

# Conversations Gone Awry: Detecting Early Signs of Conversational Failure

Justine Zhang and Jonathan P. Chang and Cristian Danescu-Niculescu-Mizil\*  
Cornell University

{jz727, jpc362}@cornell.edu, cristian@cs.cornell.edu

Lucas Dixon and Nithum Thain  
Jigsaw

{ldixon, nthain}@google.com

Yiqing Hua  
Cornell University

yh663@cornell.edu

Dario Taraborelli  
Wikimedia Foundation

dario@wikimedia.org

## Abstract

One of the main challenges online social systems face is the prevalence of antisocial behavior, such as harassment and personal attacks. In this work, we introduce the task of predicting from the very start of a conversation whether it will get out of hand. As opposed to detecting undesirable behavior after the fact, this task aims to enable early, actionable prediction at a time when the conversation might still be salvaged.

To this end, we develop a framework for capturing pragmatic devices—such as politeness strategies and rhetorical prompts—used to start a conversation, and analyze their relation to its future trajectory. Applying this framework in a controlled setting, we demonstrate the feasibility of detecting early warning signs of antisocial behavior in online discussions.

## 1 Introduction

“Or vedi l’anime di color cui vinse l’ira.”<sup>1</sup>

– Dante Alighieri, *Divina Commedia*, *Inferno*

Online conversations have a reputation for going awry (Hinds and Mortensen, 2005; Gheitsy et al., 2015): antisocial behavior (Shepherd et al., 2015) or simple misunderstandings (Churchill and Bly, 2000; Yamashita and Ishida, 2006) hamper the efforts of even the best intentioned collaborators. Prior computational work has focused on characterizing and detecting content exhibiting antisocial online behavior: trolling (Cheng et al., 2015, 2017), hate speech (Warner and Hirschberg, 2012; Davidson et al., 2017), harassment (Yin et al., 2009), personal attacks (Wulczyn et al.,

2017) or, more generally, toxicity (Chandrasekharan et al., 2017; Pavlopoulos et al., 2017b).

Our goal is crucially different: instead of identifying antisocial comments *after the fact*, we aim to detect *warning signs* indicating that a civil conversation is at risk of derailing into such undesirable behaviors. Such warning signs could provide potentially actionable knowledge at a time when the conversation is still salvageable.

As a motivating example, consider the pair of conversations in Figure 1. Both exchanges took place in the context of the Wikipedia discussion page for the article on the Dyatlov Pass Incident, and both show (ostensibly) civil disagreement between the participants. However, only one of these conversations will eventually turn awry and devolve into a personal attack (“Wow, you’re coming off as a total d\*\*k. [...] What the hell is wrong with you?”), while the other will remain civil.

As humans, we have some intuition about which conversation is more likely to derail.<sup>2</sup> We may note the repeated, direct questioning with which **A1** opens the exchange, and that **A2** replies with yet another question. In contrast, **B1**’s softer, hedged approach (“it seems”, “I don’t think”) appears to invite an exchange of ideas, and **B2** actually addresses the question instead of stonewalling. Could we endow artificial systems with such intuitions about the future trajectory of conversations?

In this work we aim to computationally capture linguistic cues that predict a conversation’s future health. Most existing conversation modeling approaches aim to detect characteristics of an observed discussion or predict the outcome after the discussion concludes—e.g., whether it involves a present dispute (Allen et al., 2014; Wang and Cardie, 2014) or contributes to the even-

\* Corresponding senior author.

<sup>1</sup>“Now you see the souls of those whom anger overcame.”

<sup>2</sup>In fact, humans achieve an accuracy of 72% on this balanced task, showing that it is feasible, but far from trivial.

---

**A1:** Why there's no mention of it here? Namely, an altercation with a foreign intelligence group? True, by the standards of sources some require it wouldn't even come close, not to mention having some really weak points, but it doesn't mean that it doesn't exist.

**A2:** So what you're saying is we should put a bad source in the article because it exists?

**B1:** Is the St. Petersburg Times considered a reliable source by wikipedia? It seems that the bulk of this article is coming from that one article, which speculates about missile launches and UFOs. I'm going to go through and try and find corroborating sources and maybe do a rewrite of the article. I don't think this article should rely on one so-so source.

**B2:** I would assume that it's as reliable as any other mainstream news source.

---

Figure 1: Two examples of initial exchanges from conversations concerning disagreements between editors working on the Wikipedia article about the Dyatlov Pass Incident. Only one of the conversations will eventually turn awry, with an interlocutor launching into a personal attack.

tual solution of a problem (Niculae and Danescu-Niculescu-Mizil, 2016). In contrast, for this new task we need to discover interactional signals of the *future* trajectory of an *ongoing* conversation.

We make a first approach to this problem by analyzing the role of politeness (or lack thereof) in keeping conversations on track. Prior work has shown that politeness can help shape the course of offline (Clark, 1979; Clark and Schunk, 1980), as well as online interactions (Burke and Kraut, 2008), through mechanisms such as softening the perceived force of a message (Fraser, 1980), acting as a buffer between conflicting interlocutor goals (Brown and Levinson, 1987), and enabling all parties to save face (Goffman, 1955). This suggests the potential of politeness to serve as an indicator of whether a conversation will sustain its initial civility or eventually derail, and motivates its consideration in the present work.

Recent studies have computationally operationalized prior formulations of politeness by extracting linguistic cues that reflect politeness strategies (Danescu-Niculescu-Mizil et al., 2013; Aubakirova and Bansal, 2016). Such research has additionally tied politeness to social factors such as individual status (Danescu-Niculescu-Mizil et al., 2012; Krishnan and Eisenstein, 2015), and the success of requests (Althoff et al., 2014) or of collaborative projects (Ortu et al., 2015). However, to the best of our knowledge, this is the first computational investigation of the relation between politeness strategies and the future trajectory of the conversations in which they are deployed. Furthermore, we generalize beyond predefined politeness strategies by using an unsupervised method to discover additional rhetorical prompts used to initiate different types of conversations that may be specific to online collaborative settings, such as coordinating work (Kittur and Kraut, 2008) or conducting factual checks.

We explore the role of such pragmatic and rhetorical devices in foretelling a particularly perplexing type of conversational failure: when participants engaged in previously civil discussion start to attack each other. This type of derailment “from within” is arguably more disruptive than other forms of antisocial behavior, such as vandalism or trolling, which the interlocutors have less control over or can choose to ignore.

We study this phenomenon in a new dataset of Wikipedia talk page discussions, which we compile through a combination of machine learning and crowdsourced filtering. The dataset consists of conversations which begin with ostensibly civil comments, and either remain healthy or derail into personal attacks. Starting from this data, we construct a setting that mitigates effects which may trivialize the task. In particular, some topical contexts (such as politics and religion) are naturally more susceptible to antisocial behavior (Kittur et al., 2009; Cheng et al., 2015). We employ techniques from causal inference (Rosenbaum, 2010) to establish a controlled framework that focuses our study on topic-agnostic linguistic cues.

In this controlled setting, we find that pragmatic cues extracted from the very first exchange in a conversation (i.e., the first comment-reply pair) can indeed provide some signal of whether the conversation will subsequently go awry. For example, conversations prompted by hedged remarks sustain their initial civility more so than those prompted by forceful questions, or by direct language addressing the other interlocutor.

In summary, our main contributions are:

- We articulate the new task of detecting early on whether a conversation will derail into personal attacks;
- We devise a controlled setting and build a labeled dataset to study this phenomenon;

- We investigate how politeness strategies and other rhetorical devices are tied to the future trajectory of a conversation.

