

How did we get here? Summarizing conversation dynamics

Anonymous ACL submission

Abstract

Throughout a conversation, the way participants interact with each other is in constant flux: their tones may change, they may resort to different strategies to convey their points, or they might alter their interaction patterns. An understanding of these *dynamics* can complement that of the actual facts and opinions discussed, offering a more holistic view of the trajectory of the conversation: how it arrived at its current state and where it is likely heading.

In this work, we introduce the task of summarizing the dynamics of conversations, by constructing a dataset of human-written summaries, and exploring several automated baselines. We evaluate whether such summaries can capture the trajectory of conversations via an established downstream task: forecasting whether an ongoing conversation will eventually derail into toxic behavior. We show that they help both humans and automated systems with this forecasting task. Humans make predictions three times faster, and with greater confidence, when reading the summaries than when reading the transcripts. Furthermore, automated forecasting systems are more accurate when constructing, and then predicting based on, summaries of conversation dynamics, compared to directly predicting on the transcripts.

1 Introduction

Conversations take place on at least two different levels (Tannen, 2005). On one level, participants directly communicate ideas, facts, and opinions, providing the topical context of the discussion. On the other level, the *dynamics* of their interactions reveal how they feel about each other, through changes in their tone—e.g., polite (Lakoff, 1973; Brown, 2015), condescending (Huckin, 2002), or sarcastic (Jorgensen, 1996)—, conversational strategies they employ—e.g., rhetorical questions (Han, 2002)—and even the patterns of their exchanges (Sacks et al., 1974; Silverstein, 1984).

Traditional summary: In this online conversation, participants discuss capitalism and its sustainability. Speaker1 argues that profit always trumps morals and ethics in business. Speaker2 disagrees, stating that unchecked capitalism is economically unsustainable and detrimental to human dignity. Speaker4 seeks clarification on the consequences of capitalism’s unsustainability, and Speaker2 suggests it could lead to economic collapse, mass unemployment, and toxic environments. They emphasize the need for checks on capitalism to prevent these issues. The conversation highlights concerns about unregulated capitalism and its impact on society.

Summary of conversation dynamics (human-written): Several users discuss regulation of capitalism. Speaker2 shares an opinion which Speaker4 questions. Speaker2 addresses Speaker4’s question *in a sincere manner*. Then, the *overall tone shifts to aggressive and confrontational* when Speaker4 *rudely invalidates* Speaker2’s response. Speaker2 *sarcastically criticizes* Speaker4’s attitude which aggravates Speaker4 more. Then, Speaker2 *genuinely apologizes and elaborates* on their point and Speaker4 *gratefully acknowledges* this and clarifies their intentions. They *continue the discussion in a calmed down, civil tone*.

Summary of conversation dynamics (generated): Four speakers engage in a discussion about the nature of capitalism and its consequences. Speaker2 expresses concerns about unchecked capitalism and argues for the need to consider human dignity. Speaker4 challenges Speaker2 to explain the consequences of an economically unsustainable capitalist system. The conversation *becomes heated*, with Speaker2 perceiving Speaker4’s questions *as confrontational*. Speaker2 defends their views and *provides examples to support their argument*. The *overall tone of the conversation remains argumentative, but civil*.

Figure 1: Traditional and dynamics summaries for the same conversation (transcript in Appendix F). Elements of conversational dynamics are colored in blue.

A holistic description of a conversation and its trajectory requires accounting for both of these communication levels. We complement prior work that has largely focused on summarizing the topical context of the discussion (Yang and Zhu, 2023), by introducing the task of generating summaries that instead capture the dynamics of the interaction between the participants. As shown in Figure 1, these cover aspects lost in a traditional summary.

051 Summaries of conversation dynamics (or SCDs
052 for short) provide a way for humans to quickly un-
053 derstand the trajectory of a discussion: what type of
054 interactions lead to its current state, and how these
055 are likely to develop? This type of understanding
056 can benefit various applications, including super-
057 vision of conversations in time-sensitive domains
058 (e.g., online community moderation and mental
059 health crisis counseling), providing context to users
060 (re)joining an online conversation, or identifying
061 and reviewing common problems in human-human
062 or human-AI conversations (Section 8).

