"We Demand Justice!": Towards Social Context Grounding of Political Texts

Anonymous ACL submission

Abstract

001Political discourse on social media often con-
tains similar language used with opposing in-
tended meanings. For example, the phrase
thoughts and prayers, is used to express sym-
pathy for mass shooting victims, as well as
satirically criticize the lack of legislative ac-
tion on gun control. Fully understanding such
discourse just by reading the text is difficult.
However, knowledge of the social context in-
formation makes it easier.

We characterize the social context required to fully understand such ambiguous discourse, by grounding the text in real-world entities, actions, and attitudes. We propose two datasets that require an understanding of social context and benchmark them using large pre-trained language models, and several novel structured models. We show that structured models, explicitly modeling social context, outperform larger models on both tasks, but still lag significantly behind human performance. Finally, we perform an extensive analysis, to obtain further insights into the language understanding challenges posed by our social grounding tasks.

1 Introduction

011

027

Over the past decade, micro-blogging websites have become the primary medium for US politicians to interact with general citizens and influence their stances for gaining support. As a result, politicians from the same party often coordinate the phrasing of their social messaging, to amplify their impact (Vaes et al., 2011; Weber and Neumann, 2021). Hence, repetitive, succinct phrases, such as "*Thoughts and Prayers*", are extensively used, although they signal more nuanced stances. Moreover, the interaction among politicians from opposing parties often leads to messaging phrased similarly, but signaling opposing real-world actions. For example, '*Thoughts and Prayers*', when used by Republicans, expresses condolences in mass



Figure 1: An example of varied *intended meanings* behind the same political message depending on the Author and Event in context

shooting events, but when used by Democrats conveys an *angry* or *sarcastic* tone as a call for action demanding "*tighter gun control measures*". Similarly, fig. 1 shows contrasting interpretations of the phrase "*We need to keep our teachers safe*!" depending on different speakers and in the context of different events.

041

042

043

044

047

050

051

053

055

056

059

060

061

062

063

064

Humans familiar with the stances of a politician and, possessing knowledge about the event from the news, can easily understand the intended meaning of political phrases. However, automatically understanding such language is challenging. Our main question in this paper is - Can an NLP model find the right meaning? From a linguistic perspective, we follow the distinction (Bach, 2008) between semantic interpretation (i.e., meaning encoded directly in the utterance and does not change based on its external context), and pragmatic interpretation (that depends on extra-linguistic information). The latter has gathered significant interest in the NLP community recently (Bender and Koller, 2020; Bisk et al., 2020), focusing on language understanding, when grounded in an external context (Fried et al., 2023). To a large extent,

Tweet Target Entity and Sentiment				Vague Text Disambiguation
Tweet: As if we needed more evidence. #kavanaugh				Vague Text: First, but not the last.
Event: Kavanaugh Supreme Court Nomination				Event: US withdraws from Paris climate agreement that enforces environmental targets after three years
Author: Earl Blumen	auer (De	emocrat Politician)		Author Party: Republican
Targets: Brett Kavanaugh (negative), Julie Swetnick (positive) Christine Ford (positive), Deborah Ramirez (positive)				Disambiguation: The withdrawal from the Paris climate agreement is the first step of many to come for the Trump administration. It will not be the last, as more positive changes are sure to follow. Incorrect Disambiguations:
Target Task Data Statistics Vague Text Data Statistics			tatistics	1) Joe Biden's inauguration marks the first day of a new era of progress
Unique Tweets	865	Unique Vague Texts	93	and prosperity, lasting positive changes are coming. (Incorrect Event)
Positive Targets	1513	Positive Examples	739	2) The Paris Climate Agreement withdrawal is the first of many
Negative Targets	1085	Negative Examples	2217	backward steps this Trump administration is sure to take in destroying
Neutral Targets	784	Total Examples	2956	our environment. (Incorrect Stance)
Non-Targets	2509	Number of Events	9	3) This is the time for America to move forward and make progress
Total Data Examples	5891	Hard Test Examples	180	without being held back by a global agreement that doesn't serve
Number of Events	3			our interests. (Doesn't match the vague text)

Table 1: Examples of Annotated Datasets and their statistics

the focus of such studies has been on grounding language in a perceptual environment (e.g., image captioning (Andreas and Klein, 2016; Sharma et al., 2018; Alikhani et al., 2020), instruction following (Wang et al., 2016; Suhr et al., 2019; Lachmy et al., 2022), and game playing (Potts, 2012; Udagawa and Aizawa, 2019) tasks). Unlike these works, in this paper, we focus on grounding language in a social context, i.e., modeling the common ground (Clark and Brennan, 1991; Traum, 1994; Stalnaker, 2002) between the author and their social media followers, that enables understanding an otherwise highly ambiguous utterances. The Social Context Understanding, needed for building successful models for such tasks, can come from a wide variety of sources. The politician's affiliation and historical stances on the issue provide can capture crucial social context. Social relationships, knowledge about the involved entities, and related prior and upcoming events form important part of the puzzle as well. In fig. 1 event #1, combining the event information (school shooting) with the speakers' gun control stances, would facilitate understanding the *intended meaning* of the text.

066

067

073

093

097

098

101

The main motivation of this paper work is to operationalize the 'Social Context Grounding' problem as a pragmatic understanding task. From a practical perspective, this would enable the creation of better NLP-CSS models that can process social media text in settings that require contextualized understanding. We suggest several datasets, designed to evaluate this ability in computational models. These task capture the intended meaning at different level of granularity. At the most basic level, providing the social context can help identify the entities targeted, and the sentiment towards them. In fig. 1, the social context (event#1,

Harris) and the text "we need to keep our teachers 102 safe" \Rightarrow "negative attitude towards guns". A more 103 nuanced account of meaning, which we formulate 104 as a separate task, captures the specific means in 105 which the negative attitude is expressed (the Inter-106 pretation in fig. 1). We additionally present two 107 datasets corresponding to these tasks, namely, 'Tar-108 get Entity and Sentiment Detection' and 'Vague 109 *Text Disambiguation*'. In the first, the goal is to 110 predict: 1) whether a given entity is the *intended* 111 target of a politician's tweet and 2) the sentiment 112 towards the intended targets. We explicitly focus 113 on tweets that do not always mention the targets 114 in their text to incentivize modeling the pragmatic 115 communicative intent of the text. In the second 116 task, given an ambiguous political message such 117 as "We demand justice" and its social context (as-118 sociated event, & the author's party affiliation), the 119 task is to identify a *plausible* unambiguous expla-120 nation of the message. Note that the ground truth 121 for all these tasks is based on human pragmatic 122 interpretation, i.e., "guns" is a negative target of 123 "we need to keep our teachers safe", despite not 124 being mentioned in the text, since it was perceived 125 in this way by a team of human annotators reading 126 the tweet and knowing social context. We show 127 examples of each task in table 1. We describe the 128 datasets in detail in section 3. 129

We evaluate the performance of various models, as a way to test the need for social context and compare different approaches for modeling it. These include pre-trained LM-based classifiers, and LLM in-context learning (Brown et al., 2020a; Black et al., 2022), which use a textual representation of the social context. We also adopt an existing graphbased discourse contextualization framework (Pujari and Goldwasser, 2021; Feng et al., 2022), to

130

131

132

133

134

135

136

137

138

2

explicitly model the social context needed to solve
the proposed tasks. Our results demonstrate that
the discourse contextualization models outperform
other models on both tasks. We present an error
analysis to gain further insights. We describe the
models in section 4 and the results in section 5.