More broadly, we show the feasibility of automatically detecting warning signs of future misbehavior in collaborative interactions. By providing a labeled dataset together with basic methodology and several baselines, we open the door to further work on understanding factors which may derail or sustain healthy online conversations. To facilitate such future explorations, we distribute the data and code as part of the Cornell Conversational Analysis Toolkit.<sup>3</sup>

## 2 Further Related Work

**Antisocial behavior.** Prior work has studied a wide range of disruptive interactions in various online platforms like Reddit and Wikipedia, examining behaviors like aggression (Kayany, 1998), harassment (Chatzakou et al., 2017; Vitak et al., 2017), and bullying (Akbulut et al., 2010; Kwak et al., 2015; Singh et al., 2017), as well as their impact on aspects of engagement like user retention (Collier and Bear, 2012; Wikimedia Support and Safety Team, 2015) or discussion quality (Arazy et al., 2013). Several studies have sought to develop machine learning techniques to detect signatures of online toxicity, such as personal insults (Yin et al., 2009), harassment (Sood et al., 2012) and abusive language (Nobata et al., 2016; Gambäck and Sikdar, 2017; Pavlopoulos et al., 2017a; Wulczyn et al., 2017). These works focus on detecting toxic behavior after it has already occurred; a notable exception is Cheng et al. (2017), which predicts future community enforcement against users in news-based discussions. Our work similarly aims to understand *future* antisocial behavior; however, our focus is on studying the trajectory of a conversation rather than the behavior of individuals across disparate discussions.

**Discourse analysis.** Our present study builds on a large body of prior work in computationally modeling discourse. Both unsupervised (Ritter et al., 2010) and supervised (Zhang et al., 2017a) approaches have been used to categorize behavioral patterns on the basis of the language that ensues in a conversation, in the particular realm of online discussions. Models of conversational behavior have also been used to predict conversation outcomes, such as betrayal in games (Niculae et al.,

2015), and success in team problem solving settings (Fu et al., 2017) or in persuading others (Tan et al., 2016; Zhang et al., 2016).

While we are inspired by the techniques employed in these approaches, our work is concerned with predicting the future trajectory of an ongoing conversation as opposed to a post-hoc outcome. In this sense, we build on prior work in modeling conversation trajectory, which has largely considered *structural* aspects of the conversation (Kumar et al., 2010; Backstrom et al., 2013). We complement these structural models by seeking to extract potential signals of future outcomes from the *linguistic discourse* within the conversation.

## 3 Finding Conversations Gone Awry

We develop our framework for understanding linguistic markers of conversational trajectories in the context of Wikipedia’s *talk page* discussions—public forums in which contributors convene to deliberate on editing matters such as evaluating the quality of an article and reviewing the compliance of contributions with community guidelines. The dynamic of conversational derailment is particularly intriguing and consequential in this setting by virtue of its collaborative, goal-oriented nature. In contrast to unstructured commenting forums, cases where one *collaborator* turns on another over the course of an initially civil exchange constitute perplexing pathologies. In turn, these toxic attacks are especially disruptive in Wikipedia since they undermine the social fabric of the community as well as the ability of editors to contribute (Henner and Sefidari, 2016).

To approach this domain we reconstruct a complete view of the conversational process in the edit history of English Wikipedia by translating sequences of revisions of each talk page into structured conversations. This yields roughly 50 million conversations across 16 million talk pages.

Roughly one percent of Wikipedia comments are estimated to exhibit antisocial behavior (Wulczyn et al., 2017). This illustrates a challenge for studying conversational failure: one has to sift through many conversations in order to find even a small set of examples. To avoid such a prohibitively exhaustive analysis, we first use a machine learning classifier to identify candidate conversations that are likely to contain a toxic contribution, and then use crowdsourcing to vet the resulting labels and construct our controlled dataset.

<sup>3</sup><http://convokit.infosci.cornell.edu>

**Job 1: Ends in personal attack.** We show three annotators a conversation and ask them to determine if its last comment is a personal attack toward someone else in the conversation.

**Job 2: Civil start.** We split conversations into snippets of three consecutive comments. We ask three annotators to determine whether any of the comments in a snippet is toxic.

Annotators	Conversations	Agreement	Annotators	Conversations	Snippets	Agreement
367	4,022	67.8%	247	1,252	2,181	87.5%

Table 1: Descriptions of crowdsourcing jobs, with relevant statistics. More details in Appendix A.

**Candidate selection.** Our goal is to analyze how the start of a *civil* conversation is tied to its potential future derailment into personal attacks. Thus, we only consider conversations that start out as ostensibly civil, i.e., where at least the first exchange does not exhibit any toxic behavior,<sup>4</sup> and that continue beyond this first exchange. To focus on the especially perplexing cases when the attacks come *from within*, we seek examples where the attack is initiated by one of the two participants in the initial exchange.

To select candidate conversations to include in our collection, we use the toxicity classifier provided by the Perspective API,<sup>5</sup> which is trained on Wikipedia talk page comments that have been annotated by crowdworkers (Wulczyn et al., 2016). This provides a toxicity score  $t$  for all comments in our dataset, which we use to preselect two sets of conversations: (a) candidate conversations that are civil throughout, i.e., conversations in which all comments (including the initial exchange) are not labeled as toxic ( $t < 0.4$ ); and (b) candidate conversations that turn toxic after the first (civil) exchange, i.e., conversations in which the  $N$ -th comment ( $N > 2$ ) is labeled toxic ( $t \geq 0.6$ ), but all the preceding comments are not ( $t < 0.4$ ).

**Crowdsourced filtering.** Starting from these candidate sets, we use crowdsourcing to vet each conversation and select a subset that are perceived by humans to either stay civil throughout (“on-track” conversations), or start civil but end with a *personal attack* (“awry-turning” conversations). To inform the design of this human-filtering process and to check its effectiveness, we start from a seed set of 232 conversations manually verified by the authors to end in personal attacks (more details about the selection of the seed set and its role in the crowd-sourcing process can be found in Appendix A). We take particular care to not over-constrain crowdworker interpretations of

what personal attacks may be, and to separate toxicity from civil disagreement, which is recognized as a key aspect of effective collaborations (Coser, 1956; De Dreu and Weingart, 2003).

We design and deploy two filtering jobs using the CrowdFlower platform, summarized in Table 1 and detailed in Appendix A. **Job 1** is designed to select conversations that contain a “rude, insulting, or disrespectful” comment towards another user in the conversation—i.e., a personal attack. In contrast to prior work labeling antisocial comments in isolation (Sood et al., 2012; Wulczyn et al., 2017), annotators are asked to label personal attacks in the *context* of the conversations in which they occur, since antisocial behavior can often be context-dependent (Cheng et al., 2017). In fact, in order to ensure that the crowdworkers read the entire conversation, we also ask them to indicate who is the target of the attack. We apply this task to the set of candidate awry-turning conversations, selecting the 14% which all three annotators perceived as ending in a personal attack.<sup>6</sup>

**Job 2** is designed to filter out conversations that do not actually start out as civil. We run this job to ensure that the *awry-turning* conversations are civil up to the point of the attack—i.e., they *turn awry*—discarding 5% of the candidates that passed Job 1. We also use it to verify that the candidate *on-track* conversations are indeed civil throughout, discarding 1% of the respective candidates. In both cases we filter out conversations in which three annotators could identify at least one comment that is “rude, insulting, or disrespectful”.

**Controlled setting.** Finally, we need to construct a setting that affords for meaningful comparison between conversations that derail and those that stay on track, and that accounts for trivial topical confounds (Kittur et al., 2009; Cheng et al., 2015). We mitigate topical confounds using matching, a technique developed for causal inference in observational studies (Rubin, 2007). Specifically, start-

<sup>4</sup>For the sake of generality, in this work we focus on this most basic conversational unit: the first comment-reply pair starting a conversation.