063 However, generating SCDs that effectively cap-
064 ture conversation trajectories presents several new
065 challenges. While prior computational work intro-
066 duced models for separately capturing individual
067 aspects of conversation dynamics (Section 6), an
068 effective and concise summary must *select* those
069 that are most relevant for understanding the trajec-
070 tory of the conversation. Additionally, an infor-
071 mative summary must not simply identify these
072 aspects separately, but should also describe how
073 they evolve and interrelate throughout a conversa-
074 tion: for example, a conversation that transitions
075 from an aggressive tone to a calmer one has a com-
076 pletely different trajectory than one that proceeds
077 in reverse order. Thus, to provide an understanding
078 of the trajectory of a conversation, an SCD must
079 *synthesize* different aspects of its dynamics across
080 multiple utterances and participants.

081 As a first step, we devise a multi-step procedure
082 for human annotators to collaboratively write SCDs.
083 Importantly, this procedure is designed to address
084 the selection and synthesis challenges described
085 above. Building on this procedure, we develop a
086 large language model prompt for generating SCDs
087 and compare them with summaries generated by
088 other baselines, including traditional summaries.¹

089 Specifically, in this paper we evaluate the use-
090 fulness of SCDs for conversation trajectory un-
091 derstanding via an established task: forecasting
092 whether an ongoing conversation will eventually
093 derail into toxic behavior (Zhang et al., 2018a; Liu
094 et al., 2018). While prior attempts at this task
095 started directly from the transcript (Section 6), we
096 explore generating SCDs as an intermediate step.
097 This approach has the potential advantage of adding
098 interpretability to automated forecasting systems

¹We include the collection of human-written summaries and sample model outputs as supplemental material and will make them and our code publicly available upon acceptance to encourage further work on this problem.

and improving efficiency for humans (such as mod- 099
erators) trying to make such judgments (Schluger 100
et al., 2022). 101

Our findings reveal the potential of SCDs to help 102
both humans and automated systems understand a 103
conversation’s trajectory, motivating further work 104
on this new task. In the downstream task of fore- 105
casting the future derailment of a conversation, hu- 106
mans make predictions three times faster, and with 107
greater confidence, when reading the SCDs than 108
when reading the transcript. Furthermore, auto- 109
mated systems are more accurate when construct- 110
ing, and then predicting based on, SCDs compared 111
to state-of-the-art systems that base their forecast 112
directly on the transcript. Finally, by comparing 113
human-written and machine-generated summaries, 114
we reveal a quality gap that motivates further com- 115
putational work on this new task. 116

In summary, this work: 117

1. introduces the task of summarizing conversa- 118
tion dynamics, together with a collection of 119
human-written summaries; 120
2. proposes a downstream evaluation method 121
that allows for comparison between methods 122
for generating them; 123
3. shows the usefulness of dynamics summaries, 124
motivating further work on this new task. 125

2 Human-written Summaries 126

To start, we introduce a procedure for writing SCDs 127
and a collection of such summaries for an existing 128
dataset of online conversations. 129

Procedure for writing summaries. To construct 130
the first collection of SCDs, we iteratively designed 131
a writing procedure that addresses the selection and 132
synthesis challenges described in the introduction. 133
In early iterations in which we asked a single an- 134
notator to both read the transcript and write its 135
SCD, we observed that they consistently omitted 136
key information that they take for granted, perhaps 137
because some aspects of the dynamics are often 138
processed non-consciously (Tannen, 2005). To ad- 139
dress this issue, we devise a procedure that uses 140
interaction between two annotators to surface key 141
elements of the conversation dynamics that readers 142
who can not see the transcript would consider rele- 143
vant. Thus, we settle on a procedure that has two 144
parts—one in which an annotator works individ- 145
ually and one in which they interact with another 146

147 annotator—which we briefly outline here (and de-
148 tail in Appendix A).

149 For the individual work, Annotator A will draft
150 several summaries for a transcript in 4 steps:

- 151 1. skim over the transcript to have an overview
152 of the topic and of the role of each speaker;
- 153 2. read the transcript utterance by utterance
154 and write a comprehensive summary, includ-
155 ing opinions and arguments expressed within
156 most utterances, turning points, and elements
157 of conversation dynamics;
- 158 3. condense the summary by selecting key points
159 and dynamics and replacing specific opinions
160 and arguments with high-level descriptions;
- 161 4. write a brief summary for each of the main
162 speaker, focusing on (the changes in their)
163 tone and on their conversational strategies.

