

Do they mean ‘us’? Interpreting Referring Expressions in Intergroup Bias

Anonymous ACL submission

Abstract

The variations between in-group and out-group speech (intergroup bias) are subtle and could underlie many social phenomena like stereotype perpetuation and implicit bias. In this paper, we model the intergroup bias as a *tagging task* on English sports comments from forums dedicated to fandom for NFL teams. We curate a unique dataset of over 6 million game-time comments from opposing perspectives (the teams in the game), each comment grounded in a non-linguistic description of the events that precipitated these comments (live win probabilities for each team). Expert and crowd annotations justify modeling the bias through tagging of implicit and explicit referring expressions and reveal the rich, contextual understanding of language and the world required for this task. For large-scale analysis of intergroup variation, we use LLMs for automated tagging, and discover that some LLMs perform best when prompted with *linguistic descriptions* of the win probability at the time of the comment, rather than numerical probability. Further, large-scale tagging of comments using LLMs uncovers **linear variations in the form of referent** across win probabilities that distinguish in-group and out-group utterances.

1 Introduction

Social bias in language is generally studied by identifying undesirable language use towards a specific demographic group (Kaneko and Bollegala, 2019; Sheng et al., 2019; Sap et al., 2020; Webson et al., 2020; Pryzant et al., 2020; Sheng et al., 2020); However, we can enrich our understanding of bias in communication by understanding it as differences in behavior situated in social relationships. Intergroup bias is the social bias stemming from the intergroup relationship between the speaker and target reference of an utterance. (Maass et al., 1989; Maass, 1999). Govindarajan et al. (2023a) modeled intergroup relationships (in-group and out-group) and interpersonal emotions in interpersonal

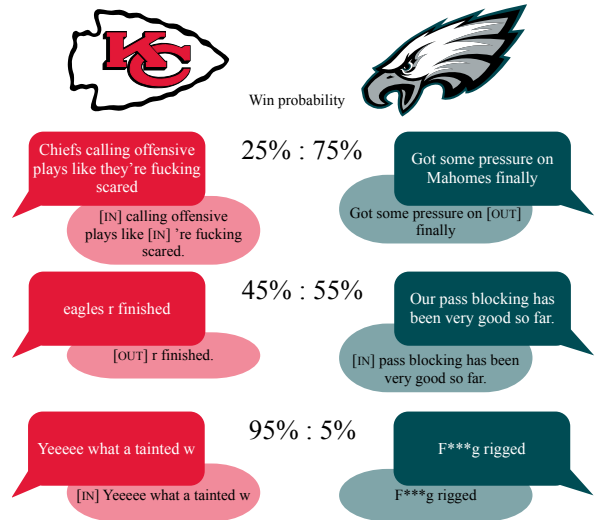


Figure 1: We construct a parallel language corpus of comments from NFL team subreddits, grounding each comment in the live win probabilities. We then tag relevant entities in each comment with intergroup tags using annotators and LLMs.

English language tweets at the utterance level, finding systematic interactions between these two parameters. While neural models based on LLMs can be trained to discriminate in-group and out-group utterances, causal probing of these models was inconclusive (Govindarajan et al., 2023b) and still left major questions unanswered:

- How does language systematically change when referring to an individual in one’s in-group versus their out-group?
- How does the state-of-the-world influence this systematic variation?

In this work, we take major steps towards answering these questions, using a task architecture in the classical NLP pipeline (Manning et al., 2014). Earlier work in the Linguistic Intergroup Bias (LIB) hypothesis (Maass, 1999) focused exclusively on the predicate, and the bias was described using an ad-hoc lexical categorization system (Semin and Fiedler, 1988). However, the **form of referencing** the in-group or out-group can reveal subtle

064 biases as well. Consider the comments in Figure 1,
065 taken from the dataset we describe in this paper.
066 Commenters refer to the in-group and out-group by
067 name, sub-groups, pronouns as well as implicitly —
068 sometimes they choose not to refer to either group
069 at all. How does the intergroup bias manifest in
070 referent forms?

071 To answer this question, we introduce a new
072 dataset of interpersonal language — comments
073 from game threads on online forums dedicated to
074 fandoms for teams in the National Football League
075 (NFL). Through careful data curation, we construct
076 a parallel corpus of sports comments, with com-
077 ments from fans of both teams in a game, *aligned*
078 in time and *grounded* in win probabilities (WP). By
079 focusing on referring expressions, we can formu-
080 late investigating the intergroup bias as a **tagging**
081 task: given a comment, the group affiliation of the
082 writer, and the state-of-the-world, return a *tagged*
083 *comment* with appropriate referring expressions
084 tagged as [IN], [OUT] or [OTHER] (see Figure 1).
085 Annotation and preliminary analysis reveal that the
086 form of the referent that speakers use when refer-
087 ring may have systematic intergroup variations.

088 We train Large Language Models (LLMs) to
089 automate large-scale tagging of our dataset, and ex-
090 amine their performance on our task. We find that
091 few-shot performance on GPT-4o is boosted using
092 linguistic descriptions of win probabilities; fine-
093 tuned Llama-3 models performed better, although
094 incorporating WP had little effect. Using our best
095 performing model to tag 100,000 comments from
096 our raw dataset, we discover two striking linguistic
097 behaviors:

- 098 1. Higher the win probability for the *in-group*,
099 the more likely commenters are to **abstract**
100 **away** from referring to the in-group. This
101 trend is remarkably linear across win proba-
102 bilities for all types of in-group references.
- 103 2. References to out-groups by commenters are
104 rarer than in-group references, and remain sta-
105 ble over all win probabilities for the in-group.

106 These findings add much needed color to the LIB
107 hypothesis — natural language is productive, and
108 commenters can express their (implicit) intergroup
109 bias in different ways. This work also lays the foun-
110 dation for future explorations of other intergroup
111 variations (in event descriptions, for example) in
112 sports-talk and other domains. We share all our
113 code, models, and data online.

2 Background and Related Work 114

Intergroup bias Linguistic Intergroup Bias (LIB) theory (Maass et al., 1989; Maass, 1999) hypothesizes that stereotypes are transmitted and persist in communication through systematic linguistic asymmetry — socially desirable in-group behaviors and socially undesirable out-group behaviors are encoded at a higher level of abstraction. The LIB has been reproduced in psychological experiments and analyses (Anolli et al., 2006; Gorham, 2006); it has also been used as an indicator for a speaker’s prejudicial attitudes (Hippel et al., 1997), and racism (Schnake and Ruscher, 1998).

Govindarajan et al. (2023a,b) take inspiration from the LIB at large to study intergroup bias as a general phenomenon in online language use. While they find regularities in its variation with emotion that neural models can ‘learn’ to identify in-group and out-group utterances more accurately than humans, probing experiments fail to describe human-observable intergroup variations in language. This work studies a much larger dataset than in their work, and by modeling the bias as a tagging task to referents, we discover characteristic lexical variations at scale that complement LIB findings.

Sports language Language use in the domain of sports has been a rich source of analyses and studies within computational linguistics, including from the perspective of quantifying *social biases*. Merullo et al. (2019) studied commentator racial biases in descriptions of football players, reaffirming previous findings illustrating clear differences in terms of sentiment descriptions (white players were more likely to be described as intelligent), and name itself (white players were more likely to be referred to by their first name). Zhang et al. (2019) focused on one aspect of language usage among (and between) fans of NBA teams: intra-group behavior with and without social contact with the out-group. They find that fans with intergroup contact are more likely to use negative language — they were more polarized than before.