<sup>5</sup><https://www.perspectiveapi.com/>

<sup>6</sup>We opted to use unanimity in this task to account for the highly subjective nature of the phenomenon.

ing from our human-vetted collection of conversations, we pair each *awry-turning* conversation, with an *on-track* conversation, such that both took place on the same talk page. If we find multiple such pairs, we only keep the one in which the paired conversations take place closest in time, to tighten the control for topic. Conversations that cannot be paired are discarded.

This procedure yields a total of 1,270 paired awry-turning and on-track conversations (including our initial seed set), spanning 582 distinct talk pages (averaging 1.1 pairs per page, maximum 8) and 1,876 (overlapping) topical categories. The average length of a conversation is 4.6 comments.

## 4 Capturing Pragmatic Devices

We now describe our framework for capturing linguistic cues that might inform a conversation’s future trajectory. Crucially, given our focus on conversations that start seemingly civil, we do not expect overtly hostile language—such as insults (Yin et al., 2009)—to be informative. Instead, we seek to identify pragmatic markers within the initial exchange of a conversation that might serve to reveal or exacerbate underlying tensions that eventually come to the fore, or conversely suggest sustainable civility. In particular, in this work we explore how politeness strategies and rhetorical prompts reflect the future health of a conversation.

**Politeness strategies.** Politeness can reflect a-priori good will and help navigate potentially face-threatening acts (Goffman, 1955; Lakoff, 1973), and also offers hints to the underlying intentions of the interlocutors (Fraser, 1980). Hence, we may naturally expect certain politeness strategies to signal that a conversation is likely to stay on track, while others might signal derailment.

In particular, we consider a set of pragmatic devices signaling politeness drawn from Brown and Levinson (1987). These linguistic features reflect two overarching types of politeness. *Positive* politeness strategies encourage social connection and rapport, perhaps serving to maintain cohesion throughout a conversation; such strategies include gratitude (“*thanks* for your help”), greetings (“*hey*, how is your day so far”) and use of “please”, both at the start (“*Please* find sources for your edit...”) and in the middle (“Could you *please* help with...?”) of a sentence. *Negative* politeness strategies serve to dampen an interlocutor’s imposition on an addressee, often through conveying

indirectness or uncertainty on the part of the commenter. Both commenters in example **B** (Fig. 1) employ one such strategy, hedging, perhaps seeking to soften an impending disagreement about a source’s reliability (“*I don’t think...*”, “*I would assume...*”). We also consider markers of *impolite* behavior, such as the use of direct questions (“*Why’s* there no mention of it?”) and sentence-initial second person pronouns (“*Your* sources don’t matter...”), which may serve as forceful-sounding contrasts to negative politeness markers. Following Danescu-Niculescu-Mizil et al. (2013), we extract such strategies by pattern matching on the dependency parses of comments.

**Types of conversation prompts.** To complement our pre-defined set of politeness strategies, we seek to capture domain-specific rhetorical patterns used to initiate conversations. For instance, in a collaborative setting, we may expect conversations that start with an invitation for working together to signal less tension between the participants than those that start with statements of dispute. We discover types of such *conversation prompts* in an unsupervised fashion by extending a framework used to infer the rhetorical role of questions in (offline) political debates (Zhang et al., 2017b) to more generally extract the rhetorical functions of comments. The procedure follows the intuition that the rhetorical role of a comment is reflected in the type of replies it is likely to elicit. As such, comments which tend to trigger similar replies constitute a particular type of prompt.

To implement this intuition, we derive two different low-rank representations of the common lexical phrasings contained in comments (agnostic to the particular topical content discussed), automatically extracted as recurring sets of arcs in the dependency parses of comments. First, we derive *reply-vectors* of phrasings, which reflect their propensities to *co-occur*. In particular, we perform singular value decomposition on a term-document matrix  $\mathcal{R}$  of phrasings and replies as  $\mathcal{R} \approx \hat{\mathcal{R}} = U_R S V_R^T$ , where rows of  $U_R$  are low-rank reply-vectors for each phrasing.

Next, we derive *prompt-vectors* for the phrasings, which reflect similarities in the subsequent replies that a phrasing *prompts*. We construct a prompt-reply matrix  $\mathcal{P} = (p_{ij})$  where  $p_{ij} = 1$  if phrasing  $j$  occurred in a reply to a comment containing phrasing  $i$ . We project  $\mathcal{P}$  into the same space as  $U_R$  by solving for  $\hat{\mathcal{P}}$  in  $\mathcal{P} = \hat{\mathcal{P}} S V_R^T$  as

Prompt Type	Description	Examples
Factual check	Statements about article content, pertaining to or contending issues like factual accuracy.	The terms <b>are used</b> interchangeably in the US. The census <b>is not talking about</b> families here.
Moderation	Rebukes or disputes concerning moderation decisions such as blocks and reversions.	<b>If</b> you continue, you may <b>be blocked</b> from editing. He’s <b>accused</b> me <b>of</b> being a troll.
Coordination	Requests, questions, and statements of intent pertaining to collaboratively editing an article.	It’s a long list so I <b>could do with</b> your <b>help</b> . <b>Let me know</b> if you agree with this and I’ll go ahead [...]
Casual remark	Casual, highly conversational aside-remarks.	<b>What’s with</b> this flag image? <b>I’m surprised</b> there wasn’t an article before.
Action statement	Requests, statements, and explanations about various editing actions.	<b>Please consider improving</b> the article to address the issues [...] The page <b>was deleted as</b> self-promotion.
Opinion	Statements seeking or expressing opinions about editing challenges and decisions.	I <b>think</b> that it <b>should be</b> the other way around. This article <b>seems to have</b> a lot of bias.

Table 2: Prompt types automatically extracted from talk page conversations, with interpretations and examples from the data. Bolded text indicate common prompt phrasings extracted by the framework. Further examples are shown in Appendix B, Table 4.

$\hat{\mathcal{P}} = \mathcal{P}V_R S^{-1}$ . Each row of  $\hat{\mathcal{P}}$  is then a prompt-vector of a phrasing, such that the prompt-vector for phrasing  $i$  is close to the reply-vector for phrasing  $j$  if comments with phrasing  $i$  tend to prompt replies with phrasing  $j$ . Clustering the rows of  $\hat{\mathcal{P}}$  then yields  $k$  conversational *prompt types* that are unified by their similarity in the space of replies. To infer the prompt type of a new comment, we represent the comment as an average of the representations of its constituent phrasings (i.e., rows of  $\hat{\mathcal{P}}$ ) and assign the resultant vector to a cluster.<sup>7</sup>

To determine the prompt types of comments in our dataset, we first apply the above procedure to derive a set of prompt types from a *disjoint* (unlabeled) corpus of Wikipedia talk page conversations (Danescu-Niculescu-Mizil et al., 2012). After initial examination of the framework’s output on this external data, we chose to extract  $k = 6$  prompt types, shown in Table 2 along with our interpretations.<sup>8</sup> These prompts represent signatures of conversation-starters spanning a wide range of topics and contexts which reflect core elements of Wikipedia, such as moderation disputes and coordination (Kittur et al., 2007; Kittur and Kraut, 2008). We assign each comment in our present dataset to one of these types.<sup>9</sup>

<sup>7</sup>We scale rows of  $U_R$  and  $\hat{\mathcal{P}}$  to unit norm. We assign comments whose vector representation has ( $\ell_2$ ) distance  $\geq 1$  to all cluster centroids to an extra, infrequently-occurring null type which we ignore in subsequent analyses.

<sup>8</sup>We experimented with more prompt types as well, finding that while the methodology recovered finer-grained types, and obtained qualitatively similar results and prediction accuracies as described in Sections 5 and 6, the assignment of comments to types was relatively sparse due to the small data size, resulting in a loss of statistical power.

