borrow yours
1383	today?" (I lent you my notes
1384	
1385 19	. Give gifts to hearer (sympathy,
1386	understanding, cooperation)
1387	• Example 1: "You've been working
1388	hard; here's a small gift."
1389	(here's a small gift)
1390	• Example 2: "Here, take this
1391	coffee-you deserve a break." (you
1392	deserve a break)
1393 II	. Negative Politeness Strategies
1394 1	. Apologizing
1395	• Example 1: "Sorry to disturb you,
1396	but I have a question." (Sorry)
1397	• Example 2: "I apologize for
1398	interrupting your meeting." (I
1399	apologize)
1400 2	. Please (sentence-medial polite form)
1401	• Example 1: "Could you please send
1402	me the document?" (please)
1403	• Example 2: "Would you please
1404	consider my suggestion?"
1405	(please)
1406 3	. Be conventionally indirect

•	Example	e 1:	"Could	you j	poss	ibly	14	07
	close	the	door?"	(Co	uld	you	14	08
	possib	ly)					14	09
•	Example	e 2:	"Wou]	Ld vo	Su	mind	14	10

handing me the pen?" (Would you 1410 mind) 1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1441

1442

1443

1444

1449

- Example 3: "By the way, do you know the time?" (By the way)
- Example 4: "Oh, by the way, did you finish the report?" (Oh, by the way)

4. Question, hedge

- Example 1: "Perhaps we could reconsider the deadline?" (Perhaps)
- Example 2: "Maybe you might find this helpful?" (Maybe)
- Example 3: "I suggest we might consider other options." (might consider)
- Example 4: "I think it's possibly better this way." (I think, possibly)

5. Be pessimistic

- Example 1: "I don't suppose you could spare a moment?" (I don't suppose)
- Example 2: "You probably wouldn't want to help, would you?" (probably wouldn't want)

6. Minimize the imposition

- Example 1: "I just need a quick 1438 moment of your time." (just need 1439 a quick moment) 1440
- Example 2: "This will take only a second, I promise." (only a second)

7. Give deference

- Example 1: "Professor, could you 1445 clarify this point?" (Professor) 1446
- Example 2: "Excuse me, sir, may 1447 I interrupt?" (Excuse me, sir) 1448

8. Impersonalize speaker and hearer

• Example 1: "It seems this task 1450 needs attention." (It seems) 1451

1452 1453 1454	 Example 2: "There appears to be a misunderstanding." (There appears) 	 Example 1: "Oh no, I forgot my wallet!" (forgot my wallet) - hint to pay for them 	1495 1496 1497
1404		• Example 2: "My phone just died."	1498
1455	9. State the FIA as a general rule	(phone just died) – hint to borrow	1499
1456	• Example 1: "Visitors are	a phone	1500
1457	requested not to use cell		
1458	phones." (Visitors are	3. Presuppose	1501
1459	requested)	• Example 1: "I cleaned it again	1502
1460	• Example 2: "Eating is not allowed	today." (again) – presupposes	1503
1461	in the library." (is not allowed)	someone else didn't	1504
1/60	10 Nominalize	• Example 2: "Did you check the	1505
1402		oven?" (Did you check) – implies	1506
1463	 Example 1: "Your participation is 	concern or oversight	1507
1464	required." (participation)	1 Understate	1500
1465	• Example 2: "Submission of	4. Under State	1500
1466	your paper is expected soon."	• Example 1: "The movie was not	1509
1467	(Submission)	exactly thrilling." (not exactly	1510
1468	11 Go on record incurring a debt	thrilling)	1511
1400		• Example 2: "His speech was	1512
1469	• Example 1: "I'd greatly	somewhat unclear." (somewhat	1513
1470	appreciate it if you helped	unclear)	1514
1471	me." (I'd greatly appreciate it)	5 Overstate	1515
1472	• Example 2: "I'll owe you one if	J. Overstate	1515
1473	you can cover my shift." (I'll	• Example 1: "I've waited forever	1516
1474	owe you one)	for your reply!" (waited forever)	1517
1475	12. Counterfactual modal forms	• Example 2: "I'm starving!"	1518
1476	(could/would)	(starving)	1519
1/77	• Example 1. "Could you assist me	6. Tautologies	1520
1478	with this?" (Could you)		
1/70	• Example 2: "Would you mind	• Example 1: "Business is	1521
1479	checking this for me?" (Would you	business. (Business is	1522
1481	mind)	business)	1523
		• Example 2: "It is what it is."	1524
1482	13. Indicative modal forms (can/will)	(It IS what It IS)	1929
1483	• Example 1: "Can you help me with	7. Contradictions	1526
1484	these files?" (Can you)	• Example 1. "It's good but at the	1507
1485	• Example 2: "Will you be able to	same time not good " (good but	1527
1486	come by later?" (Will you)	not good)	1520
		• Example 2: "I'm happy and not	1520
1487	III. Off-Record (Indirect) Strategies	happy about this " (happy and not	1530
1488	1 Give hints	happy about this. (happy and hot	1532
1400		heppy	1002
1489	• Example 1: "It's chilly in	8. Be ironic	1533
1490	here" (chilly in here) –	• Example 1: "Lovely day we're	1534
1491	nint to close the window	having!" (Lovely day) - during	1535
1492	• Example 2: "I'm thirsty." (I'm	bad weather	1536
1493	thirsty) - hint to offer a drink	• Example 2: "That went well!"	1537