Our work differs from previous work in two major ways. Firstly, we focus on the intergroup bias — how do fans talk about their team (in-group), versus the opponent (out-group)? Secondly, this paper **grounds** the analysis of intergroup bias in numerical descriptions of the state-of-the-world. The state-of-the-world in a sports game at any moment can be described using the scoreboard, thus

providing grounding for utterances follow. Non-linguistic, numerical descriptions of the events that precipitate an utterance overcome the drawbacks of using ad-hoc, proximally derived metrics like *social desirability* (in LIB) or *affect* (in Govindarajan et al. (2023b)) as an axis to study linguistic variation. As we shall describe in §3, sports games, and in particular NFL games, are rich with statistical information amenable to describe the state-of-the-world on a well-calibrated numerical scale.

3 Data & Annotation

3.1 Dataset

Data & Preprocessing Our new dataset of intergroup language comes from Reddit — specifically subreddits dedicated to fandoms for each of the 32 teams in the NFL. During the NFL season, each subreddit has *game threads* — posts created by moderators on which fans can comment in tandem with the live game involving their team. Crucially, since every subreddit has their own thread, we effectively have a **parallel intergroup language dataset**; two teams and their fans commenting on *the same game events*. Further, these subreddits are dedicated to individual team fandoms, so we can fairly assume that the team the subreddit represents is the in-group for all commenters.¹

We focus on all completed games from the 2021–22 and 2022–23 NFL seasons, and attempted to scrape all comments from the game threads for both teams involved in every game. Within comments from game threads, we filtered it down to comments that happened during active game-time, and removed comments that were only URL links. Overall, our raw data has over 6 million comments from 768 game threads on 32 subreddits, grounded in 491 NFL games.

Grounding football comments American Football has some attractive features as a sport considering that our interest is in the *language surrounding the events* in a game — it is highly strategic, and outcomes are heavily dependent on a coach’s strategies and plays in a (relatively) small number of discrete events (called *plays*, Pelechrinis and Papalexakis, 2016). The state-of-the-world at any moment in a football game is determined by a variety of factors — seconds remaining in half (and game), yard line, score differential, down, yards to go, home

advantage, timeouts remaining, betting odds lines from Vegas, and so many more (Horowitz et al., 2017; Yurko et al., 2018). Baldwin (2021) modeled the **Win Probability** (henceforth **WP**) of a team at any point during the game using a decision trees over the aforementioned features, building a well-calibrated model with low error. We chose WP as a succinct, non-linguistic description of the events preceding an utterance.

Using the `nflFastR` (Baldwin, 2021) package, we can obtain WPs for individual plays in each game, as well as the time of completion of a play. Combined with the timestamps at which the comments were submitted (obtained from the Reddit API), we build our parallel corpus of intergroup language grounded in win probabilities. The WP cleverly models the complexities of a real-world sporting event into one number that accurately models how **desirable** the state-of-the-world is to the in-group (see Figure 1).

3.2 Tagging

As we motivated in §1 and Figure 1, tagging references to entities enables us to perform analyses at scale *and* discover individual lexical variations. Consider the following examples:

- (1) a. Rams are gifting us a chance to win and we can’t take advantage. The f***!!!!
- b. if the ravens and chiefs beat these dudes by double digits then damn it so should we!

Even without contextual information about the game for the above comments, we see *multiple* readily identifiable references to the in-group and out-group, within the same utterance. The words or phrases that refer to relevant individuals can now be tagged with in-group ([IN]) or out-group ([OUT]) For instance, (1) would be tagged thus:

- (2) a. [OUT] are gifting [IN] a chance to win and [IN] can’t take advantage. The f***!!!!
- b. if [OTHER] and [OTHER] beat [OUT] by double digits then damn it so should [IN]!

We define the in-group ([IN]) as the team the commenter supports (and its fans), and the out-group ([OUT]) as the opponent in that particular game (and its fans). The spans ‘the ravens’ and ‘chiefs’ in (1-b) are clearly not a reference to the in-group nor the opponent of the game. However, they are a reference to *a group of interest in this domain* — another NFL team and/or its fans. We consider these references to be [OTHER], and a special case of out-group references.

¹Note that we focus on language of online commenters (fans) on Reddit, not *commentators* for the game.

261 Sometimes, the references to the in-group, out- 311
262 group or other are not explicit. However, we can 312
263 infer based on common-sense reasoning that the 313
264 comment as whole, or a sentence in the comment, 314
265 is **implicitly referring** to a relevant group: 315

266 (3) What a conservative play call 316

267 There is no explicit word/phrasal reference to any 317
268 team in the above comment. However, it is clear in 318
269 context (the fan’s team is losing, with WP of 9%) 319
270 that the commenter is referring to the in-group. 320
271 To facilitate these implicit annotations, we sen- 321
272 tence tokenize the comments in our dataset using 322
273 Stanza (Qi et al., 2020), append a sentence-level 323
274 token [SENT] before each sentence in every com- 324
275 ment in our dataset. If the sentence as a whole is 325
276 judged to implicitly refer to a relevant group, the 326
277 [SENT] token is replaced with the relevant tag. 327

278 3.3 Annotation 328

279 Annotators are presented with a comment from our 329
280 dataset, the source subreddit (team) for the com- 330
281 ment, the parent comment (if the comment is a 331
282 reply in a thread), and the live score at the time 332
283 of the comment. The task of tagging words and 333
284 phrases from comments in our dataset with inter- 334
285 group tags can be highly involved — in addition 335
286 to knowledge of American Football, commonsense 336
287 reasoning over the meaning of an utterance in con- 337
288 text of the live game, one needs knowledge of the 338
289 teams and its players. For instance, in (4), one 339
290 needs to know that the commenter supports the Sea- 340
291 hawks, and that there is a prominent player named 341
292 Wilson, to accurately tag in context that Wilson 342
293 indeed is an in-group reference. 343

294 (4) Our oline should start holding since apparently 344
295 it ’s okay now . Maybe Wilson can actually get 345
296 some time to throw . 346

297 Implicit annotations on the [SENT] token require 347
298 a higher bar of reference, since all comments are 348
299 about the game at hand and will involve both teams 349
300 to some extent. For example, we judge the fol- 350
301 lowing comments to not have explicit or implicit 351
302 references to any relevant groups of interest even 352
303 though they are about the game: 353

304 (5) a. Fair enough ! 354

305 b. winning cures all lmao 355

306 c. turning the game off , have a good day yall 356

307 In case it is impossible to verify an explicit or im- 357
308 plicit reference, annotators are instructed to not 358
309 highlight any parts of the comment. All annota- 359
310 tors were free to search the web for names or ex-

311 pressions they were unfamiliar with, as well as 312
313 refer to reports of the game to understand the utter- 314
315 ance completely, and accurately tag all references. 316
317 All annotation experiments were carried out using 318
319 the thresh. tools annotation interface (Heineman 320
321 et al., 2023). Annotators highlight spans within a 322
323 comment and select from one of 3 tags, and select 324
325 a confidence level from a five-point scale. 326