<sup>9</sup>While the particular prompt types we discover are spe-

## 5 Analysis

We are now equipped to computationally explore how the pragmatic devices used to start a conversation can signal its future health. Concretely, to quantify the relative propensity of a linguistic marker to occur at the start of awry-turning versus on-track conversations, we compute the log-odds ratio of the marker occurring in the initial exchange—i.e., in the first or second comments—of awry-turning conversations, compared to initial exchanges in the on-track setting. These quantities are depicted in Figure 2A.<sup>10</sup>

Focusing on the **first** comment (represented as  $\diamond$ s), we find a rough correspondence between linguistic *directness* and the likelihood of future personal attacks. In particular, comments which contain *direct questions*, or exhibit *sentence-initial you* (i.e., “2<sup>nd</sup> person start”), tend to start awry-turning conversations significantly more often than ones that stay on track (both  $p < 0.001$ ).<sup>11</sup> This effect coheres with our intuition that directness signals some latent hostility from the conversation’s initiator, and perhaps reinforces the forcefulness of contentious impositions (Brown and Levinson, 1987). This interpretation is also sug-

cific to Wikipedia, the methodology for inferring them is unsupervised and is applicable in other conversational settings.

<sup>10</sup>To reduce clutter we only depict features which occur a minimum of 50 times and have absolute log-odds  $\geq 0.2$  in at least one of the data subsets. The markers indicated as statistically significant for Figure 2A remain so after a Bonferroni correction, with the exception of factual checks, hedges (lexicon,  $\diamond$ ), gratitude ( $\diamond$ ), and opinion.

<sup>11</sup>All  $p$  values in this section are computed as two-tailed binomial tests, comparing the proportion of awry-turning conversations exhibiting a particular device to the proportion of on-track conversations.

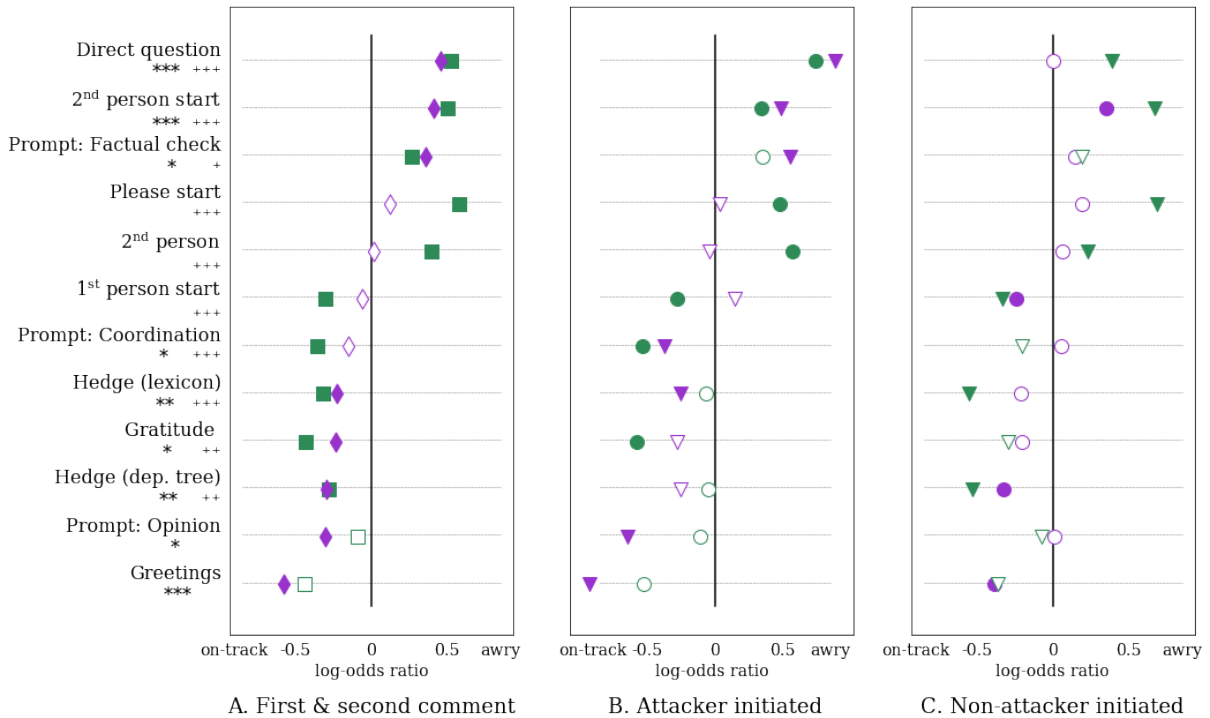


Figure 2: Log-odds ratios of politeness strategies and prompt types exhibited in the first and second comments of conversations that turn awry, versus those that stay on-track. **All:** Purple and green markers denote log-odds ratios in the first and second comments, respectively; points are solid if they reflect significant ( $p < 0.05$ ) log-odds ratios with an effect size of at least 0.2. **A:**  $\diamond$ s and  $\square$ s denote **first** and **second** comment log-odds ratios, respectively; \* denotes statistically significant differences at the  $p < 0.05$  (\*),  $p < 0.01$  (\*\*), and  $p < 0.001$  (\*\*\*) levels for the first comment (two-tailed binomial test); + denotes corresponding statistical significance for the second comment. **B** and **C:**  $\nabla$ s and  $\circ$ s correspond to effect sizes in the comments authored by the **attacker** and **non-attacker**, respectively, in **attacker initiated** (B) and **non-attacker initiated** (C) conversations.

gested by the relative propensity of the *factual check* prompt, which tends to cue disputes regarding an article’s factual content ( $p < 0.05$ ).

In contrast, comments which initiate on-track conversations tend to contain *gratitude* ( $p < 0.05$ ) and *greetings* ( $p < 0.001$ ), both positive politeness strategies. Such conversations are also more likely to begin with *coordination* prompts ( $p < 0.05$ ), signaling active efforts to foster constructive teamwork. Negative politeness strategies are salient in on-track conversations as well, reflected by the use of *hedges* ( $p < 0.01$ ) and *opinion* prompts ( $p < 0.05$ ), which may serve to soften impositions or factual contentions (Hübler, 1983).

These effects are echoed in the **second** comment—i.e., the **first reply** (represented as  $\square$ s). Interestingly, in this case we note that the difference in pronoun use is especially marked. First replies in conversations that eventually de-

rail tend to contain more *second person pronouns* ( $p < 0.001$ ), perhaps signifying a replier pushing back to contest the initiator; in contrast, on-track conversations have more *sentence-initial I/We* (i.e., “1<sup>st</sup> person start”,  $p < 0.001$ ), potentially indicating the replier’s willingness to step into the conversation and work with—rather than argue against—the initiator (Tausczik and Pennebaker, 2010).

**Distinguishing interlocutor behaviors.** Are the linguistic signals we observe solely driven by the eventual attacker, or do they reflect the behavior of both actors? To disentangle the attacker and non-attackers’ roles in the initial exchange, we examine their language use in these two possible cases: when the *future* attacker initiates the conversation, or is the first to reply. In **attacker-initiated** conversations (Figure 2B, 608 conversations), we see that both actors exhibit a propensity for the linguistically direct markers (e.g., *direct questions*)

that tend to signal future attacks. Some of these markers are used particularly often by the **non-attacking replier** in awry-turning conversations (e.g., *second person pronouns*,  $p < 0.001$ , ○s), further suggesting the dynamic of the replier pushing back at—and perhaps even escalating—the attacker’s initial hint of aggression. Among conversations initiated instead by the **non-attacker** (Figure 2C, 662 conversations), the non-attacker’s linguistic behavior in the first comment (○s) is less distinctive from that of initiators in the on-track setting (i.e., log-odds ratios closer to 0); markers of future derailment are (unsurprisingly) more pronounced once the eventual attacker (▽s) joins the conversation in the second comment.<sup>12</sup>

More broadly, these results reveal how different politeness strategies and rhetorical prompts deployed in the initial stages of a conversation are tied to its future trajectory.