1539	9. Use metaphors	2. Direct commands (imperatives)	1580
1540	• Example 1: "He's got a heart of	• Example 1: "Stop right now!"	1581
1541	stone." (heart of stone)	(Stop)	1582
1542	• Example 2: "She's a ray of	• Example 2: "Bring it to me	1583
1543	sunshine." (ray of sunshine)	immediately." (Bring it)	1584
1544	10. Rhetorical questions	3. Sentence-initial imperative forms	1585
1545	• Example 1: "How many times must	("Please" start-less polite)	1586
1546	I tell you?" (How many times)	• Example 1: "Please move out of my	1587
1547	• Example 2: "Do I look like I'm	way." (Please move)	1588
1548	joking?" (Do I look)	 Example 2: "Please finish your work guickly." (Please finish) 	1589 1590
1549	11. Be ambiguous		
1550	• Example 1: "Something feels off	4. Sentence-initial second-person	1591
1551	about this " (feels off)	statements (less polite)	1592
1551	• Example 2: "It coome unucual	• Example 1: "You need to fix this."	1593
1552	• Example 2: It seems unusual	(You need to)	1594
1000	Someriow (Seems unusual)	• Example 2: "You've misunderstood	1595
1554	12. Be vague	me." (You've misunderstood)	1596
1555	• Example 1: "I'm a bit upset." (a	5. Factuality (direct assertions. less	1597
1556	bit)	polite)	1598
1557	• Example 2: "I kind of disagree."		4 500
1558	(kind of)	• Example 1: "Actually, you did	1599
		incorrectly)	1601
1559	13. Over-generalize	• Evennle 2: "The truth is you	1001
1560	 Example 1: "Everyone knows it's 	• Example 2: "The truth is you failed to deliver " (you failed	1602
1561	not true." (Everyone knows)	to deliver)	160/
1562	• Example 2: "Nobody likes that."		1004
1563	(Nobody)	 Negative lexicon (negative sentiment, impolite) 	1605 1606
1564	14. Displace hearer	po00)	1000
1565	• Example 1. "I wish someone would	• Example 1: "You're always messing	1607
1566	help " (someone)	things up!" (always messing	1608
1500	• Evample 2: "It'd be great if	things up)	1609
1568	someone cleaned up " (someone)	• Example 2: "If you're going to	1610
1500	someone creanca ap. (someone)	accuse me" (accuse me)	1611
1569	15. Be incomplete (ellipsis)	C.5 A comprehensive list of politeness	1612
1570	• Example 1: "If only you knew "	strategies with examples	1613
1571	(If only you knew)	See Table 10 for a comprehensive list of politeness	1614
1572	• Example 2: "Well if you could	strategies used in both human and LLM responses.	1615
1572	iust " (if you could just)	The examples and strategies shown are based on	1616
1010		golden labels from our collected responses.	1617
1574	IV. Bald-on-Record Strategies		
1575	1. Direct questions/statements		
1576	 Example 1: "What are you doing?" 		
1577	(What are you doing?)		
1578	• Example 2: "Where did you put it?"		
1579	(Where did you put it?)		

Category	Politeness Strategy	Example
	1. Assert speaker's knowledge of hearer's wants	<i>I know</i> you are up to the challenge!
	2. Avoid disagreement	Pretty decent for a beginner
	3. Be optimistic	You can even make them better next
		time!
	4. Exaggerate interest, approval, sympathy	Your dance greatly exceeded all ex-
		pectations.
	5. Give gifts to hearer	I am so proud of you.
Positive	6. Give or ask for reasons	I have tasted some really good cakes,
Politeness		and yours
	7. Gratitude	I'm so grateful
	8. Greeting	Hey, I read your review
	9. Include both speaker and hearer	Let's go through it together
	10. Intensify interest in hearer	You were born to be on stage
	11. Notice, attend to hearer's interests	I can see you put in lots of effort
	12. Offer, promise	Let me know if you need any tips.
	13. Positive lexicon	It was absolutely amazing!
	14. Presuppose/assert common ground	With other artists of your caliber
	15. Seek agreement	Suns your first time baking?
	To. Use III-group Identity markers	Tour app is pretty good, <i>Henry</i> !
	17. Apologizing	I didn't like it, sorry!
	18. Be conventionally indirect	If you would like
	19. Counterfactual modal forms	Could/Would you
Negative	20. Give deference	<i>In my expert opinion</i> , your painting is
Politeness		Tabulous.
	21. Impersonalize speaker and nearer	<i>There are</i> a few places to improve.
	22. Minimize the imposition 23. Nominaliza	I have a jew suggestions (0 +social)
	25. Nommanze	I would not be the best person to eval-
	24 Question hadge	Maybe try adding some different fla
	24. Question, neuge	voring ingradiants payt time?
	25. Be ironic	It was horrible, my eyes are bleeding.
	26. Be vague	It was interesting (0-rating case)
	27. Contradictions	It was great, nowever, it needs im-
		provement Variable de la constante de la constante
Off-Record	28. Displace hearer	You looked so confident and elegant!
	20 Cive association alway	(when commenting on performance)
	29. Give hints	It could be better if you adjusted the
	50. Give mints	sweetness!
	31 Overstate	The cookies tasted great $(1 \pm social)$
	32 Presuppose	Pretty decent for a beginner
	33 Understate	It was not good (0-hearts rating)
	34. Use metaphors	Your singing was like music to my
		ears!
	35 Direct commands	Try practicing with precise measure-
	55. Direct commands	ments
Bald-on-	36. Factuality	I didn't like the cookie at all
Record	37. Negative lexicon	It was terrible
	38. Sentence-initial 2nd-person statements	You need to work on that.
	39. Sentence-initial imperative forms	Please for gods sake improve on these
	*	areas

Table 10: A comprehensive list of politeness strategies with examples from our collected responses. We consider all the politeness strategies and politeness markers in the golden annotation results.