327 **Expert annotated dataset** We gather expert 328
329 annotations for constructing a ‘gold’ annotated 329
330 dataset to evaluate crowd annotations and modeling 330
331 moving forward. The first author of the paper an- 331
332 notated 1499 comments (randomly sampled from 332
333 all game-time comments) for intergroup references 333
334 based on a pre-defined, written protocol (described 334
335 in detail in Appendix B). 26.7% of comments were 335
336 judged to have no relevant intergroup reference, 336
337 and in the remaining comments, references to the 337
338 in-group (76.3%) vastly out-number references to 338
339 the out-group (14.6%) or other. This is not surpris- 339
340 ing, since these are comments from forums dedi- 340
341 cated to fandom of teams — people are much more 341
342 likely to talk about their team over the opponent. 342
343 We partition our gold dataset into a test set of 318 343
344 datapoints, and a training set of 1181 datapoints. 344

345 **Crowd annotation** To understand our dataset fur- 345
346 ther, we recruited three undergraduates to annotate 346
347 the test split of our expert dataset. Our goals were 347
348 to understand where disagreements arose, as well 348
349 as how and when knowledge of the events in the 349
350 game helped in disambiguating references. Annota- 350
351 tors were given similar instructions, and were free 351
352 to search the web and lookup statistics and reports 352
353 on the game in question. We found in pilot experi- 353
354 ments that the live-score was more interpretable to 354
355 humans than WP, influencing our choice to provide 355
356 that as context to annotators. 356

347 4 Preliminary Analysis 348

349 4.1 Annotation Analysis 350

351 **Inter-annotator agreement** Average Fleiss 351
352 κ (Fleiss, 1971) among crowd annotators is 0.69, 352
353 indicating moderate agreement. In addition to the 353
354 inter-annotator score, by counting exact-matches 354
355 and weighting partial matches between individual 355
356 crowd annotators and gold annotations, we 356
357 calculate an ‘accuracy’ score of 0.65 ± 0.005 . 357
358 This can be interpreted as a human ceiling for 358
359 performance on this task, and characterizes its 359
360 inherent subjectivity and difficulty. 360

Disagreement can be a signal Looking at the source of disagreements among annotators (and between crowd annotators and experts) can give us insights into the nature of the task itself (Atwell et al., 2022), as well as why differences in judgements of intergroup affiliation can come down to annotator biases or judgement given context. For example, annotators disagree sometimes on what counts as a ‘reference’:

(6) a. ...**Lambeau** has the second worst bathrooms .

b. Can’t do that against **an offense** this good.

Lambeau in (6-a) was judged by the expert annotator to be a reference to the out-group in this context — the opponent/out-group is the Green Bay Packers. However, to tag this referent [OUT], annotators would need to know or deduce that Lambeau Field is the Green Bay Packers stadium, and judge that this constitutes a relevant intergroup reference. Thus, disambiguating some references can be time-consuming and hard. *an offense* in (6-b) was judged by some annotators to refer to the out-group in context. However the generic nature of the referent lead other annotators to judge that this was an overall statement about the game, rather than an explicit reference.

Whether or not examples in (6) contain references to the in/out-group is not simply a consequence of the difficulty of our task, or the inability for annotators to transparently describe the mental state of commenters. Rather, we need to analyze them as possibly another subtle influence of the intergroup bias itself — demonstrated by questioning why commenters chose the forms in (6) rather than in (7), which convey the same meaning, and would be uncontroversial in annotation:

(7) a. ...**the Packer’s** stadium has the second worst bathrooms .

b. Can’t do that against **a Packers** offense. . .

4.2 Qualitative Analysis & Trends

Mereology of referring expressions Expert annotation revealed that commenters refer to groups of interest in a myriad of different ways. In the previous section, we liberally defined the annotation protocol for highlighting references to *individuals* in the in-group, out-group and other. Using insights from mereology (Varzi, 2019), we derive a taxonomy of parthood in intergroup relations, that defines what it means for a reference to constitute a reference towards the in-group/out-group/other:

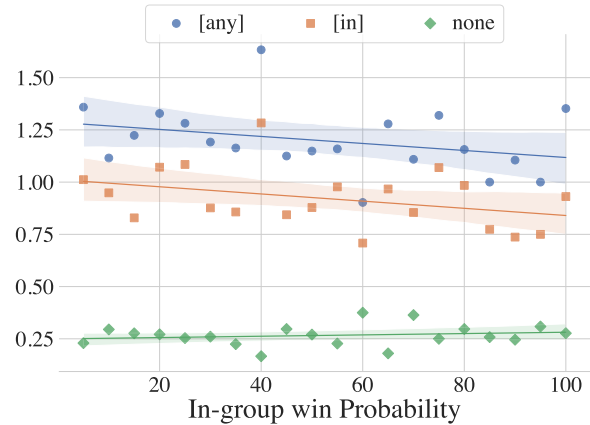


Figure 2: Per-comment frequency of in-group, any and ‘none’ references in gold dataset over WP.

1. **People:** Names, nicknames, shirt numbers, initials, pronouns, etc. : *Tua, TK87, he/him. . .*
2. **Subset of the team:** This refers to groups of players, or coaches, rather than just one player: *the offense, our defense, o-line, . . .*
3. **Team:** Name of the team (*rams, bills, cowboys*), nicknames (*lambs, cowgirls*), city names (*LA, Buffalo, Dallas*), pronominal expressions like *our boys* for the in-group, pronouns like *they/them* for the in-group and out-group, and many more.
4. **Team plus supporters:** The first person pronouns *we* and *us*, but can also be done with the third person pronouns *they* and *them*. The latter of course, could also refer to out-group or other, and require context to disambiguate.

The taxonomy above is ordered in order of increasing coverage of the whole group, by the referring part — the size of the reference gets larger from people to the entire group. Thus, players are the smallest unit of reference within a group, and the team/organization plus its supporters constitute the largest possible reference to the group itself.

Trends The annotated dataset enables us to study qualitative trends, that will guide quantitative modeling analyses presented in §6. We specifically focus on two phenomenon that are directly observable in the data and illustrated with examples — diversity in form of referring expression, and trends over WP. Within the gold dataset, we can observe two clear trends by plotting the frequency of a feature of interest over comments that fall within a win probability (WP) window. Figure 2 plots the frequency of any reference, in-group references, and ‘None’ references over all 5% WP windows:

1. References to the in-group, and references to

any group overall, go down with WP.
 2. ‘None’ references increase steadily with WP.
 The steady increase in number of ‘None’ references in higher WP windows is interesting, but requires robust analysis. While the trends observed in this section are not statistically significant, this can be attributed to the small sample size of only 1499 comments. The intergroup bias is a social phenomenon, and like many social phenomenon, we can make clear inferences at scale. Obtaining human annotated data at scale would be prohibitively hard and expensive in this setting — we use LLMs, to automate this task, thus allowing us make inferences about trends in the intergroup bias as a function of WP.