## 6 Predicting Future Attacks

We now show that it is indeed feasible to predict whether a conversation will turn awry based on linguistic properties of its very first exchange, providing several baselines for this new task. In doing so, we demonstrate that the pragmatic devices examined above encode signals about the future trajectory of conversations, capturing some of the intuition humans are shown to have.

We consider the following balanced prediction task: given a pair of conversations, which one will eventually lead to a personal attack? We extract all features from the very first exchange in a conversation—i.e., a comment-reply pair, like those illustrated in our introductory example (Figure 1). We use logistic regression and report accuracies on a leave-one-page-out cross validation, such that in each fold, all conversation pairs from a given talk page are held out as test data and pairs from all other pages are used as training data (thus preventing the use of page-specific information). Prediction results are summarized in Table 3.

**Language baselines.** As baselines, we consider several straightforward features: word count (which performs at chance level), sentiment lexicon (Liu et al., 2005) and bag of words.

**Pragmatic features.** Next, we test the predictive power of the **prompt types** and **politeness**

<sup>12</sup>As an interesting avenue for future work, we note that some markers used by non-attacking initiators potentially still anticipate later attacks, suggested by, e.g., the relative prevalence of *sentence-initial you* ( $p < 0.05$ , ○s).

Feature set	# features	Accuracy
Bag of words	5,000	56.7%
Sentiment lexicon	4	55.4%
<b>Politeness strategies</b>	38	60.5%
<b>Prompt types</b>	12	59.2%
<b>Pragmatic (all)</b>	50	61.6%
<i>Interlocutor features</i>	5	51.2%
<i>Trained toxicity</i>	2	60.5%
<i>Toxicity + Pragmatic</i>	52	64.9%
<i>Humans</i>		72.0%

Table 3: Accuracies for the balanced future-prediction task. Features based on pragmatic devices are **bolded**, reference points are *italicized*.

**strategies** features introduced in Section 4. The 12 prompt type features (6 features for each comment in the initial exchange) achieve 59.2% accuracy, and the 38 politeness strategies features (19 per comment) achieve 60.5% accuracy. The **pragmatic** features combine to reach 61.6% accuracy.

**Reference points.** To better contextualize the performance of our features, we compare their predictive accuracy to the following reference points: *Interlocutor features*: Certain kinds of interlocutors are potentially more likely to be involved in awry-turning conversations. For example, perhaps newcomers or anonymous participants are more likely to derail interactions than more experienced editors. We consider a set of features representing participants’ experience on Wikipedia (i.e., number of edits) and whether the comment authors are anonymous. In our task, these features perform at the level of random chance.

*Trained toxicity*: We also compare with the toxicity score of the exchange from the Perspective API classifier—a perhaps unfair reference point, since this supervised system was trained on additional human-labeled training examples from the same domain and since it was used to create the very data on which we evaluate. This results in an accuracy of 60.5%; combining trained toxicity with our pragmatic features achieves 64.9%.

*Humans*: A sample of 100 pairs were labeled by (non-author) volunteer human annotators. They were asked to guess, from the initial exchange, which conversation in a pair will lead to a personal attack. Majority vote across three annotators was used to determine the human labels, resulting in an accuracy of 72%. This confirms that humans have



some intuition about whether a conversation might be heading in a bad direction, which our features can partially capture. In fact, the classifier using pragmatic features is accurate on 80% of the examples that humans also got right.

**Attacks on the horizon.** Finally, we seek to understand whether cues extracted from the first exchange can predict future discussion trajectory beyond the immediate next couple of comments. We thus repeat the prediction experiments on the subset of conversations in which the first personal attack happens after the fourth comment (282 pairs), and find that the pragmatic devices used in the first exchange maintain their predictive power (67.4% accuracy), while the sentiment and bag of words baselines drop to the level of random chance.

Overall, these initial results show the feasibility of reconstructing some of the human intuition about the future trajectory of an ostensibly civil conversation in order to predict whether it will eventually turn awry.

## 7 Conclusions and Future Work

In this work, we started to examine the intriguing phenomenon of conversational derailment, studying how the use of pragmatic and rhetorical devices relates to future conversational failure. Our investigation centers on the particularly perplexing scenario in which one participant of a civil discussion later attacks another, and explores the new task of predicting whether an initially healthy conversation will derail into such an attack. To this end, we develop a computational framework for analyzing how general politeness strategies and domain-specific rhetorical prompts deployed in the initial stages of a conversation are tied to its future trajectory.

Making use of machine learning and crowdsourcing tools, we formulate a tightly-controlled setting that enables us to meaningfully compare conversations that stay on track with those that go awry. The human accuracy on predicting future attacks in this setting (72%) suggests it is feasible at least at the level of human intuition. We show that our computational framework can recover some of that intuition, hinting at the potential of automated methods to identify signals of the future trajectories of online conversations.

Our approach has several limitations which open avenues for future work. Our correlational analyses do not provide any insights into *causal*

mechanisms of derailment, which randomized experiments could address. Additionally, since our procedure for collecting and vetting data focused on precision rather than recall, it might miss more subtle attacks that are overlooked by the toxicity classifier. Supplementing our investigation with other indicators of antisocial behavior, such as editors blocking one another, could enrich the range of attacks we study. Noting that our framework is not specifically tied to Wikipedia, it would also be valuable to explore the varied ways in which this phenomenon arises in other (possibly non-collaborative) public discussion venues, such as Reddit and Facebook Pages.

While our analysis focused on the very first exchange in a conversation for the sake of generality, more complex modeling could extend its scope to account for conversational features that more comprehensively span the interaction. Beyond the present binary classification task, one could explore a sequential formulation predicting whether the next turn is likely to be an attack as a discussion unfolds, capturing conversational dynamics such as sustained escalation.

Finally, our study of derailment offers only one glimpse into the space of possible conversational trajectories. Indeed, a manual investigation of conversations whose eventual trajectories were misclassified by our models—as well as by the human annotators—suggests that interactions which initially seem prone to attacks can nonetheless maintain civility, by way of level-headed interlocutors, as well as explicit acts of reparation. A promising line of future work could consider the complementary problem of identifying pragmatic strategies that can help bring uncivil conversations back on track.

**Acknowledgements.** We are grateful to the anonymous reviewers for their thoughtful comments and suggestions, and to Maria Antoniak, Valts Blukis, Liye Fu, Sam Havron, Jack Hessel, Ishaan Jhaveri, Lillian Lee, Alex Niculescu-Mizil, Alexandra Schofield, Laure Thompson, Andrew Wang, Leila Zia and the members of the Wikimedia Foundation anti-harassment program for extremely insightful (on-track) conversations and for assisting with data annotation efforts. This work is supported in part by NSF CAREER Award IIS-1750615, NSF Grant SES-1741441, a Google Faculty Award, a WMF gift and a CrowdFlower AI for Everyone Award.