145

146

147

148

149

150

151

152

153

154

156

157

158

159

160

161

162

163

165

166

168

169

170

171

172

174

175

176

178

179

180

181 182

183

187

We also present a qualitative visualization of a political event, *Brett Kavanaugh Supreme Court Nomination* (section 6.4), from target entitysentiment perspective. It showcases a unique summary of the event discourse. We perform human evaluation on our '*Vague Text Disambiguation*' dataset, and observe that humans find this task much easier than the evaluated models. We also present observations of human vs. LLM errors in disambiguation. In summary, our contributions are:

- 1. Defining and operationalizing the 'Social Context Grounding' task in political discourse
- 2. Evaluating various state-of-the-art context representation models on the task. We adopt existing discourse contextualization framework for the proposed tasks, and evaluate GPT-3's in-context learning performance, as well.
- 3. Performing human studies to benchmark the dataset difficulty and GPT-3 generation performance, when compared to human workers.¹

2 Related Work

Pragmatic Language Grounding gained significant focus recently (Bender and Koller, 2020; Bisk et al., 2020) following the rise of Pretrained Language Models (Devlin et al., 2019; Liu et al., 2019; Brown et al., 2020a) as unified NLP models. Most grounding tasks address multi-modal or physical environment descriptions (Barnard et al., 2003; Vogel and Jurafsky, 2010; Chen and Mooney, 2011; Tellex et al., 2011; Mitchell et al., 2012; Anderson et al., 2018). We refer the reader to (Fried et al., 2023) for a thorough overview. In contrast, we focus on grounding language in a *social* context. Social Context Modeling Hovy and Yang (2021) show that modeling social context is necessary for human-level NLU. As political messages are often targeted at the voter base aware of the political context (Weber and Neumann, 2021; Vaes et al., 2011), they are vague by design. Several previous works model social context for entity linking (Yang et al., 2016), social media connections relationship for fake news detection (Baly et al., 2018; Mehta et al., 2022) and, political bias detection (Li and

Goldwasser, 2019; Baly et al., 2020). These works model partial aspects of social context, relevant to their tasks. Two recent frameworks aim to capture social context holistically (Pujari and Goldwasser, 2021; Feng et al., 2022). Evaluation tasks presented in both works show interesting social context understanding but are not fully representative of the challenges of *Social Context Grounding*. Zhan et al. (2023) propose a dataset for dialogue understanding addressing general social commonsense.

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

Related Semantic and Pragmatic tasks closest to our Target Entity Sentiment Identification task is Stance Detection in social media (Mohammad et al., 2016; AlDayel and Magdy, 2020). To clarify our contribution, Mohammad et al. (2016), a popular SemEval task, looks at sentiment towards 5 targets, while our data has 362 unique targets. Allaway and McKeown (2020) and Zhang et al. (2022) also propose stance datasets on tweets. But, they focus mainly on semantic understanding of text that allows them to predict agreement or disagreement with well-defined statements. Our Vague Text Disambiguation task is related to recent works that study implicit inferences (Hoyle et al., 2023), and pragmatic understanding (Hu et al., 2023). However, our tasks evaluate pragmatic understanding using an explicit context, absent in those tasks.

3 Social Context Grounding Tasks

We design and collect two datasets for *Social Context Grounding* evaluation, and define three pragmatic interpretation tasks. In the *Tweet Target Entity and Sentiment* dataset, we collect annotations of opinionated tweets from known politicians for their intended targets and sentiments towards them. We focus on three political events for this task. The dataset and its collection are described below in section 3.1. In the *Vague Text Disambiguation Task*, we collect plausible explanations of vague texts, given the social context, consisting of *author affiliation* and *specific event*. We focus on eight political events. This dataset is detailed in section 3.2. Examples and data statistics are shown in table 1.

3.1 Tweet Target Entity and Sentiment Task

In this task, given a tweet T, its context, and an entity E, the objective is to predict whether or not Eis a target of T and the sentiment towards E. Political discourse often contains opinionated discourse about world events and social issues. We collect tweets that don't directly mention the target entities.

¹Our data and code will be released under MIT license

Thus, connecting the text with the event details and the author's general perspectives is necessary to solve this task effectively. We pick the focal entities for the given event and let human annotators expand on that initial set, based on their interpretation of the contextualized text. A *target* entity is conceptualized as an entity present in the full intended interpretation of the tweet.

237

238

241

243

245

246

247

251

253

254

We focus our tweet collection on three recent divisive events: George Floyd Protests, 2021 US Capitol Attacks, and Brett Kavanaugh's Supreme Court Nomination. We identify relevant participating entities for each of the three events. Examples of the involved entities for the event George Floyd Protests were George Floyd, United States Police, Derek Chauvin, Donald Trump, Joe Biden, United States Congress, Black people, Democratic Party, Republican Party, BLM, Antifa.

3.1.1 Target-Sentiment Data Collection

We filter 3, 454 tweets for the *three* events using 256 257 hashtags, keyword-based querying, and the dates of the event-based filtering from the Congress Tweets repository corpus². We collect a subset of 1,779tweets that contain media (images/video) to increase the chances of the tweet text not containing 261 the target entity mentions. Then, we use 6 in-house human annotators and Amazon Mechanical Turk (AMT) workers who are familiar with the event context for annotation. We ask them to annotate 265 the targeted entities and sentiments towards the tar-266 gets. The authors of this paper also participated 267 in the annotation process. We provide them with entity options based on the event in the focus of 269 the tweet. Annotators are allowed to add additional options if needed. We also ask the annotators to mark non-targets for each tweet. We instruct them 272 to keep the non-targets as relevant to the event as 273 possible to create harder negative examples. Each 274 tweet is annotated by three annotators. We filter 275 865 unique tweets with 5,891 annotations, with majority agreement on each tweet. All the AMT 277 annotations were additionally verified by in-house 278 annotators for correctness. AMT workers were paid USD 1 per tweet. It took 3 minutes on average for each assignment, resulting in an hourly 281 pay of USD 20. We include screenshots of the collection task GUIs in the appendix. We split the train, and test sets by events, authors, and targets to incentivize testing the general social grounding

capabilities of the models. The test set also consists of authors, targets, and events not seen in the training set. We use *Capitol Riots* event for the test set of *Target Entity and Sentiment* Task. We split the examples into 4, 370 train, 511 development, and 1,009 test examples. We compute the mean Cohen's kappa score for annotations and report inter-annotator agreement for annotated targets (0.47) and sentiment (0.73)