5 Modeling intergroup bias with LLMs

Large Language Models (LLMs) have shown remarkable abilities in various domains over the last few years (Srivastava et al., 2023; Brown et al., 2020). Our novel tagging framework to model intergroup bias requires linguistic understanding, knowledge of the NFL and its teams, as well as complex reasoning over why a commenter might choose certain word forms compatible with the state-of-the-world — making LLMs well suited to this task. In this section, we design modeling experiments to tag comments from our dataset with intergroup labels towards two objectives:

- Understand how LLMs statistically ‘reason’ over meaning in context of an utterance and game state (WP) to tag comments.
- Discover hidden intergroup variations in referring expressions by tagging a large sample of comments from our raw, scraped data.

5.1 Modeling conditions

We focus on two specific models — Llama-3-8B (AI@Meta, 2024) and GPT-4o (OpenAI, 2024). Both are decoder based models that perform best at a wide variety of benchmarks, and allow us to compare and contrast the performance of an open-weights model with finetuning, versus a larger closed model with few-shot prompting. Building upon previous work, we prompt both models with a combination of instructions, chain-of-thought explanations, and few-shot examples (Wei et al., 2022). Llama-3 is prompted with the same input format, but we also finetune the model on the train split of our gold dataset. See Appendix D for further details on training and inference.

CoT Explanations We finetune Llama-3 with GPT-4o generated CoT explanations (Wadhwa et al., 2023). We first generate a explanation from GPT-4o for each comment in our gold dataset using instructions, few-shot examples, the target tagged comment and list of referring expressions provided as input to GPT-4o. All few-shot explanations were written by the first author, and examples were drawn from outside the gold dataset.

Our task is framed end-to-end as the model receiving the untagged comment as input with some contextual information (in-group, out-group, WP), and being asked to generate the comment with relevant words/phrases replaced with the appropriate tags. To understand the impact of WP on model performance we design 3 conditions

Numeric WP The model receives WP as a numeric input — a percentage between 0 and 100 that is WP for the in-group.

No WP WP is not provided as input to the model, and the instructions nor few-shot explanations neither use nor mention it.

Linguistic WP We experiment with providing WP as a scalar description of game state, from ‘Team A is very likely to win’ to ‘Team B is very likely to win’ based on the numeric WP corresponding to the comment.

We also experimented with utilizing the WP to modify the temperature when decoding (Atwell et al., 2022). When temperature scaling (TS) is used, we set the temperature to $\sin(\pi \cdot WP)$ — this pushes the LM to choose less likely words when the game’s outcome is more uncertain.

Evaluation To evaluate the performance of a model on the test dataset, we report **micro-F1** scores for each of the three tags, and a weighted macro-F1 score overall. To give partial credit for the model’s tagged output slightly overlapping with the gold tagged spans, we assign partial scores (0.5 and 0.25) for being within 3 and 5 characters of the correct tagged spans respectively.

5.2 Results

Table 1 shows the results for both of our models on all conditions. While both models exceed the human baseline performance that we calculated in §4, Llama-3 nudges GPT-4o overall. GPT-4o performs better at identifying out-group and other references by names or nicknames due to its much larger size and more parametric knowledge.

	Model	Random Baseline	Numeric	Numeric WP+TS	No WP	No WP +TS	Ling. WP +TS	Ling WP
GPT-4o	[IN]	35.6(3.2)	66.6(1.4)	67.4(2.0)	67.7(1.9)	69.9(1.8)	71.2(0.7)	71.7(0.8)
	[OUT]	20.1(1.1)	64.6(3.1)	67.1(2.0)	63.6(2.7)	66.6(1.0)	63.4(1.9)	63.7(2.5)
	[OTHER]	14.0(5.9)	54.1(1.1)	53.9(2.1)	49.0(1.6)	47.5(2.1)	48.4(4.7)	51.6(8.8)
	Overall	30.8(2.7)	64.9(1.3)	65.9(1.5)	65.0(1.4)	66.9(1.3)	67.4(1.3)	68.2(0.3)
Llama-3-8b	[IN]	35.6(3.2)	72.0(1.5)	72.0(1.5)	72.5(0.6)	72.1(0.4)	72.6(1.9)	72.6(1.9)
	[OUT]	20.1(1.1)	60.2(2.3)	57.9(4.3)	59.9(0.7)	58.3(4.3)	58.7(3.5)	57.7(4.7)
	[OTHER]	14.0(5.9)	59.2(8.8)	58.8(8.2)	64.0(4.2)	57.6(6.8)	59.8(4.5)	59.1(5.4)
	Overall	30.8(2.7)	68.8(2.4)	68.4(2.0)	69.6(0.9)	68.4(1.0)	68.9(2.5)	68.8(2.6)

Table 1: Results from few-shot experiments on GPT-4o (top), and finetuning Llama-3-8b (bottom).

WP helps... sometimes? Including WP did not change the performance of Llama-3 noticeably. As we observed in annotation, there are few examples of comments being ambiguous enough that the state-of-the-world is enough to disambiguate what a reference could be. Entire classes of references (from our taxonomy in §4.2) are quite unambiguous even without whole-sentence context.

We do observe however that providing WP in language form boosts the in-group tagging performance of GPT-4o in few-shot settings, besting numeric WPs ($p < 0.005$ with a bootstrap test). Analysis of model’s outputs reveal GPT-4o’s fickleness and inability to reason over numerical scales — for instance it reasons that WPs ranging from 1% to as high as 41% are ‘low’ in its explanations. Further, it rarely uses the numbers to infer the WP for the out-group in explanations. Since we re-write low WPs with the name of the out-group (as winning) in the linguistic WP condition, this might explain the model’s slight boost in performance.

While Llama-3’s performance is better through finetuning, it does not benefit from incorporating WP in training or inference. We attempted scaling the loss during training with WP and expert annotator confidence ratings where available, but these didn’t boost performance. Whether larger LLMs exhibit similar behaviors to GPT-4o when finetuned, we leave to future work.

Error analysis While phrasing the WP with words improves tagging performance on in-group referents (over numerical WP as Table 1 shows), especially with GPT-4o, performance on out-group references remains stable. However, we do observe the model making similar ‘errors’ like annotators that we described in §4.1 — for instance, GPT-4o occasionally tagged *a single WR* (an indefinite/generic entity) as out-group in (8-a), which

was not judged to be a relevant referent in expert annotation. However, GPT-4o also makes some basic errors, such as tagging *a catch* as out-group in (8-b) as well.

- (8) a. Wait wait wait. Did I hear that right? **They** don’t have **a single WR** with a catch today??
b. Wait wait wait .Did I hear that right ? [OUT] don’t have [OUT] with [OUT]??

6 Analysis of model-tagged comments

Our novel tagging framework is amenable to application on a large sample of our raw data, facilitating us to observe and analyze variations in how members refer to different groups as a function of WP. We sample 100,000 comments from a larger raw dataset, and apply our best performing finetuned Llama-3 model towards this task. Since WP could not be effectively incorporated to improve performance, we used the model finetuned with no WP information. Further, we verified there was no correlation between the model’s accuracy and WP — accuracy mostly followed comment density across WP (see Appendix A). After inference from our finetuned LLM, we use regular expressions to ensure that any obvious words (names and nicknames of teams, *we/us*) were also tagged appropriately, and to count different referring expressions accounting for inflections.

Figure 3 plots the frequency of different references over WP. We divide WP into 5% windows, and count the number of comments that contain a specific tag (or tag-lexical item pair), and divide it by the number of comments within that window in the entire sample. Figure 4 plots a few more reference variables of interest and is similar except the variables are normalized by number of comments that contain *any* reference. There are two findings we wish to highlight.