## References

- Yavuz Akbulut, Yusuf Levent Sahin, and Bahadir Eristi. 2010. Cyberbullying victimization among Turkish online social utility members. *Journal of Educational Technology & Society*.
- Kelsey Allen, Giuseppe Carenini, and Raymond T Ng. 2014. Detecting disagreement in conversations using pseudo-monologic rhetorical structure. In *Proceedings of EMNLP*.
- Tim Althoff, Cristian Danescu-Niculescu-Mizil, and Dan Jurafsky. 2014. How to ask for a favor: A case study on the success of altruistic requests. In *Proceedings of ICWSM*.
- Ofer Arazy, Lisa Yeo, and Oded Nov. 2013. Stay on the Wikipedia task: When task-related disagreements slip into personal and procedural conflicts. *Journal of the Association for Information Science and Technology*.
- Malika Aubakirova and Mohit Bansal. 2016. Interpreting neural networks to improve politeness comprehension. In *Proceedings of EMNLP*.
- Lars Backstrom, Jon Kleinberg, Lillian Lee, and Cristian Danescu-Niculescu-Mizil. 2013. Characterizing and curating conversation threads: Expansion, focus, volume, re-entry. In *Proceedings of WSDM*.
- Penelope Brown and Stephen Levinson. 1987. *Politeness: Some universals in language usage*. Cambridge University Press.
- Moira Burke and Robert Kraut. 2008. Mind your Ps and Qs: The impact of politeness and rudeness in online communities. In *Proceedings of CSCW*.
- Eshwar Chandrasekharan, Umashanthi Pavalanathan, Anirudh Srinivasan, Adam Glynn, Jacob Eisenstein, and Eric Gilbert. 2017. You can't stay here: The efficacy of Reddit's 2015 ban examined through hate speech. In *Proceedings of CSCW*.
- Despoina Chatzakou, Nicolas Kourtellis, Jeremy Blackburn, Emiliano De Cristofaro, Gianluca Stringhini, and Athena Vakali. 2017. Measuring #GamerGate: A tale of hate, sexism, and bullying. In *Proceedings of WWW*.
- Justin Cheng, Michael Bernstein, Cristian Danescu-Niculescu-Mizil, and Jure Leskovec. 2017. Anyone can become a troll: Causes of trolling behavior in online discussions. In *Proceedings of CSCW*.
- Justin Cheng, Cristian Danescu-Niculescu-Mizil, and Jure Leskovec. 2015. Antisocial behavior in online discussion communities. In *Proceedings of ICWSM*.
- Elizabeth F Churchill and Sara Bly. 2000. Culture virtues: Considering culture and communication in virtual environments. *SIGGroup Bulletin*.
- Herbert H Clark. 1979. Responding to indirect speech acts. *Cognitive psychology*.
- Herbert H Clark and Dale H Schunk. 1980. Polite responses to polite requests. *Cognition*.
- Benjamin Collier and Julia Bear. 2012. Conflict, criticism, or confidence: An empirical examination of the gender gap in Wikipedia contributions. In *Proceedings of CSCW*.
- Lewis A Coser. 1956. *The Functions of Social Conflict*. Routledge.
- Cristian Danescu-Niculescu-Mizil, Lillian Lee, Bo Pang, and Jon Kleinberg. 2012. Echoes of power: Language effects and power differences in social interaction. In *Proceedings of WWW*.
- Cristian Danescu-Niculescu-Mizil, Moritz Sudhof, Dan Jurafsky, Jure Leskovec, and Christopher Potts. 2013. A computational approach to politeness with application to social factors. In *Proceedings of ACL*.
- Thomas Davidson, Dana Warmusley, Michael Macy, and Ingmar Weber. 2017. Automated hate speech detection and the problem of offensive language. In *Proceedings of ICWSM*.
- Carsten K De Dreu and Laurie R Weingart. 2003. Task versus relationship conflict, team performance, and team member satisfaction: A meta-analysis. *Journal of Applied Psychology*.
- Bruce Fraser. 1980. Conversational mitigation. *Journal of Pragmatics*.
- Liye Fu, Lillian Lee, and Cristian Danescu-Niculescu-Mizil. 2017. When confidence and competence collide: Effects on online decision-making discussions. In *Proceedings of WWW*.
- Björn Gambäck and Utpal Kumar Sikdar. 2017. Using convolutional neural networks to classify hate-speech. In *Proceedings of the Workshop on Abusive Language Online*.
- Ali Ghehtasy, José Abdelnour-Nocera, and Bonnie Nardi. 2015. Socio-technical gaps in online collaborative consumption (OCC): An example of the Etsy community. In *Proceedings of ICDC*.
- Erving Goffman. 1955. On face-work: An analysis of ritual elements in social interaction. *Psychiatry*.
- Christophe Henner and Maria Sefidari. 2016. [Wikimedia Foundation Board on healthy Wikimedia community culture, inclusivity, and safe spaces](#). Wikimedia Blog.
- Pamela J Hinds and Mark Mortensen. 2005. Understanding conflict in geographically distributed teams: The moderating effects of shared identity, shared context, and spontaneous communication. *Organization Science*.
- Axel Hübler. 1983. *Understatements and Hedges in English*. John Benjamins Publishing.

- Joseph M Kayany. 1998. Contexts of uninhibited on-line behavior: Flaming in social newsgroups on usenet. *Journal of the Association for Information Science and Technology*.
- Aniket Kittur, Ed H Chi, and Bongwon Suh. 2009. What’s in Wikipedia?: Mapping topics and conflict using socially annotated category structure. In *Proceedings of CHI*.
- Aniket Kittur and Robert E Kraut. 2008. Harnessing the wisdom of crowds in Wikipedia: Quality through coordination. In *Proceedings of CSCW*.
- Aniket Kittur, Bongwon Suh, Bryan A Pendleton, and Ed H Chi. 2007. He says, she says: Conflict and coordination in Wikipedia. In *Proceedings of CHI*.
- Vinodh Krishnan and Jacob Eisenstein. 2015. “You’re Mr. Lebowksi, I’m the Dude”: Inducing address term formality in signed social networks. In *Proceedings of NAACL*.
- Ravi Kumar, Mohammad Mahdian, and Mary McGlohon. 2010. Dynamics of conversations. In *Proceedings of KDD*.
- Haewoon Kwak, Jeremy Blackburn, and Seungyeop Han. 2015. Exploring cyberbullying and other toxic behavior in team competition online games. In *Proceedings of CHI*.
- Robin T Lakoff. 1973. The logic of politeness: Minding your P’s and Q’s. In *Proceedings of the Chicago Linguistic Society*.
- Bing Liu, Minqing Hu, and Junsheng Cheng. 2005. Opinion observer: Analyzing and comparing opinions on the web. In *Proceedings of WWW*.
- Vlad Niculae and Cristian Danescu-Niculescu-Mizil. 2016. Conversational markers of constructive discussions. In *Proceedings of NAACL*.
- Vlad Niculae, Srijan Kumar, Jordan Boyd-Graber, and Cristian Danescu-Niculescu-Mizil. 2015. Linguistic harbingers of betrayal: A case study on an online strategy game. In *Proceedings of ACL*.
- Chikashi Nobata, Joel Tetreault, Achint Thomas, Yashar Mehdad, and Yi Chang. 2016. Abusive language detection in online user content. In *Proceedings of WWW*.
- Marco Ortu, Bram Adams, Giuseppe Destefanis, Paras-tou Tourani, Michele Marchesi, and Roberto Tonelli. 2015. Are bullies more productive? Empirical study of affectiveness vs. issue fixing time. In *Proceedings of MSR*.
- John Pavlopoulos, Prodromos Malakasiotis, and Ion Androutsopoulos. 2017a. Deep learning for user comment moderation. In *Proceedings of the Workshop on Abusive Language Online*.
- John Pavlopoulos, Prodromos Malakasiotis, and Ion Androutsopoulos. 2017b. Deeper attention to abusive user content moderation. In *Proceedings of EMNLP*.
- Alan Ritter, Colin Cherry, and Bill Dolan. 2010. Unsupervised modeling of Twitter conversations. In *Proceedings of NAACL*.
- Paul R Rosenbaum. 2010. *Design of observational studies*. Springer.
- Donald B Rubin. 2007. The design versus the analysis of observational studies for causal effects: Parallels with the design of randomized trials. *Statistics in Medicine*.
- Tamara Shepherd, Alison Harvey, Tim Jordan, Sam Srauy, and Kate Miltner. 2015. Histories of hating. *Social Media + Society*.
- Vivek K Singh, Marie L Radford, Qianjia Huang, and Susan Furrer. 2017. “They basically like destroyed the school one day”: On newer app features and cyberbullying in schools. In *Proceedings of CSCW*.
- Sara Owsley Sood, Elizabeth F Churchill, and Judd Antin. 2012. Automatic identification of personal insults on social news sites. *Journal of the American Society for Information Science and Technology*.
- Chenhao Tan, Vlad Niculae, Cristian Danescu-Niculescu-Mizil, and Lillian Lee. 2016. Winning arguments: Interaction dynamics and persuasion strategies in good-faith online discussions. In *Proceedings of WWW*.
- Yla R Tausczik and James W Pennebaker. 2010. The psychological meaning of words: LIWC and computerized text analysis methods. *Journal of Language and Social Psychology*.
- Kal Turnbull. 2018. “Thats Bullshit” – Rude Enough for Removal? A Multi-Mod Perspective. Change My View Blog.
- Jessica Vitak, Kalyani Chadha, Linda Steiner, and Zahra Ashktorab. 2017. Identifying women’s experiences with and strategies for mitigating negative effects of online harassment. In *Proceedings of CSCW*.
- Lu Wang and Claire Cardie. 2014. A piece of my mind: A sentiment analysis approach for online dispute detection. In *Proceedings of ACL*.
- William Warner and Julia Hirschberg. 2012. Detecting hate speech on the World Wide Web. In *Proceedings of the Workshop on Language in Social Media*.
- Wikimedia Support and Safety Team. 2015. [Harassment survey](#). Wikimedia Foundation.
- Ellery Wulczyn, Nithum Thain, and Lucas Dixon. 2016. [Wikipedia talk labels: Toxicity](#).