287

288

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

329

330

331

332

333

334

335

3.2 Vague Text Disambiguation Task

The task of *Vague Text Disambiguation* is designed to capture pragmatic interpretation at a finergrained level. It can be viewed as a variant of the well known paraphrase task, adapted for the social context settings. The model is evaluated on its ability to identify plausible interpretations (i.e., a sentence explicitly describing the author's intent) of an ambiguous quote given the event context and author's affiliation. E.g., "*protect our children from mass shootings*" could easily be disambiguated as either "*ban guns*" or "*arm teachers*" when the author's stance on the issue of '*gun rights*' is known.

Our data collection effort is designed to capture different aspects of social context grounding and facilitate detailed error analysis. Defined as a binary classification task over tuples (Party, Event, Vague text, Explicit text), we create negative examples by flipping tuple elements values of positive examples. This allows us to evaluate whether models can capture event relevance, political stance, or constrain the interpretation based on the vague text. For example, in the context of Event #1 in fig. 1, we can test if models simply capture the correlation between Democrats and negative stance towards guns access by replacing the vague text to "let your voice be heard", which would make the interpretation in fig. 1 implausible despite being consistent with that stance, while other consistent interpretations would be plausible (e.g., "go outside and join the march for our lives").

3.2.1 Vague Text Data Collection

Data collection was done in several steps. (1)**Vague Texts Collection**. We collected vague text candidates from tweets by US politicians (i.e. senators and representatives) between the years 2019 to 2021 from Congress Tweets corpus. We identified a list of 9 well-known events from that period and identified event-related tweets using their time frame and relevant hashtags. We used a pre-trained BERT-based (Devlin et al., 2019) NER model to

²https://github.com/alexlitel/congresstweets

collect tweets that do not contain any entity men-336 tions to identify potential candidates for vague texts. 337 We manually identified examples that could have contrasting senses by flipping their social context. We obtain 93 vague text candidates via this process. (2)In-Context Plausible Meaning Annotation. 341 We match the 93 ambiguous tweets with different events that fit them. We use both Democrat and Republican as the author party affiliation. We obtain 600 context-tweet pairs for AMT annotation. For each tweet, we ask AMT workers to annotate the following two aspects: 1) sentiment towards 347 the three most relevant entities in the event (sanity check) and 2) a detailed explanation of the intended meaning given the event and author's party affiliation. We obtain 469 reasonable annotations. After this step, each annotation was screened by in-house annotators. We ask three in-house annotators to vote on the correctness, appropriateness, and plau-354 sibility of the annotation given the context. Thus, we create a total of 374 examples.

(3)LLM-based Data Expansion. Using these examples, we further generate candidates for the task using LLM few-shot prompting. We use the examples from the previous step as in-context few-shot examples in the prompt. We use GPT-NeoX (Black et al., 2022) and GPT-3 (Brown et al., 2020a) for candidate generation. For each generated answer, manual inspection by three in-house annotators is performed to ensure data quality. We generate 928 candidates using GPT-NeoX and GPT-3. Human expert filtering results in 650 generations that pass the quality check. After removing redundant samples, we obtain 365 examples. Thus, we obtain a total of 739 annotations for this task. Then, for each of the 739 examples, we ask in-house anno-371 tators to select 3 relevant negative options from 372 the pool of explanations. We instruct them to pick hard examples that might contain similar entities 374 as the correct interpretation. This results in 2,956 375 binary classification data samples. We analyze and 376 discuss the results of human validation of large LM generations in section 6.

> Similar to the previous task, we split the train, test sets by events, and vague text to test the general social understanding capabilities of the model. We reserve *Donald Trump's second impeachment verdict* event for the test set. We also reserve Democratic examples of 2 events and Republican examples of 2 events exclusively for the test set. We split the dataset into 1, 916 train, 460 development, and 580 test examples. 180 of the test examples

384

387

are from events/party contexts unseen in train data.

391

392

393

394

395

396

397

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

4 Modeling Social Context

The key technical question this paper puts forward is how to model the social context, such that the above tasks can be solved with high accuracy. We observe that humans can perform this task well (section 6.3), and evaluate different context modeling approaches in terms of their ability to replicate human judgments. These correspond to No Context, Text-based context representation (e.g., Twitter Bio, relevant Wikipedia articles), and Graph-based context representation, simulating the social media information that human users are exposed to when reading the vague texts.

We report the results of all our baseline experiments in table 2 and table 3. The first set of results evaluate fine-tuned pre-trained language models (PLM), namely BERT (Devlin et al., 2019) and RoBERTa (Liu et al., 2019), with three stages of modeling context. Firstly, we evaluate no contextual information setting. Second, we include the authors' Twitter bios as context. Finally, we evaluate the information from the author, event, and target entity Wikipedia pages as context (models denoted PLM Baselines {No, Twitter Bio, Wikipedia} Context, respectively).

We evaluate GPT-3³ in *zero-shot* and *four-shot* in-context learning paradigm on both tasks. We provide contextual information in the prompt as short event descriptions and authors' affiliation descriptions. Note that GPT-3 is trained on news data until Sep. 2021 which includes the events in our data (models denoted LLM Baseline).

We evaluate the performance of politician embeddings from Political Actor Representation (PAR) (Feng et al., 2022) and Discourse Contextualization Framework (DCF) (Pujari and Goldwasser, 2021) models. (models denoted Static Contextutalized Embeddings). We use PAR embeddings available on their GitHub repository⁴. For DCF model, we use released pre-trained models from GitHub repository⁵ to generate author, event, text, and target entity embeddings. We evaluate the embeddings on both tasks. We briefly review these models in section 4.1 & section 4.2.