164 In the interactive part, Annotator B will write the
165 SCD, by interacting with Annotator A with a goal
166 of understanding the conversation trajectory. In this
167 process, Annotator B may read the summaries writ-
168 ten in the previous steps by Annotator A, but not the
169 transcript, and may make inquiries on details they
170 deem important to understand the trajectory, such
171 as ‘was this said neutrally, or is there something
172 about the tone that I should note?’ or ‘is the com-
173 ment overtly rude, or is it just passive-aggressive
174 or blunt?’, surfacing key aspects that were not ex-
175 plicitly mentioned in Annotator A’s summaries.

176 **Conversation transcripts data.** We apply this
177 procedure to summarize conversations from the
178 Conversations Gone Awry (CGA) dataset (Chang
179 and Danescu-Niculescu-Mizil, 2019),² a conversa-
180 tion corpus collected from the ‘Change My View’
181 subreddit, where people actively seek to have oth-
182 ers challenge their views on controversial topics.
183 This community has been studied extensively in
184 part because of their explicit norms against toxic
185 behavior, and corresponding labels inferred from
186 the moderators’ interventions.

187 In the CGA corpus, conversations are paired
188 such that every conversation that derailed—i.e.,
189 ended in a toxic comment removed by moderators—
190 is matched with another conversation on the same
191 topic that did not. For us, these labels provide an
192 opportunity to test the extent to which SCDs pro-
193 vide an intuition about the future trajectory of the

194 conversation (i.e., will it derail or not). To focus
195 on the *future* trajectory, we remove the last 3 utter-
196 ances from every conversation (in addition to the
197 toxic comment, if there was one). Since our inter-
198 est is in summarization, we only keep pairs where
199 both conversations are longer than 10 utterances.

200 **Collection of human-written summaries.** We
201 produce human-written summaries for 50 conver-
202 sations from the train split of CGA. The summary
203 writing process took roughly 240 annotator-hours.³
204 Summaries are on average 66 words long (annota-
205 tors are instructed to keep them under 80); for com-
206 parison, the transcripts are on average 984 words
207 long. An example summary is shown in Figure 1,
208 and a qualitative analysis is provided in Section 5.

209 **Informativeness check.** Before we proceed, we
210 check whether the summaries are actually informa-
211 tive. Given their highly abstractive nature, there
212 is a risk that they become so general as to not
213 distinguish between different conversations (e.g.,
214 ‘Speaker1 disagreed with Speaker2.’ would apply
215 to most of the conversations in the data). We devise
216 a procedure for systematically checking whether
217 this is the case, avoiding subjective interpretations
218 that are prone to apophenia.

219 We ask new annotators to read a transcript, and
220 then present them with a multiple-choice question.
221 Each choice corresponds to a summary segment
222 involving two speakers. One of the choices is
223 from the actual summary of the provided transcript,
224 while the other two are distractors: one from the
225 summary of the paired conversation (thus, on the
226 same topic, but with the opposite derailment label)
227 and the other from the summary of another conver-
228 sation with the same label as the transcript, but on
229 a different topic. This way, neither the topic nor the
230 label fully reveals the answer: to be identified cor-
231 rectly, the segment must contain information that
232 matches the transcript better than the distractors.

233 For example, for our introductory example, three
234 choices could be: “SpeakerX sarcastically criti-
235 cizes SpeakerY’s attitude which aggravates Speak-
236 erY more.” (actual segment), “SpeakerX poses
237 a rhetorical question, which SpeakerY contradicts
238 sarcastically, raising the tension and causing Speak-
239 erX to disagree rudely.” (same-pair distractor),
240 “SpeakerX first shares their opinion and later poses
241 rhetorical questions, and SpeakerY disagrees in a
242 matter-of-fact manner.” (same-label distractor).

³For each conversation transcript, the individual part takes about 2 hours and the interactive part takes about 20 minutes.

²Accessed via the ConvoKit library (Chang et al., 2020).

243 Though we designed this procedure to avoid
244 excessive workload when evaluating informativeness,⁴ each question still requires reading one tran-
245 script and carefully checking the segment choices
246 against it. Therefore, we limit our total number of
247 questions to 10, covering 30 conversations through
248 distractors. (Further details in Appendix B.)