- Ellery Wulczyn, Nithum Thain, and Lucas Dixon. 2017. Ex machina: Personal attacks seen at scale. In *Proceedings of WWW*.
- Naomi Yamashita and Toru Ishida. 2006. Automatic prediction of misconceptions in multilingual computer-mediated communication. In *Proceedings of IUI*.
- Dawei Yin, Zhenzhen Xue, Liangjie Hong, Brian D Davison, April Kontostathis, and Lynne Edwards. 2009. Detection of harassment on Web 2.0. In *Proceedings of the Workshop on Content Analysis in the Web 2.0*.
- Amy X Zhang, Bryan Culbertson, and Praveen Paritosh. 2017a. Characterizing online discussion using coarse discourse sequences. In *Proceedings of ICWSM*.
- Justine Zhang, Ravi Kumar, Sujith Ravi, and Cristian Danescu-Niculescu-Mizil. 2016. Conversational flow in Oxford-style debates. In *Proceedings of NAACL*.
- Justine Zhang, Arthur Spirling, and Cristian Danescu-Niculescu-Mizil. 2017b. Asking too much? The rhetorical role of questions in political discourse. In *Proceedings of EMNLP*.

## Appendix

### A Details on annotation procedure

The process of constructing a labeled dataset for personal attacks was challenging due to the complex and subjective nature of the phenomenon, and developed over several iterations as a result. In order to guide future work, here we provide a detailed explanation of this process, expanding on the description in Section 3.

Our goal in this work was to understand linguistic markers of conversations that go awry and devolve into personal attacks—a highly subjective phenomenon with a multitude of possible definitions.<sup>13</sup> To enable a concrete analysis of conversational derailment that encompasses the scale and diversity of a setting like Wikipedia talk pages, we therefore needed to develop a well-defined conceptualization of conversational failure, and a procedure to accurately discover instances of this phenomenon at scale.

Our approach started from an initial qualitative investigation that resulted in a seed set of example conversational failures. This seed set then informed the design of the subsequent crowdsourced filtering procedure, which we used to construct our full dataset.

#### Initial qualitative investigation

To develop our task, we compiled an initial sample of potentially awry-turning conversations by applying the candidate selection procedure (detailed in Section 3) to a random subset of Wikipedia talk pages. This procedure yielded a set of conversations which the underlying trained classifier deemed to be initially civil, but with a later toxic comment. An informal inspection of these candidate conversations suggested many possible forms of toxic behavior, ranging from personal attacks (‘Are you that big of a coward?’), to uncivil disagreements (‘Read the previous discussions before bringing up this stupid suggestion again.’), to generalized attacks (‘Another left wing inquisition?’) and even to outright vandalism (‘Wikipedia SUCKS!’) or simply unnecessary use of foul language.

Through our manual inspection, we also identified a few salient points of divergence between the classifier and our (human) judgment of toxic-

ity. In particular, several comments which were machine-labeled as toxic were clearly sarcastic or self-deprecating, perhaps employing seemingly aggressive or foul language to bolster the collegial nature of the interaction rather than to undermine it. These false positive instances highlight the necessity of the subsequent crowdsourced vetting process—and point to opportunities to enrich the subtle linguistic and interactional cues such classifiers can address.

**Seed set.** Our initial exploration of the automatically discovered candidate conversations and our discussions with the members of the Wikimedia Foundation anti-harassment program pointed to a particularly salient and perplexing form of toxic behavior around which we centered our subsequent investigation: personal attacks *from within*, where one of the two participants of the ostensibly civil initial exchange turns on another interlocutor. For each conversation where the author of the toxic-labeled comment also wrote the first or second comment, the authors manually checked that the interaction started civil and ended in a personal attack. The combined automatic and manual filtering process resulted in our seed set of 232 awry-turning conversations.

We additionally used the candidate selection procedure to obtain on-track counterparts to each conversation in the seed set that took place on the same talk-page; this pairing protocol is further detailed in Section 3.

**Human performance.** We gaged the feasibility of our task of predicting future personal attacks by asking (non-author) volunteer human annotators to label a 100-pair subset of the seed set. In this informal setting, also described in Section 6, we asked each annotator to guess which conversation in a pair will lead to a personal attack on the basis of the initial exchange. Taking the majority vote across three annotators, the human guesses achieved an accuracy of 72%, demonstrating that humans indeed have some systematic intuition for a conversation’s potential for derailment.

**Informing the crowdsourcing procedure.** To scale beyond the initial sample, we sought to use crowdworkers to replicate our process of manually filtering automatically-discovered candidates, enabling us to vet machine-labeled awry-turning and on-track conversations across the entire dataset. Starting from our seed set, we adopted an iterative approach to formulate our crowdsourcing tasks.

<sup>13</sup>Refer to [Turnbull \(2018\)](#) for examples of challenges community moderators face in delineating personal attacks.

In particular, we designed an initial set of task instructions—along with definitions and examples of personal attacks—based on our observations of the seed set. Additionally, we chose a subset of conversations from the seed set to use as *test questions* that crowdworker judgements on the presence or absence of such behaviors could be compared against. These test questions served both as anchors to ensure the clarity of our instructions, and as quality controls. Mismatches between crowdworker responses and our own labels in trial runs then motivated subsequent modifications we made to the task design. The crowdsourcing jobs we ultimately used to compile our entire dataset are detailed below.