Finally, we use tweets of politicians from related previous events and build context graphs for each

³gpt-3.5-turbo-1106 via OpenAI API

⁴https://github.com/BunsenFeng/PAR

⁵https://github.com/pujari-rajkumar/ compositional_learner

Model -		Target Identification				Sentiment Identification			
		Prec	Rec	Macro-F1	Acc	Prec	Rec	Macro-F1	Acc
No Context	BERT-large	69.09	72.35	68.83	70.56	58.74	60.17	58.95	58.37
Baselines	RoBERTa-base	66.58	69.54	65.14	66.40	61.68	61.27	61.36	60.65
PLMs	BERT-large + user-bio	69.03	71.86	69.34	71.66	60.02	60.44	60.13	59.86
+Twitter Bio Context	RoBERTa-base + user-bio	65.83	68.65	64.79	66.30	60.06	59.91	59.94	59.46
PLMs	BERT-large + wiki	63.58	65.78	60.33	61.05	53.48	56.44	53.9	53.32
+Wikipedia Context	RoBERTa-base + wiki	69.02	72.32	68.62	70.27	57.62	59.10	58.07	58.28
LIMa	GPT-3 0-shot	69.25	70.58	69.77	73.78	56.20	55.04	54.18	56.80
LLIMS	GPT-3 4-shot	69.81	72.99	66.45	67.03	58.12	57.10	55.00	57.51
	RoBERTa-base + PAR Embs	68.38	71.63	67.67	69.18	55.01	56.89	55.51	55.40
Static Contextutalized	BERT-large + PAR Embs	65.40	67.33	60.25	60.56	55.24	57.54	55.89	55.80
Embedding Models	RoBERTa-base + DCF Embs	72.89	75.95	73.56	75.82	63.05	63.52	62.90	63.03
	BERT-large + DCF Embs	68.76	72.02	68.32	69.97	61.59	63.25	61.22	60.75
Discourse	BERT-large + DCF	71.12	74.61	71.17	72.94	65.81	65.25	65.34	65.31
Contextualized Models RoBERTa-base + DCF		70.44	73.86	70.39	72.15	63.45	63.34	63.37	63.23

Table 2: Results of baseline experiments on *Target Entity* (binary task) and *Sentiment* (4-classes) test sets. We report macro-averaged Precision, macro-averaged Recall, macro-averaged F1, and Accuracy metrics.

Madal	Vag	Vague Text Disambiguation					
widdei	Prec	Rec	Macro-F1	Acc			
No Context Baselines							
BERT-large	52.24	55.58	50.28	53.75			
RoBERTa-base	55.3	51.82	54.53	56.08			
PLMs + Wikipedia Conte	ext						
BERT-large + wiki	52.31	46.90	66.87	76.03			
BERT-base + wiki	51.85	38.62	64.36	75.69			
LLMs							
GPT-3 0-shot	63.10	62.92	62.58	63.5			
GPT-3 4-shot	62.05	62.29	61.86	62.04			
Static Contextutalized Er	nbeddin	g Mode	ls				
BERT-large + PAR	47.68	49.66	65.53	73.79			
BERT-base + PAR	45.93	54.48	65.49	72.59			
BERT-large + DCF Embs	47.18	63.45	67.55	73.10			
BERT-base + DCF Embs	56.58	59.31	71.71	78.45			
Discourse Contextualization Models							
BERT-large + DCF	52.76	59.31	69.94	76.55			
BERT-base + DCF	52.73	60.00	70.06	76.55			

Table 3: Results of baseline experiments on *Vague Text Disambiguation* dataset test split, a binary classification task. We report macro-averaged Precision, macro-averaged Recall, macro-averaged F1, and Acc. metrics

data example as proposed in Pujari and Goldwasser (2021). We use Wikipedia pages of authors, events, and target entities to add social context information to the graph. Then, we train the Discourse Contextualization Framework (DCF) for each task and evaluate its performance on both tasks (models denoted Discourse Contextualization Model). Further details of our baseline experiments are presented in subsection section 4.3. Results of our baseline experiments are discussed in section 5.

435

436

437

438

439

440

441

442

443

444

445

4.1 Discourse Contextualization Framework

446Discourse Contextualization Framework (DCF)447(Pujari and Goldwasser, 2021) leverages relations448among social context components to learn contex-449tualized representations for text, politicians, events,450and issues. It consists of *encoder* and *composer*451modules that compute holistic representations of

the context graph. The encoder creates an initial representation of nodes. Composer propagates the information within the graph to update node representations. They define link prediction learning tasks over context graphs to train the model. They show that their representations significantly outperform several PLM-based baselines trained using the same learning tasks. 452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

4.2 Political Actor Representation

Feng et al. (2022) propose the *Political Actor Representation* (PAR) framework, a graph-based approach to learn more effective politician embeddings. They propose three learning tasks, namely, 1) Expert Knowledge Alignment 2) Stance Consistency training & 3) Echo chamber simulation, to infuse social context into the politician representations. They show that PAR representations outperform SOTA models on *Roll Call Vote Prediction* and *Political Perspective Detection*.

4.3 Experimental Setup

Target Entity Detection is binary classification with $\langle author, event, tweet, target-entity \rangle$ as input and target/non-target label as output. Sentiment Detection is set up as 4-way classification. Input is the same as the target task and output is one of: {positive, neutral, negative, non-target}. Vague Text Disambiguation is a binary classification task with $\langle party-affiliation, event, vague-text, explanation-text \rangle$ and a match/no-match label as output.

In phase 1 no-context baselines, we use the author, event, tweet, and target embeddings generated by PLMs. We concatenate them for input. In Twitter-bio models, we use the author's Twitter bio embeddings to represent them. Wiki context models receive Wikipedia page embeddings of author,

event, and target embeddings. It is interesting to 487 note that the Wikipedia context models get all the 488 information needed to solve the tasks.. In phase 489 2 LLM experiments, we use train samples as in-490 context demonstrations. We provide task and event 491 descriptions in the prompt. In phase 3 PAR mod-492 els, we use politician embeddings released on the 493 PAR GitHub repository to represent authors. We re-494 place missing authors with their wiki embeddings. 495 For the Vague Text task, we average PAR embed-496 dings for all politicians of the party to obtain party 497 embeddings. For DCF embedding models, we gen-498 erate representations for all the inputs using context 499 graphs. We also use authors' tweets from relevant past events. We build graphs using author, event, tweet, relevant tweets, and target entity as nodes and edges as defined in the original DCF paper. In phase 4, we use the same setup as the DCF embedding model and additionally back-propagate to 505 DCF parameters. This allows us to fine-tune the DCF context graph representation for our tasks.