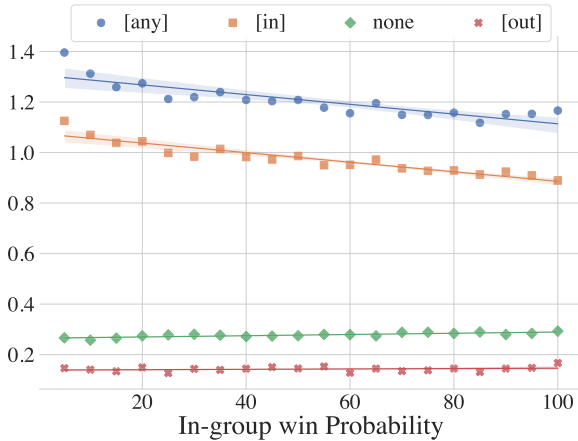


Figure 3: Per-comment frequency of various reference variables over all 5% WP windows. A 95% CI regression line is fit separately for each variable.

Winning trumps all Figure 3 clearly shows a linear, decreasing tendency for commenters to refer to any entity or the in-group the more likely the team is to win. Observing a sample of high WP (9)) comments reveals an increase in positive excitement, but also increased terseness closer to victory:

- (9) a. HOLY S**T
- b. WHAT A THROW

Figure 4 illustrates the tendency to refer less to the in-group, is compounded with an increased tendency to refer to the out-group, or to refer implicitly over a sentence (the [SENT] lemma) when referring at all over WP. Overall, this paints a clear picture — the more likely the in-group is to win, commenters prefer to refer to the out-group or to refer to entities implicitly if they refer at all; They prefer to abstract away from specific events and express excitement the closer the in-group is to winning.

WP as a well calibrated predictor A striking feature of Figures 3 and 4 is the linear relationships between reference variables and WP. Table 2 in Appendix E estimates coefficients and R-squared for linear fits, but we can observe visually that with increasing WP, commenters are more likely to refer to the out-group with ‘they’ than the in-group. From Figure 4, the slope for in-group references ($-2.8e^{-4}$) is larger than the slope for references to the in-group using first person singular pronouns ($-2e^{-4}$); Commenters are more likely to refer to in-group using the most inclusive term at higher WPs, when referring to the in-group at all.

These findings add to the subtle ways we perpetuate bias in our linguistic behavior, especially towards **in-group protection** (Maass, 1999). While

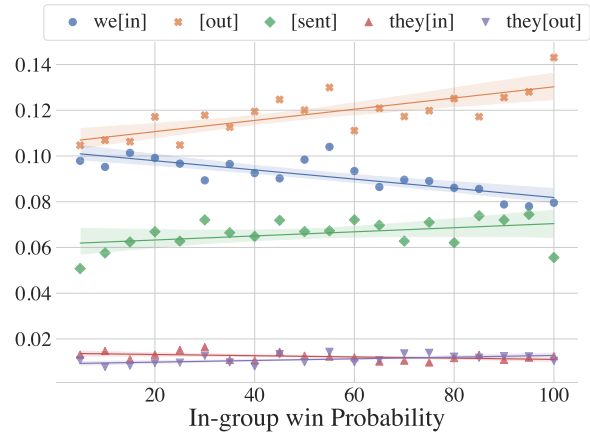


Figure 4: Normalized per-comment frequency of various reference variables over all 5% WP windows. A 95% CI regression line is fit separately for each variable.

commenters are more than willing to criticize the in-group across WP, the self-protective instinct is evident in the way they choose to refer to the in-group using *we/us* more often when losing, the reduced tendency to refer to the in-group using *they/them*, or to not refer to the in-group at all. Thus, how commenters choose the form of reference to an in-group constitutes just as subtle a bias as their choice of predicate.

7 Conclusion

We expand the study of intergroup bias by building a parallel corpus of sports comments grounded in win probabilities from live games. Annotation experiments reveal that modeling the bias as a tagging problem over words can reveal unobserved variations, as well as make it amenable to large-scale modeling. Through few-shot and finetuning experiments, we find that LLMs can out-match human baseline performance at this task, but struggle to reason over win probabilities, or use it meaningfully towards tagging. Tagging a large sample of our dataset reveals linear trends between various referring expressions and WP, showing that intergroup bias can manifest in commenter’s choice of who to refer to when commenting on a game and how. Careful data curation and understanding, combined with focused usage of LLMs as statistical information processing tools can thus reveal linguistic variations in social language use online at scale. In future work we plan to exploit the parallel nature of our corpus further to understand team differences in language variation, as well as how WP can be effectively incorporated into a model of social meaning.

688 Limitations

689 Our work expands the study of intergroup bias in
690 language by focusing on natural language use in
691 online conversations on the Reddit platform. Fur-
692 ther, our focus on grounding the utterances lead us
693 to focus on sports talk, specifically conversations
694 around NFL games. Biases in demographics of
695 users on Reddit, or demographics of NFL fans are
696 thus inherent in our data and analysis. Future work
697 needs to study the prevalence of our findings in
698 other sports with similar statistics that enables ef-
699 ficient grounding of utterances, as well as in more
700 general speech.

701 We identify that both few-shot performance by
702 GPT-4o and finetuned performance by Llama-3 are
703 close to, or out-perform the human ceiling perfor-
704 mance. Human ceiling performance is simply the
705 average accuracy of crowd annotators against ex-
706 pert annotators. As we note in the paper, this is
707 a difficult and inherently subjective task. Our re-
708 sults do not mean that models (finetuned or not)
709 have a better understanding of what constitutes in-
710 tergroup references, nor that they are more aligned
711 with the task. Llama-3 was trained on the training
712 split of the expert annotated gold data-set. While
713 GPT-4o was exposed to the same set of examples
714 as human annotators, it is a very large (possibly a
715 mixture of trillions of parameters) model that con-
716 tains a multitude of statistical associations that aids
717 in instruction following.

718 Ethics

719 We downloaded comments from Reddit threads
720 using the official Reddit API, and will disseminate
721 our data in accordance with the Reddit terms of
722 service. We will only release the comment and
723 post ids for the raw data, and usernames will be
724 anonymized. We will release the annotated data in
725 full with the same precautions. We have censored
726 some of the profanity in the comments when used
727 as examples in this paper, since our focus isn't on
728 abusive/negative language exclusively.

729 All created artifacts from this work (code, anno-
730 tated data) will be released under the MIT License.

731 Crowd-sourced annotations were collected from
732 three undergraduates employed by one of the au-
733 thors for 15\$ an hour.