### Crowdsourced filtering

Based on our experiences in constructing and examining the seed set, we designed a crowdsourcing procedure to construct a larger set of personal attacks. Here we provide more details about the crowdsourcing tasks, outlined in Section 3. We split the crowdsourcing procedure into two jobs, mirroring the manual process used to construct the seed set outlined above. The first job selected conversations ending with personal attacks; the second job enforced that awry-turning conversations start civil, and that on-track conversations remain civil throughout. We used the CrowdFlower platform to implement and deploy these jobs.

**Job 1: Ends in personal attack.** The first crowdsourcing job was designed to select conversations containing a personal attack. In the annotation interface, each of three annotators was shown a candidate awry-turning conversation (selected using the procedure described in Section 3). The suspected toxic comment was highlighted, and workers were asked whether the highlighted comment contains a personal attack—defined in the instructions as a comment that is “rude, insulting, or disrespectful towards a person/group or towards that person/group’s actions, comments, or work.” We instructed the annotators not to confuse personal attacks with civil disagreement, providing examples that illustrated this distinction.

To control the quality of the annotators and their responses, we selected 82 conversations from the seed set to use as *test questions* with a known label. Half of these test questions contained a personal attack and the other half were known to be civil. The CrowdFlower platform’s quality control

tools automatically blocked workers who missed at least 20% of these test questions.

While our task sought to identify personal attacks towards other interlocutors, trial runs of Job 1 suggested that many annotators construed attacks directed at other targets—such as groups or the Wikipedia platform in general—as personal attacks as well. To clarify the distinction between attack targets, and focus the annotators on labeling personal attacks, we asked annotators to specify *who* the target of the attack is: (a) someone else in the conversation, (b) someone outside the conversation, (c) a group, or (d) other. The resultant responses allowed us to filter annotations based on the reported target. This question also played the secondary role of ensuring that annotators read the entire conversation and accounted for this additional context in their choice.

In order to calibrate annotator judgements of what constituted an attack, we enforced that annotators saw a reasonable balance of awry-turning and on-track conversations. By virtue of the candidate selection procedure, a large proportion of the conversations in the candidate set contained attacks. Hence, we also included 804 candidate on-track conversations in the task.

Using the output of Job 1, we filtered our candidate set to the conversations where *all three annotations* agreed that a personal attack had occurred. We found that unanimity produced higher quality labels than taking a majority vote by omitting ambiguous cases (e.g., the comment “It’s our job to document things that have received attention, however ridiculous we find them.” could be insulting towards the things being documented, but could also be read as a statement of policy).<sup>14</sup>

**Job 2: Civil start.** The second crowdsourcing job was designed to enforce that candidate awry-turning conversations start civil, and candidate on-track conversations remain civil throughout. Each of three annotators was shown comments from both on-track and awry-turning conversations that had already been filtered through Job 1. They were asked whether any of the displayed comments were toxic—defined as “a rude, insulting, or disrespectful comment that is likely to make someone leave a discussion, engage in fights, or give up on sharing their perspective.” This definition was adapted from previous efforts to annotate toxic be-

---

<sup>14</sup>This choice further sacrifices recall for the sake of label precision, an issue that is also discussed in Section 7.

havior (Wulczyn et al., 2016) and intentionally targets a broader spectrum of uncivil behavior.

As in Job 1, we instructed annotators to not confound civil disagreement with toxicity. To reinforce this distinction, we included an additional question asking them whether any of the comments displayed disagreement, and prompted them to identify particular comments.

Since toxicity can be context-dependent, we wanted annotators to have access to the full conversation to help inform their judgement about each comment. However, we were also concerned that annotators would be overwhelmed by the amount of text in long conversations, and might be deterred from carefully reading each comment as a result. Indeed, in a trial run where full conversations were shown, we received negative feedback from annotators regarding task difficulty. To mitigate this difficulty without entirely omitting contextual information, we divided each conversation into snippets of three comments each. This kept the task fairly readable while still providing some local context. For candidate awry-turning conversations, we generated the snippets from all comments except the last one (which is known from Job 1 to be an attack). For on-track conversations, we generated the snippets from all comments in the conversation.

We marked conversations as toxic if at least three annotators, across all snippets of the conversation, identified at least one toxic comment. As in Job 1, we found that requiring this level of consensus among annotators produced reasonably high-quality labels.

**Overall flow.** To compile our full dataset, we started with 3,218 candidate awry-turning conversations which were filtered using Job 1, and discarded all but 435 conversations which all three annotators labeled as ending in a personal attack towards someone else in the conversation. These 435 conversations, along with paired on-track conversations, were then filtered using Job 2. This step removed 30 pairs: 24 where the awry-turning conversation was found to contain toxicity before the personal attack happened, and 6 where the on-track conversation was found to contain toxicity. We combined the crowdsourced output with the seed set to obtain a final dataset of 1,270 paired awry-turning and on-track conversations.

## B Further examples of prompt types

Table 4 provides further examples of comments containing the prompt types we automatically extracted from talk page conversations using the unsupervised methodology described in Section 4; descriptions of each type can be found in Table 2. For additional interpretability, we also include examples of typical *replies* to comments of each prompt type, which are also extracted by the method.

Prompt Type	Example comments	Example replies
Factual check	I <b>don't see</b> how this <b>is</b> relevant. That <b>does not mean</b> you can use this abbreviation everywhere. "Techniques" <b>refer</b> specifically to his fighting.	I <b>don't understand</b> your dispute. This <b>means</b> he <b>is</b> unlikely to qualify as an expert. They did <b>not believe</b> he will return. I <b>disagree</b> .
Moderation	<b>Please stop</b> making POV edits to the article. Your edits appear to be vandalism and have <b>been reverted</b> . I <b>have removed</b> edits which seem nationalistic. These mistakes should <b>not be allowed</b> to remain in the article.	I've <b>reverted</b> your change [...] I've <b>asked</b> them to stop. The next occurrence will <b>result in</b> a block. <b>Do not remove</b> my question.
Coordination	I <b>have been working on</b> creating an article. <b>Feel free</b> to correct my mistake. I <b>expanded</b> the article from a stub. I'll <b>make sure</b> to include a plot summary.	If you can do it I <b>would appreciate it</b> . I <b>have to go</b> but I'll be back later. <b>Ok, thanks</b> . <b>Hopefully</b> it will <b>be</b> fixed in a week.
Casual remark	<b>Just to</b> save you any issue in the future [...] <b>Remember</b> that badge I gave you? <b>Oh</b> , that's fabulous, <b>love</b> the poem! <b>Not sure</b> how that last revert came in there.	<b>Yeah</b> , this has <b>gotten</b> out of hand. <b>Anyway</b> , it's <b>nice</b> to see you took the time [...] <b>Yep</b> , that's <b>cool</b> . I <b>just thought</b> your comment was no longer needed.
Action statement	<b>If</b> you have <b>uploaded</b> other media, <b>consider</b> checking the criteria. <b>Could</b> somebody <b>please explain</b> how they differ? The info <b>was placed</b> in the appropriate section. Could you <b>undelete</b> my article?	That article has been <b>tagged for</b> deletion. I've <b>fixed</b> the wording. <b>Replaced with</b> free picture for all pages. It <b>has been deleted by</b> an admin.
Opinion	I've been <b>thinking of</b> setting up a portal. I <b>am wondering</b> if he is not supposed to be editing here. It's <b>hard</b> to combine these disputes.	It <b>seems</b> very much <b>in</b> the Wiki spirit. <b>Sounds like</b> a good idea. I <b>also think</b> we <b>need</b> to clarify this.

Table 4: Further examples of representative comments in the data for each automatically-extracted prompt type, along with examples of typical replies prompted by each type, produced by the methodology outlined in Section 4. Bolded text indicate common prompt and reply phrasings identified by the framework in the respective examples; note that the comment and reply examples in each row do not necessarily correspond to one another.