5 Results

508

509 The results of our baseline experiments are described in Tab. 2 and 3. We evaluate our models using macro-averaged precision, recall, F1, and ac-511 curacy metrics (due to class imbalance, we focus 512 on macro-F1). Several patterns, consistent across 513 all tasks, emerge. First, modeling social context is 514 still an open problem. None of our models were 515 able to perform close to human level. Second, 516 adding context can help performance, compared 517 to the No-Context baselines, models incorporating 518 519 context performed better, with very few exceptions. Third, LLMs are not the panacea for social-context 520 pragmatic tasks. Despite having access to a textual 521 context representation as part of the prompt, and having access to relevant event-related documents 523 during their training phase, these models under-524 perform compared to much simpler models that 525 were fine-tuned for this task. Finally, explicit con-526 text modeling using the DCF model consistently leads to the best performance. The DCF model mainly represents the social context in the form of text documents for all nodes. Further symbolic 530 addition of other types of context such as social 532 relationships among politicians and relationships between various nodes could further help in achieving better performance on these tasks. In the Target Entity task, RoBERTa-base + DCF embeddings obtain 73.56 F1 vs. 68.83 for the best no-context 536

baseline. Twitter bio and wiki-context hardly improve, demonstrating the effectiveness of modeling contextual information explicitly vs. concatenating context as text documents. No context performance well above the random performance of 50 F1 indicates the bias in the target entity distribution among classes. We discuss this in section 6.4. In Sentiment Identification task, we see that BERT-large + DCF back-propagation outperforms all other models. Vague Text Disambiguation task results in table 3 show that DCF models outperform other models significantly. 71.71 F1 is obtained by BERT-base + DCF embeddings. BERTbase performing better than bigger PLMs might be due to DCF model's learning tasks being trained using BERT-base embeddings.

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

6 Analysis and Discussion

6.1 Ablation Analysis on Vague Text Task

We report ablation studies in table 5 on the Vague Text task test set. We consider 5 splits: (1) Unseen Party: (party, event) not in the train set but (opposing-party, event) is present, (2) Unseen Event: $\langle party \rangle$ not in train set, (3) Flip Event: negative samples with corresponding 'event flippedparty/vague tweet matched' positive samples in train set and analogous (4) Flip Party and (5) Flip Tweet splits. We observe the best model in each category. They obtain weaker performance on unseen splits, as expected, unseen events being the hardest. Contextualized models achieve higher margins. DCF gains 7.6(13.2%) and DCF embeddings attain 8.12(20.42%) macro-F1 improvement over BERT-base+wiki compared to respective margins of 8.86% and 11.42% on the full test set. In the flip splits with only negative examples, accuracy gain over random baseline for all splits is seen. This indicates that models learn to jointly condition on context information rather than learn spurious correlations over particular aspects of the context. Specifically, flip-tweet split results indicate that models don't just learn party-explanation mapping.

6.2 Vague Text LLM Generation Quality

We look into the quality of our LLM-generated disambiguation texts. While GPT-NeoX (Black et al., 2022) produced only 98 good examples out of the 498 generated instances with the rest being redundant, GPT-3 (Brown et al., 2020a) performed much better. Among the 430 generated instances, 315 were annotated as good which converts to an accep-

Democrat Only Entities		Common Entities					Republican Only Entities	
Target Sentiment		Agreed-Upon Entities		Divisive Entities			Torgot	Sontimont
		Target	Sentiment	Sentiment (D)	Target	Sentiment (R)	Target	Sentiment
				Positive	Christine Blasey Ford	Negative	Sucon Colline	Positivo
Anita Hill	Positive	US Supreme Court	Neutral	Positive	Deborah Ramirez	Negative	Chuelt Creeseler	Positive
Patty Murray	Positive	US Senate	Neutral	Positive	Julie Swetnick	Negative	Diona Esinctain	Negotive
Merrick Garland	Positive	FBI	Neutral	Negative	Brett Kavanaugh	Positive	Chuelt Sehumen	Negative
Jeff Flake	Negative	Judiciary Committee	Neutral	Negative	Donald Trump	Positive	Chuck Schuller	Neutrol
				Negative	Mitch McConnell	Positive	Sean Hamily	Ineutral

Table 4: Target Entity-Sentiment centric view of Kavanaugh Supreme Court Nomination discourse

	Unseen	Unseen	Flip	Flip	Flip
Data Split	Party	Event	Tweet	Event	Party
	Ma-F1	Ma-F1	Acc	Acc	Acc
Random	44.70	29.69	75	75	75
BERT-base+wiki	57.58	39.76	88.14	89.77	87.77
BERT-base +DCF Embs	61.79	47.88	86.10	93.18	84.57
BERT-base+DCF	65.18	45.65	82.03	89.77	84.04

Table 5: Ablation Study Results on Vague Text Task

tance rate of 20.04% for GPT-NeoX and 73.26% for GPT-3 respectively. In-house annotators evaluated the quality of the generated responses for how well they aligned with the contextual information. They rejected examples that were either too vague, align with the wrong ideology, or were irrelevant. In the prompt, we condition the input examples in all the few shots to the same event and affiliation as the input vague text. In comparison, the validation of AMT annotations for the same task yielded 79.8% good examples even after extensive training and qualification tests. Most of the rejections from AMT were attributed to careless annotations.

588

589

591

592

593

594

595

599

600

604

607

610

611

612

613

614

615

6.3 Vague Text Human Performance

We look into how humans perform on the *Vague Text Disambiguation* task. We randomly sample 97 questions and ask annotators to answer them as multiple-choice questions. Each vague text-context pair was given 4 choices out of which only one was correct. We provide a brief event description along with all the metadata available to the annotator. Each question was answered by 3 annotators. Among the 97 answered questions, the accuracy was 94.85%, which shows this task is easy for humans who understand the context. Respective performance of best models on this subset of data for BERT-base+wiki (54.89%), BERT-base+DCFembs (63.38%), BERT-base+DCF (64.79%) is much lower than human performance.

6.4 Target Entity Visualization

The main goal of this analysis is to demonstrate the usefulness and inspire modeling research in the direction of entity-sentiment-centric view of political events. table 4 visualizes one component of how partisan discourse is structured in these events. We study *Kavanaugh Supreme Court Nomination*. We identify discussed entities and separate them into divisive and agreed-upon entities. This analysis paints an accurate picture of the discussed event. We observe that the main entities of Trump, Dr. Ford, Kavanaugh, Sen. McConnell, and other accusers/survivors emerge as divisive entities. Entities such as Susan Collins and Anita Hill who were vocal mouthpieces of the respective party stances but didn't directly participate in the event emerge as partisan entities. Supreme Court, FBI, and other entities occur but only as neutral entities. 619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

6.5 DCF Context Understanding

We look into examples that are incorrectly predicted using Wikipedia pages but correctly predicted by the DCF model in the appendix (table 6). In examples 1 & 2 of *Target Entity-Sentiment* task, when the entity is not explicitly mentioned in the tweet, the Wiki-Context model fails to identify them as the targets. We posit that while the Wikipedia page of each relevant event will contain these names, explicit modeling of entities in the DCF model allows correct classification. Examples 1 - 3 of *Vague Text Disambiguation* task show that when no clear terms indicate the sentiment towards a view, the Wiki-Context model fails to disambiguate the tweet text. Explicit modeling of politician nodes seems to help the DCF model.