734 References

735 AI@Meta. 2024. [Llama 3](#).

- Luigi Anolli, Valentino Zurloni, and Giuseppe Riva. 2006. Linguistic Intergroup Bias in Political Communication. *The Journal of General Psychology*, 133:237 – 255. 736
737
738
739
- Katherine Atwell, Remi Choi, Junyi Jessy Li, and Malihe Alikhani. 2022. The role of context and uncertainty in shallow discourse parsing. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 797–811, Gyeongju, Republic of Korea. International Committee on Computational Linguistics. 740
741
742
743
744
745
746
- Ben Baldwin. 2021. Open source football: nffastr ep, wp, cp xyac, and xpass models. 747
748
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. *ArXiv*, abs/2005.14165. 749
750
751
752
753
754
755
756
757
758
759
760
- Joseph L Fleiss. 1971. [Measuring nominal scale agreement among many raters](#). *Psychological bulletin*, 76(5):378. 761
762
763
- Bradley W. Gorham. 2006. News Media’s Relationship With Stereotyping: The Linguistic Intergroup Bias in Response to Crime News. *Journal of Communication*, 56(2):289–308. Place: United Kingdom Publisher: Blackwell Publishing. 764
765
766
767
768
- Venkata Subrahmanyam Govindarajan, Katherine Atwell, Barea Sinno, Malihe Alikhani, David I. Beaver, and Junyi Jessy Li. 2023a. [How people talk about each other: Modeling generalized intergroup bias and emotion](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2496–2506, Dubrovnik, Croatia. Association for Computational Linguistics. 769
770
771
772
773
774
775
776
777
- Venkata Subrahmanyam Govindarajan, David Beaver, Kyle Mahowald, and Junyi Jessy Li. 2023b. [Counterfactual probing for the influence of affect and specificity on intergroup bias](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12853–12862, Toronto, Canada. Association for Computational Linguistics. 778
779
780
781
782
783
784
- David Heineman, Yao Dou, and Wei Xu. 2023. Thresh: A unified, customizable and deployable platform for fine-grained text evaluation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 336–345, Singapore. Association for Computational Linguistics. 785
786
787
788
789
790
791

792	W. Hippel, Denise Sekaquaptewa, and P. Vargas. 1997.	Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A. Smith, and Yejin Choi. 2020.	847
793	The Linguistic Intergroup Bias As an Implicit Indicator of Prejudice. <i>Journal of Experimental Social Psychology</i> , 33:490–509.	Social bias frames: Reasoning about social and power implications of language. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5477–5490, Online. Association for Computational Linguistics.	848
794			849
795			850
796	Maksim Horowitz, Ron Yurko, and S Ventura. 2017.		851
797	nflscraper: Compiling the nfl play-by-play api for easy use in r .		852
798			853
799	Masahiro Kaneko and Danushka Bollegala. 2019.	Sherry B Schnake and Janet B Ruscher. 1998.	854
800	Gender-preserving Debiasing for Pre-trained Word Embeddings. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1641–1650, Florence, Italy. Association for Computational Linguistics.	Modern racism as a predictor of the linguistic intergroup bias. <i>Journal of Language and Social Psychology</i> , 17(4):484–491.	855
801			856
802			857
803		G. R. Semin and K. Fiedler. 1988.	858
804		The cognitive functions of linguistic categories in describing persons: Social cognition and language. <i>Journal of Personality and Social Psychology</i> , 54:558–568. Publisher: American Psychological Association.	859
805	Anne Maass. 1999. Linguistic Intergroup Bias: Stereotype Perpetuation Through Language. In Mark P. Zanna, editor, <i>Advances in Experimental Social Psychology</i> , volume 31, pages 79–121. Academic Press.		860
806			861
807			862
808		Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2020.	863
809	Anne Maass, Daniel Anthony Salvi, Luciano Arcuri, and Gün R. Semin. 1989. Language use in intergroup contexts: the linguistic intergroup bias. <i>Journal of personality and social psychology</i> , 57 6:981–93.	Towards Controllable Biases in Language Generation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 3239–3254, Online. Association for Computational Linguistics.	864
810			865
811			866
812			867
813	Christopher Manning, Mihai Surdeanu, John Bauer, Jenny Finkel, Steven Bethard, and David McClosky. 2014. The Stanford CoreNLP natural language processing toolkit . In <i>Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> , pages 55–60, Baltimore, Maryland. Association for Computational Linguistics.	Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. 2019.	868
814		The Woman Worked as a Babysitter: On Biases in Language Generation. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3407–3412, Hong Kong, China. Association for Computational Linguistics.	869
815			870
816			871
817			872
818			873
819			874
820			875
821	Jack Merullo, Luke Yeh, Abram Handler, Alvin Grisom II, Brendan O’Connor, and Mohit Iyyer. 2019.	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, and et al. 2023.	876
822	Investigating sports commentator bias within a large corpus of American football broadcasts. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 6355–6361, Hong Kong, China. Association for Computational Linguistics.	Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. <i>Transactions on Machine Learning Research</i> .	877
823			878
824			879
825			880
826			881
827		Achille Varzi. 2019. Mereology. In Edward N. Zalta, editor, <i>The Stanford Encyclopedia of Philosophy</i> , Spring 2019 edition. Metaphysics Research Lab, Stanford University.	882
828			883
829			884
830			885
831	OpenAI. 2024. Hello gpt-4o .	Somin Wadhwa, Silvio Amir, and Byron C. Wallace. 2023.	886
832	Konstantinos Pelechrinis and Evangelos Papalexakis. 2016.	Revisiting relation extraction in the era of large language models. <i>Proceedings of the conference. Association for Computational Linguistics. Meeting</i> , 2023:15566–15589.	887
833	The anatomy of american football: evidence from 7 years of nfl game data. <i>PLoS one</i> , 11(12):e0168716.		888
834			889
835			890
836	Reid Pryzant, Richard Diehl Martinez, Nathan Dass, Sadao Kurohashi, Dan Jurafsky, and Diyi Yang. 2020.	Albert Webson, Zhizhong Chen, Carsten Eickhoff, and Ellie Pavlick. 2020.	891
837	Automatically Neutralizing Subjective Bias in Text. <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , 34(01):480–489.	Are “Undocumented Workers” the Same as “Illegal Aliens”? Disentangling Denotation and Connotation in Vector Spaces. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4090–4105, Online. Association for Computational Linguistics.	892
838			893
839			894
840			895
841	Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza: A Python natural language processing toolkit for many human languages . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> .		896
842			897
843			898
844		Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022.	899
845		Chain-of-thought prompting elicits reasoning in large language models. In	900
846			901
			902

903 *Advances in Neural Information Processing Systems*,
904 volume 35, pages 24824–24837. Curran Associates,
905 Inc.

906 Ronald Yurko, Samuel L. Ventura, and Maksim
907 Horowitz. 2018. nflwar: a reproducible method for
908 offensive player evaluation in football. *Journal of*
909 *Quantitative Analysis in Sports*, 15:163 – 183.

910 Jason Shuo Zhang, Chenhao Tan, and Qin Lv. 2019.
911 [Intergroup contact in the wild: Characterizing lan-](#)
912 [guage differences between intergroup and single-](#)
913 [group members in nba-related discussion forums.](#)
914 *Proc. ACM Hum.-Comput. Interact.*, 3(CSCW).

915 A Data

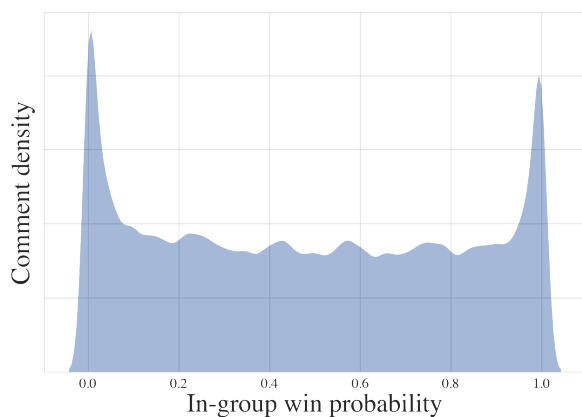


Figure 5: Comment density against WP.