7 Conclusion and Future Work

In this paper, we motivate, define, and operationalize 'Social Context Grounding' for political text. We build two novel datasets to evaluate social context grounding in NLP models that 'are easy for humans' when the relevant social context is provided. We experiment with many types of contextual models. We show that explicit modeling of social context outperforms other models while lacking behind humans.

Limitations

659

660 Our work only addresses English language text in US political domain. We also build upon large language models and large PLMs which are trained upon huge amounts of uncurated data. Although we employed human validation at each stage, biases could creep into the datasets. We also don't account for the completeness of our datasets as it is a pioneering work on a new problem. Social context is vast and could have a myriad of components. We only take a step in the direction of social context grounding in this work. The perfor-670 mance on these datasets might not indicate full so-671 cial context understanding but they should help in 672 sparking research in the direction of models that ex-673 plicitly model such context. Although we tuned our 674 prompts a lot, better prompts and evolving models 675 might produce better results on the LLM baselines. 676 Our qualitative analysis is predicated on a handful of examples. They are attempts to interpret the re-678 sults of large neural models and hence don't carry as much confidence as our empirical observations. We believe the insights from our findings will encourage more research in this area. For example, the development of discourse contextualized models that aim to model human-style understanding of background knowledge, emotional intelligence, and societal context understanding is a natural next step of our research.

Ethics Statement

704

705

In this work, our data collection process consists of using both AMT and GPT-3. For the *Target Entity and Sentiment* task, we pay AMT workers \$1 per HIT and expect an average work time of 3 minutes. This translates to an hourly rate of \$20 which is above the federal minimum wage. For the *Vague Text Disambiguation* task, we pay AMT workers \$1.10 per HIT and expect an average work time of 3 minutes. This translated to an hourly rate of \$22.

We recognize collecting political views from AMT and GPT-3 may come with bias or explicit results and employ expert gatekeepers to filter out unqualified workers and remove explicit results from the dataset. Domain experts used for annotation are chosen to ensure that they are fully familiar with the events in focus. Domain experts were provided with the context related to the events via their Wikipedia pages, background on the general issue in focus, fully contextualized quotes, and authors' historical discourse obtained from ontheissues.org. We have an annotation quid-pro-quo system in our lab which allows us to have a network of in-house annotators. In-house domain experts are researchers in the CSS area with familiarity with a range of issues and stances in the US political scene. They are given the information necessary to understand the events in focus in the form of Wikipedia articles, quotes from the politicians in focus obtained from ontheissues.org, and news articles related to the event. We make the annotation process as unambiguous as possible. In our annotation exercise, we ask the annotators to mark only high-confidence annotations that can be clearly explained. We use a majority vote from 3 annotators to validate the annotations for the target entity task.

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

Our task is aimed at understanding and grounding polarized text in its intended meaning. We take examples where the intended meaning is clearly backed by several existing real-world quotes. We do not manufacture the meaning to the vague statements, we only write down unambiguous explanations where context clearly dictates the provided meaning. Applications of our research as we envision would be adding necessary context to short texts by being able to identify past discourse from the authors that are relevant to the particular text in its context. It would also be able to ground the text in news articles that expand upon the short texts to provide full context.

References

- Abeer AlDayel and Walid Magdy. 2020. Stance detection on social media: State of the art and trends. *CoRR*, abs/2006.03644.
- Malihe Alikhani, Piyush Sharma, Shengjie Li, Radu Soricut, and Matthew Stone. 2020. Cross-modal coherence modeling for caption generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6525–6535, Online. Association for Computational Linguistics.
- Emily Allaway and Kathleen McKeown. 2020. Zero-Shot Stance Detection: A Dataset and Model using Generalized Topic Representations. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 8913– 8931, Online. Association for Computational Linguistics.
- Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton Van Den Hengel. 2018. Visionand-language navigation: Interpreting visuallygrounded navigation instructions in real environments. In *Proceedings of the IEEE conference on*

818

819

820

computer vision and pattern recognition, pages 3674–3683.

761

762

763

765

771

774

775

776

777

778

779

782

783

790

796

799

801

807

810

811

812

813

814

815

816

817

- Jacob Andreas and Dan Klein. 2016. Reasoning about pragmatics with neural listeners and speakers. In Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pages 1173– 1182, Austin, Texas. Association for Computational Linguistics.
- Kent Bach. 2008. *Pragmatics and the Philosophy of Language*, pages 463 487. Wiley Online Library.
- Ramy Baly, Giovanni Da San Martino, James Glass, and Preslav Nakov. 2020. We can detect your bias: Predicting the political ideology of news articles. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 4982–4991.
 - Ramy Baly, Georgi Karadzhov, Dimitar Alexandrov, James Glass, and Preslav Nakov. 2018. Predicting factuality of reporting and bias of news media sources. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3528–3539, Brussels, Belgium. Association for Computational Linguistics.
 - Kobus Barnard, Pinar Duygulu, David Forsyth, Nando De Freitas, David M Blei, and Michael I Jordan. 2003.
 Matching words and pictures. *The Journal of Machine Learning Research*, 3:1107–1135.
 - Emily M. Bender and Alexander Koller. 2020. Climbing towards NLU: On meaning, form, and understanding in the age of data. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5185–5198, Online. Association for Computational Linguistics.
 - Yonatan Bisk, Ari Holtzman, Jesse Thomason, Jacob Andreas, Yoshua Bengio, Joyce Chai, Mirella Lapata, Angeliki Lazaridou, Jonathan May, Aleksandr Nisnevich, Nicolas Pinto, and Joseph Turian. 2020.
 Experience grounds language. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 8718–8735, Online. Association for Computational Linguistics.
- Sid Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, USVSN Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. 2022. GPT-NeoX-20B: An opensource autoregressive language model. In Proceedings of the ACL Workshop on Challenges & Perspectives in Creating Large Language Models.
- 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, Jeffrey 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. 2020a. Language models are few-shot learners.

- 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, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam Mc-Candlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020b. Language models are few-shot learners. CoRR, abs/2005.14165.
- David Chen and Raymond Mooney. 2011. Learning to interpret natural language navigation instructions from observations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 25, pages 859–865.
- Herbert H Clark and Susan E Brennan. 1991. Grounding in communication. *Perspectives on Socially Shared Cognition*, pages 127 – 149.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Shangbin Feng, Zhaoxuan Tan, Zilong Chen, Ningnan Wang, Peisheng Yu, Qinghua Zheng, Xiaojun Chang, and Minnan Luo. 2022. Par: Political actor representation learning with social context and expert knowledge. *arXiv preprint arXiv:2210.08362*.
- Daniel Fried, Nicholas Tomlin, Jennifer Hu, Roma Patel, and Aida Nematzadeh. 2023. Pragmatics in language grounding: Phenomena, tasks, and modeling approaches. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12619– 12640, Singapore. Association for Computational Linguistics.
- Dirk Hovy and Diyi Yang. 2021. The importance of modeling social factors of language: Theory and practice. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 588–602, Online. Association for Computational Linguistics.
- Alexander Hoyle, Rupak Sarkar, Pranav Goel, and Philip Resnik. 2023. Natural language decompositions of implicit content enable better text representations. In *Proceedings of the 2023 Conference*

- 874 875
- 877

890

- 895
- 896 897
- 899
- 900 901
- 902
- 903
- 907
- 908 909

910 911 912

914 915

913

- 917
- 918 919
- 920
- 921 922

926

927 928

929 930

- on Empirical Methods in Natural Language Processing, pages 13188–13214, Singapore. Association for Computational Linguistics.
- Jennifer Hu, Sammy Floyd, Olessia Jouravlev, Evelina Fedorenko, and Edward Gibson. 2023. A finegrained comparison of pragmatic language understanding in humans and language models. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4194-4213, Toronto, Canada. Association for Computational Linguistics.
- Royi Lachmy, Valentina Pyatkin, Avshalom Manevich, and Reut Tsarfaty. 2022. Draw me a flower: Processing and grounding abstraction in natural language. Transactions of the Association for Computational Linguistics, 10:1341-1356.
 - Chang Li and Dan Goldwasser. 2019. Encoding social information with graph convolutional networks forPolitical perspective detection in news media. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 2594– 2604, Florence, Italy. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. ArXiv, abs/1907.11692.
- Nikhil Mehta, María Leonor Pacheco, and Dan Goldwasser. 2022. Tackling fake news detection by continually improving social context representations using graph neural networks. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1363-1380.
- Margaret Mitchell, Jesse Dodge, Amit Goyal, Kota Yamaguchi, Karl Stratos, Xufeng Han, Alyssa Mensch, Alexander Berg, Tamara Berg, and Hal Daumé III. 2012. Midge: Generating image descriptions from computer vision detections. In Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics, pages 747-756.
- Saif M. Mohammad, Parinaz Sobhani, and Svetlana Kiritchenko. 2016. Stance and sentiment in tweets. CoRR, abs/1605.01655.
- Christopher Potts. 2012. Goal-driven answers in the cards dialogue corpus. In Proceedings of the 30th west coast conference on formal linguistics, pages 1-20. Cascadilla Proceedings Project Somerville, MA.
- Rajkumar Pujari and Dan Goldwasser. 2021. Understanding politics via contextualized discourse processing. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 1353-1367, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. 2018. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2556–2565, Melbourne, Australia. Association for Computational Linguistics.

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

- Robert Stalnaker. 2002. Common ground. Linguistics and philosophy, 25(5/6):701-721.
- Alane Suhr, Claudia Yan, Jack Schluger, Stanley Yu, Hadi Khader, Marwa Mouallem, Iris Zhang, and Yoav Artzi. 2019. Executing instructions in situated collaborative interactions. 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 2119–2130, Hong Kong, China. Association for Computational Linguistics.
- Stefanie Tellex, Thomas Kollar, Steven Dickerson, Matthew Walter, Ashis Banerjee, Seth Teller, and Nicholas Roy. 2011. Understanding natural language commands for robotic navigation and mobile manipulation. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 25, pages 1507–1514.
- David Traum. 1994. A computational theory of grounding in natural language conversation.
- Takuma Udagawa and Akiko Aizawa. 2019. A natural language corpus of common grounding under continuous and partially-observable context. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 33, pages 7120–7127.
- Jeroen Vaes, Maria Paola Paladino, and Chiara Magagnotti. 2011. The human message in politics: The impact of emotional slogans on subtle conformity. The Journal of Social Psychology, 151(2):162–179. PMID: 21476460.
- Adam Vogel and Dan Jurafsky. 2010. Learning to follow navigational directions. In Proceedings of the 48th annual meeting of the association for computational linguistics, pages 806-814.
- Sida I. Wang, Percy Liang, and Christopher D. Manning. 2016. Learning language games through interaction. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 2368–2378, Berlin, Germany. Association for Computational Linguistics.
- Derek Weber and Frank Neumann. 2021. Amplifying influence through coordinated behaviour in social networks. Social Network Analysis and Mining, 11.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame,

987 Quentin Lhoest, and Alexander Rush. 2020. Trans988 formers: State-of-the-art natural language processing.
989 In Proceedings of the 2020 Conference on Empirical
990 Methods in Natural Language Processing: System
991 Demonstrations, pages 38–45, Online. Association
992 for Computational Linguistics.

993

994 995

996

997

998 999

1000

1001

1002

1003 1004

1005

- Yi Yang, Ming-Wei Chang, and Jacob Eisenstein. 2016. Toward socially-infused information extraction: Embedding authors, mentions, and entities. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1452–1461, Austin, Texas. Association for Computational Linguistics.
 - Haolan Zhan, Zhuang Li, Yufei Wang, Linhao Luo, Tao Feng, Xiaoxi Kang, Yuncheng Hua, Lizhen Qu, Lay-Ki Soon, Suraj Sharma, Ingrid Zukerman, Zhaleh Semnani-Azad, and Gholamreza Haffari. 2023. Socialdial: A benchmark for socially-aware dialogue systems.
- 1006Xinliang Frederick Zhang, Nick Beauchamp, and
Lu Wang. 2022. Generative entity-to-entity stance de-
tection with knowledge graph augmentation. In Pro-
ceedings of the 2022 Conference on Empirical Meth-
ods in Natural Language Processing, pages 9950–
9969, Abu Dhabi, United Arab Emirates. Association
for Computational Linguistics.