916 B Annotation

917 **Protocol** Annotators were given the following
918 instructions:

- 919 1. All comments are from game threads corre-
920 sponding to specific NFL games between two
921 teams. You will be given the source of the
922 comment — this is the team the writer of the
923 comment supports, the opponent in that game,
924 and the live score at the time of making the
925 comment.
- 926 2. Highlight any words and phrases that refer to
927 individuals (people, teams, sub-groups within
928 the team, organizations).
- 929 3. If the reference is to the same group as the
930 source subreddit of the comment, tag this high-
931 light as **in-group** ([IN]).
- 932 4. If the reference is towards the opponent in this
933 specific game for which the comment is writ-
934 ten, tag this highlight as **out-group** ([OUT]).
- 935 5. If the reference is towards any other team
936 in the NFL apart from the two teams in-
937 volved in this game, tag this highlight as **other**
938 ([OTHER]).

- 939 6. Some comments will not have an obvious ref-
940 erence to an in-group/out-group/other entity.
941 Leave these comments un-annotated. If you're
942 unsure of an annotation, you can indicate your
943 confidence, but only use the confidence scale
944 if you're not very confident with your annota-
945 tion. I will take an empty confidence annota-
946 tion as full confidence.
- 947 7. Do not annotate a [SENT] token if there is
948 a word in the sentence that can be annotated
949 with the same label.

950 They were also given the following examples.
951 Models were finetuned with the following as few-
952 shot examples — they were provided WP over
953 live score for a more holistic representation of the
954 game, and explanations were modified depending
955 on whether WP was provided or not.

956 Example 1

957 COMMENT: [SENT] Defense getting
958 absolutely bullied by a dude that
959 looks like he sells solar panels
960 IN-GROUP: Jets
961 OUT-GROUP: Bears
962 LIVE SCORE: Jets 7 - 3 Bears
963 TARGET: [SENT] [IN] getting
964 absolutely bullied by [OUT] that
965 looks like [OUT] sells solar panels.
966 EXPLANATION: The commenter is
967 probably talking about the in-group,
968 since 'Defense' is said without
969 qualification, and the description of
970 the offensive player is disparaging
971 ('he sells solar panels'). 'Defense'
972 should be tagged [IN] since it refers
973 to in-group, and 'a dude' and 'he'
974 should be tagged [OUT] since it refers
975 to an out-group offensive player.

976 Example 2

977 COMMENT: [SENT] Hasn't really been him .
978 [SENT] Receivers have been missing a lot
979 of easy catches.
980 IN-GROUP: Dolphins
981 OUT-GROUP: Chargers
982 LIVE SCORE: Dolphins 0 - Chargers 0
983 TARGET: [SENT] Hasn't really been [IN] .
984 [SENT] [IN] have been missing a lot of
985 easy catches .
986 EXPLANATION: The second sentence is
987 complaining about the receivers missing
988 a lot of catches, thus absolving another

989	player of some blame, which is something	dead . [SENT] Suck it , [OUT] !	1038
990	fans would only do for the in-group team	EXPLANATION: This is a game between the	1039
991	they support. Thus 'him' in first sentence,	Chiefs and the Chargers, and the commenter	1040
992	and 'Receivers' in second sentence should	is a supporter of the Chiefs, so 'the	1041
993	be tagged with [IN].	chiefs' in the first sentence and 'KC'	1042
994	Example 3	in the second sentence should be tagged	1043
995	COMMENT: [SENT] Cards and rams are gonna	[OUT]. Herbert is a player for the	1044
996	be in the post-season regardless, so I	Chargers, and should be tagged with [IN]	1045
997	don't really care about them losing	since he is a member of the in-group with	1046
998	unless they play us.	respect to the commenter.	1047
999	IN-GROUP: 49ers	Example 6	1048
1000	OUT-GROUP: Jaguars	COMMENT: [SENT] Need points but 7 would	1049
1001	LIVE SCORE: 49ers 30 - 10 Jaguars	be HUGE momentum	1050
1002	TARGET: [SENT] [OTHER] and [OTHER] are	IN-GROUP: Bengals	1051
1003	gonna be in the post-season regardless,	OUT-GROUP: Chiefs	1052
1004	so I don't really care about [OTHER]	LIVE SCORE: Bengals 3 - 13 Chiefs	1053
1005	losing unless they play [IN].	TARGET: [IN] Need points but 7 would be	1054
1006	EXPLANATION: The game is between the	HUGE momentum	1055
1007	49ers and Jaguars, while the words 'Cards'	EXPLANATION: The in-group team is losing	1056
1008	and 'rams' refers to other teams in the	currently as the score shows, so this	1057
1009	NFL. Thus they should be tagged [OTHER]	comment is implicitly about the in-group	1058
1010	since they are neither in-group nor	needing points to gain momentum. Thus	1059
1011	out-group, as should the word 'them'.	'[SENT]' should be tagged with '[IN]'	1060
1012	'us' should be tagged [IN] since it	since there is no explicit word/phrase	1061
1013	refers to the in-group team the player	that refers to the in-group, but the	1062
1014	supports.	comment is referring to the in-group	1063
1015	Example 4	implicitly.	1064
1016	COMMENT: [SENT] How are we this shit on	C Prompts	1065
1017	defense	Below is the prompt provided to both GPT-4o and	1066
1018	IN-GROUP: Steelers	Llama-3. Examples are the same as the ones pro-	1067
1019	OUT-GROUP: Eagles	vided to human annotators, listed in the previous	1068
1020	LIVE SCORE: Steelers 7 - 21 Eagles	section. The following prompt does not use win	1069
1021	TARGET: [SENT] How are [IN] this shit on	probabilities; The prompts which do use WP are	1070
1022	defense	the same as below, except they include a definition	1071
1023	EXPLANATION: 'we' here, and almost always,	of WP as 'the probability of the in-group winning	1072
1024	refers to the in-group since they don't like	the game at the time of the comment - if the win	1073
1025	their team's defense, which is reflected in	probability is high, the in-group team is probably	1074
1026	the score. 'we' should therefore be tagged	doing well and going to win.' in the prompt text.	1075
1027	with [IN] since it refers to in-group.	Tag references to entities as in-group	1076
1028	Example 5	([IN]), out-group ([OUT]) or other	1077
1029	COMMENT: [SENT] The chiefs got	([OTHER]) in live, online sports comments	1078
1030	straight fucked with that Herbert INT	during NFL games. The input is the	1079
1031	getting called dead .	comment, the in-group team the commenter	1080
1032	[SENT] Suck it , KC !	supports and the out-group opponent team	1081
1033	IN-GROUP: Chargers	during that game. Using knowledge of	1082
1034	OUT-GROUP: Chiefs	American football and contextual language	1083
1035	LIVE SCORE: Chargers 28 - 28 Chiefs	understanding, identify words and phrases	1084
1036	TARGET: [SENT] [OUT] got straight	denoting entities (players, teams, city	1085
1037	fucked with that [IN] INT getting called	names, sub-groups within the team) that	1086