A Reproducibility

1013

We use the HuggingFace Transformers (Wolf et al., 1014 2020) library for PLMs. We use GPT-NeoX im-1015 plementation by ElutherAI (Black et al., 2022) and 1016 GPT-3 (Brown et al., 2020b) via OpenAI API for 1017 our LLM baselines. We run 100 epochs for all 1018 experiments. We use 10 NVIDIA GeForce 1080i 1019 GPUs for our experiments. We use the train, de-1020 velopment, and test splits detailed in section 3 for 1021 our experiments. We use the development macro-1022 F1 for early stopping. We run all our experiments 1023 using random seeds to ensure reproducibility. We 1024 experiment with a random seed value set to $\{13\}$. 1025 We report the results of the 3 fold cross-validation. 1026 We report only the mean on all the cross-fold validation results for clarity. All our code, datasets, 1028 and result logs will be released publicly upon ac-1029 ceptance. We experiment with 3, 5, and 10 fold 1030 cross-validation. As the results on the development 1031 data are almost identical, we report the results of 3 1032 fold cross-validation in all our experiments. 1033

1034	B Error Analysis
1035	C Annotation Interfaces
1036	D GPT Prompts
1037	GPT-3 prompts used for target-entity task:
1038	Event: <event></event>
1039	Event background:
1040	<background-description></background-description>
1041	Tweet: <tweet-text></tweet-text>
1042	Author: <author-name></author-name>
1043	Author Party: <party-affiliation></party-affiliation>
1044	Author background: <first sentences<="" td="" two=""></first>
1045	of author-wiki-page>
1046	Target Entity: <entity-name></entity-name>
1047	Entity background: <first sentence<="" th="" two=""></first>
1048	of entity-wiki-page> Task: Identify if
1049	the given entity is a target of the
1050	tweet. A target entity is defined as an
1051	entity that would be present in the full
1052	unambiguous explanation of the tweet.
1053	Is the given entity a target entity of
1054	the tweet? Answer yes or no.
1055	GPT-3 prompts used for target-sentiment task:
1056	Event: <event></event>
1057	Event background:
1058	<background-description></background-description>
1059	Tweet: <tweet-text></tweet-text>
1060	Author: <author-name></author-name>
1061	Author Party: <party-affiliation></party-affiliation>

Author background: <first sentences<="" th="" two=""><th>1062</th></first>	1062
of author-wiki-page>	1063
Target Entity: <entity-name></entity-name>	1064
Entity background: <first sentence<="" th="" two=""><th>1065</th></first>	1065
of entity-wiki-page> Task: Identify the	1066
sentiment of the tweet towards the given	1067
target entity. Consider that the tweet is	1068
ambiguous and the entity might be implied	1069
without being explicitly mentioned.	1070
What is the sentiment of the tweet	1071
towards the target entity? Answer with	1072
positive, negative, or neutral.	1073
GPT-3 prompts used for vague text disambigua-	1074
tion task:	1075
Event: <event></event>	1076
Event background:	1077
<background-description></background-description>	1078
Vague message: <vague-text></vague-text>	1079
Author Party: <party-affiliation></party-affiliation>	1080
Author background: <first sentences<="" th="" two=""><th>1081</th></first>	1081
of party-wiki-page>	1082
Task: Given the event, vague message, and	1083
party affiliation of the author, explain	1084
unambiguously the intended meaning of	1085
the vague message.	1086
Generate an unambiguous explanation	1087
for the vague message given the party	1088
affiliation of the author and the event	1089
in context.	1090

Target Entity and Sentiment Task	Vague Text Disambiguation Task
Tweet: Republicans held Justice Scalia's seat open for more	
than 400 days. Justice Kennedy's seat has been vacant for	
less than two months. It's more important to investigate a	Tweet: Thanks for this.
serious allegation of sexual assault than to rush Kavanaugh	
onto the Supreme Court for a lifetime appointment.	
Author: Adam Schiff (Democrat)	Affiliation: Democrat
Event: Brett Kavanaugh Supreme Court nomination	Event: United States withdrawal from the Paris Agreement
	Paraphrase: There's nothing surprising in withdrawing from
Entity: Christine Blasey Ford	the Paris agreement. Thanks for not caring our environment and
	future generations.
Wiki-Context Prediction: Not Target DCF Prediction: Target (correct)	Wiki-Context Prediction: No DCF Prediction: Yes (correct)
Tweet: We will not be intimidated. Democracy will not be	
intimidated. We must hold the individuals responsible for the	
Jan. 6th attack on the U.S. Capitol responsible. Thank you	Tweet: Let us say enough. Enough.
@RepAOC for tonight's Special Order Hour and we will	
continue our efforts to #HoldThemAllAccountable.	
Author: Adriano Espaillat (Democrat)	Affiliation: Democrat
Event: January 6 United States Capitol attack	Event : Second impeachment of Donald Trump ended with not guilty
	Paraphrase : The failure of the Democrats to impeach Donald Trump is a strong moment for our legislature which can get
Entity: Donald Trump	back to its work helping the American people. Today we've been able to tell the American people what we have known all along,
	that Donald Trump was not guilty of these charges.
Wiki-Context Predicted: Not Target DCF Prediction: Target (correct)	Wiki-Context Predicted: Yes DCF Prediction: No (correct)
Tweet: #GeorgeFloyd #BlackLivesMatter #justiceinpolicing	
QI @OmarJimenez Former Minneapolis police officer Derek	
Chauvin is in the process of being released from the Hennepin	Tweet: Lots of honking and screaming from balconies.
County correctional facility his attorney tells us. He is one of	Something must be going on.
the four officers charged in the death of George Floyd. He	
faces murder and manslaughter charges.	
Author: Adriano Espaillat (Democrat)	Affiliation: Democrat
Event: George Floyd protests	Event: Presidential election of 2020
Entity: Derek Chauvin	Paraphrase: I'm sure that the people are celebrating the election results.
Wiki-Context Predicted Sentiment: Positive DCF Prediction: Negative (correct)	Wiki-Context Prediction: No DCF Prediction: Yes (correct)

Table 6: Examples where baseline model fails but DCF works



Figure 2: An example of Tweet Target Entity and Sentiment Annotation GUI

Task	Example 1 Example 2					
Instru	Instructions					
Backg	round					
	General Context Short Event Description					
	Date Published (MM/DD/YYYY): 11/04/2020 After a three-year delay, the US has become the first nation in the world to formally with dearwise accompany					
	Author: Republican - Anonymous withdraw from the Paris climate agreement.					
	Event Happened: United States withdrawal from the Paris Agreement					
Tweet						
Ju	st the beginning					
▲ Pie Ev	▲ Please read the example, instruction and background information carefully before proceed. Even if you have seen the tweet before, the context has changed. Not reading it carefully may result in rejection.					
Sentime Entities m	ent Analysis ay not be present in the tweet text, please try your best to inference from the	context.				
Sentiment	t towards Donald Trump					
O Posi	O Positive O Neutral O Negative					
Sentiment	t towards Republicans					
O Posi	O Positive O Neutral O Negative					
Sentiment towards Democrats						
O Posi	O Positive O Neutral O Negative					
Paraphr	Paraphrase					
Provide an paraphrase for the given tweet and context						
Refer to the examples given and try to think about what the author want to say explicitly.						
Write the paraphrase here as if you are the original author						

Submit

Figure 3: An example of Vague Text Disambiguation GUI