1087	refer to the in-group ([IN] - team the	references to tag.	1138
1088	commenter supports), out-group ([OUT] -	TARGET: [SENT] I thought so. [SENT]	1139
1089	the opponent) or other teams ([OTHER] -	Wish I could say the same ;)	1140
1090	some other team in the NFL that is not the		1141
1091	in-group or the opponent), with respect	Now tag only the relevant words/phrases	1142
1092	to the commenter. Return the list of	in the following comment as either	1143
1093	words/phrases that are to be tagged	in-group ([IN]), out-group ([OUT]), or	1144
1094	(REF_EXPRESSIONS), an EXPLANATION	other ([OTHER]), if any. First return	1145
1095	reasoning over why these words and phrases	the list of words to be tagged, then	1146
1096	in COMMENT should be tagged and with what	explain your reasoning as to why these	1147
1097	tag, and the TARGET comment itself with	words/phrases should be tagged from	1148
1098	relevant words/phrases replaced with the	COMMENT and with which tags, and finally	1149
1099	respective tags ([IN], [OUT] or [OTHER])	return the tagged comment in that order.	1150
1100	in your final output.		1151
1101		The explanations in the prompt with WP are sim-	1152
1102	Each sentence in a comment is separated by	ilar to the explanations provided previously. Here	1153
1103	a [SENT] token. Sometimes a sentence in	is the explanation for Example 1 from the previous	1154
1104	the comment will be about the in/out/other	section using win probabilities:	
1105	group but not have an explicit word/phrase	The commenter is probably talking about	1155
1106	that refers to the group; In such cases,	the in-group, since 'Defense' is said	1156
1107	tag the [SENT] token for that sentence with	without qualification, and the	1157
1108	the corresponding tag label.	description of the offensive player is	1158
1109		disparaging ('he sells solar panels').	1159
1110	Here are 6 examples, with	This is in spite of the win probability	1160
1111	REF_EXPRESSIONS being the list of	being relatively high for the in-group	1161
1112	words/phrases to be tagged from COMMENT,	- sometimes commenters choose to focus	1162
1113	EXPLANATION being a reasonable reason	on immediate plays rather than the	1163
1114	for why these words/phrases should be	overall state of the game, so perhaps	1164
1115	tagged with appropriate tags, and TARGET	this comment was in response to a bad	1165
1116	being the correct tagged output for	showing by the defense. 'Defense' should	1166
1117	COMMENT.	be tagged [IN] since it refers to	1167
1118		in-group, and 'a dude' and 'he' should	1168
1119	[EXAMPLES 1-6 FOLLOW HERE]	be tagged [OUT] since it refers to an	1169
1120		out-group offensive player.	1170
1121	Some comments will have no explicit or	Here is an example model input after the prompt:	1171
1122	implicit reference to the in-group,	COMMENT: [SENT] Defense getting	1172
1123	out-group, or other, or it could be	absolutely bullied by a dude that	1173
1124	extremely hard to disambiguate any	looks like he sells solar panels.	1174
1125	references based on given information.	IN-GROUP: Jets	1175
1126	In such cases, return Target as a copy	OUT-GROUP: Bears	1176
1127	of Comment, justify this with the	WIN PROBABILITY: 71.5%	1177
1128	Explanation, "No explicit or implicit	REF_EXPRESSIONS:	1178
1129	references to tag.", and return [] for		1179
1130	REF_EXPRESSIONS. Here is an example:	and here is the model's expected output:	
1131		['Defense', 'a dude', 'he']	1180
1132	COMMENT: [SENT] I thought so. [SENT]	EXPLANATION: The commenter is	1181
1133	Wish I could say the same ;)	probably talking about the in-group,	1182
1134	IN-GROUP: Jaguars	since 'Defense' is said without	1183
1135	OUT-GROUP: Titans	qualification, and the description	1184
1136	REF_EXPRESSIONS: []	of the offensive player is disparaging	1185
1137	EXPLANATION: No explicit or implicit	('he sells solar panels'). 'Defense'	1186
		should be tagged [IN] since it refers	1187

1188 to in-group, and 'a dude' and 'he'
 1189 should be tagged [OUT] since it refers
 1190 to an out-group offensive player.
 1191 TARGET: [SENT] [IN] getting absolutely
 1192 bullied by [OUT] that looks like [OUT]
 1193 sells solar panels .

1194 Explanations from GPT-4o were generated using
 1195 the same prompt format — the model is tasked with
 1196 generating EXPLANATION instead of TARGET.
 1197 The same examples were used, with EXPLANA-
 1198 TION at the end rather than TARGET. GPT-4o was
 1199 tasked to generate explanations of maximum length
 1200 256 tokens with seed of 1.

1201 **Linguistic WP** In the linguistic WP condition,
 1202 we replace the percentage WP value with a text
 1203 string like below:

- 1204 • 0–25: *Team A* is very likely to lose.
- 1205 • 25–45: *Team A* is likely to lose.
- 1206 • 45–55: Both teams are equally likely to win.
- 1207 • 55–75: *Team B* is likely to win.
- 1208 • 75–100: *Team B* is very likely to win.

1209 D Modeling implementation

1210 **GPT-4o** All few-shot experiments were run with
 1211 gpt-4o-2024-05-13. Temperature was set to 1 if
 1212 temperature scaling wasn't used, else it is dynami-
 1213 cally set to $\sin(\pi \times WP)$.

1214 **Llama-3-8B** We fine-tuned the base llama-3-8b
 1215 model from Meta's Huggingface model space². We
 1216 used the Axolotl³ framework for all fine-tuning
 1217 experiments with the following hyper-parameter
 1218 settings:

- 1219 • batch size of 4 for training and inference.
- 1220 • sample packing and padding to sequence
 1221 length were enabled, with a max sequence
 1222 length of 2560. None of our inputs exceeded
 1223 this limit.
- 1224 • Cosine learning rate scheduler with warmup
 1225 of 10 steps, learning rate set to $1e - 5$, weight
 1226 decay of 0.1, and a minimum learning rate
 1227 ratio of 0.1
- 1228 • Maximum of 2 train epochs with early stop-
 1229 ping, and patience set to 3.
- 1230 • The model is evaluated and saved every 59
 1231 steps for a maximum of 595 steps.
- 1232 • Flash attention and gradient checkpointing
 1233 were enabled.

²huggingface.co/meta-llama/Meta-Llama-3-8B

³<https://github.com/OpenAccess-AI-Collective/axolotl>

All finetuning experiments were done on 2
 Nvidia A40 GPUs, and each fine-tuning run took
 approximately 1.5 hours.

E Modeling Analysis

Feature	Slope($\times 10^{-4}$)	R-squared
Any reference	-19.3	0.72
No reference	2.4	0.65
In-group	-2.8	0.31
we	-2	0.61
out-group	2.5	0.56
they_in	-0.3	0.15
they_out	0.4	0.25

Table 2: Table of slopes of feature of interest against increasing WP, alongside the r-squared showing how much of the variance is explained by the linear regression fit. The slopes for Any and no reference are calculated with frequencies normalized by total number of referents in a WP window. All other slopes for referent variables are measured with frequencies normalized by comments with references in that WP window.