249 Two annotators completed the task, one of which
250 answered 10 out of 10 questions correctly and the
251 other answered 8 of them correctly (noting low
252 confidence on the 2 answers they got wrong), sug-
253 gesting that our summaries indeed pass this basic
254 informativeness check.
255

256 3 Machine-generated Summaries

257 We now turn to explore several simple baselines for
258 generating SCDs, setting the stage for developing
259 more specialized methods in future work. The GPT-
260 family models have shown remarkable results in
261 various summarization benchmarks (Zhang et al.,
262 2023; Yang et al., 2023). Among them, ChatGPT
263 is particularly suitable for adapting to new tasks
264 like ours without demanding a sizable train set.
265 Thus, for the first group of baselines, we query Ope-
266 nAI’s ChatGPT (GPT-3.5-turbo-0613) API with
267 default parameters using different prompts, from
268 the most common prompt for traditional summa-
269 rization tasks to prompts inspired by the procedure
270 we developed for humans:⁵

271 **Traditional prompt.** After experimenting with
272 several prompts on a development set, we use a
273 concise prompt for our traditional summarization
274 baseline: ‘briefly summarize the following online
275 conversation in 80 words.’ Figure 1 includes a
276 traditional summary generated by this prompt.

277 **Zeroshot prompt.** We devise a prompt that explic-
278 itly integrates our goal of generating summaries
279 that can help people understand the conversation
280 trajectory. After experimenting with several word
281 choices for referring to trajectory, dynamics, and
282 specific dynamics elements, we settle on a con-

⁴An equivalent check could be implemented by provid-
ing one summary segment and three transcripts to pick
from. This method corresponds to the existing literature in
communication-based evaluations for natural language gen-
eration (Newman et al., 2019), and specifically the idea that
an informative summary should capture the salient informa-
tion that makes the source text stand out with respect to other
related texts (Zhang et al., 2018c). However, this equivalent
method would require substantially longer time due to the
lengths of the transcripts.

⁵We prompt the model to generate summaries of at most 80
words and set the max new token limit to 128 (corresponding
to approximately 96 words) as a hard limit .

283 cise prompt, ‘write a short summary capturing the
284 trajectory of the online conversation’ with addi-
285 tional constraints such as not including specific
286 arguments and capturing elements of tone and con-
287 versation strategies (Figure 3 in the Appendix).

288 **Procedural prompt.** We build on the insights we
289 gathered from developing the procedure for human
290 annotators (Section 2) to construct a more elaborate
291 prompt. This prompt (Figure 3 in the Appendix)
292 thus includes instructions adapted from those pro-
293 vided to the annotators, together with examples
294 that they found useful for understanding the in-
295 structions. Because we only include segments of
296 summary examples instead of complete transcript
297 and summary pairs, the procedural prompt can be
298 positioned in-between zeroshot and few-shot in-
299 context learning. Figure 1 shows the procedural
300 prompt summary for our introductory example.

301 We also experimented with few-shot in-context
302 learning on a small subset of the train set, but man-
303 ual inspection did not reveal an increase in quality.
304 Thus, due to significantly higher API costs, we did
305 not pursue this path. Appendix D includes more
306 discussion on our prompt engineering.

307 **Finetuning.** Finally, we experimented with fine-
308 tuning GPT-3.5-turbo as well as with smaller dia-
309 log summarization systems (BART-large and Di-
310 alogLED) (Lewis et al., 2020; Zhong et al., 2022)
311 using the 50 human-written summaries and the
312 corresponding transcripts. We provide details in
313 Appendix D.

314 4 Downstream Evaluation: 315 Forecasting Derailment

316 Popular metrics for summarization—e.g., ROUGE
317 (Lin, 2004), BERTScore (Zhang et al., 2020),
318 and QA-based metrics—are notoriously unreli-
319 able when evaluating LLM-generated summaries or
320 summaries of long documents (Goyal et al., 2023;
321 Koh et al., 2022). We thus follow recommenda-
322 tions of Deutsch et al. (2021) and perform a down-
323 stream evaluation, in which we quantify the extent
324 to which SCDs provide an understanding of the
325 conversation trajectory.

326 Specifically, we choose the task of forecasting
327 whether a conversation will eventually derail into
328 toxic behavior (Zhang et al., 2018a). Unlike previ-
329 ous work in which the prediction was made based
330 on a truncated transcript of the conversation (for a
331 comprehensive discussion of prior models see Sec-
332 tion 6), here we aim to make the prediction directly

on the SCD of that truncated transcript. In addition to providing means to evaluate and compare current and future models for generating SCDs, this derailment forecasting task is also important in itself, as it was shown to enable important practical applications: automated forecasts can be used to inform users during ongoing discussions (Chang et al., 2022) while human forecasts are made by moderators in their everyday workflow (Schluger et al., 2022) (see Section 8 for practical and ethical considerations of real-world deployment).

We first compare the usefulness of SCDs for automated forecasting systems. Then we devise an experiment to estimate their usefulness for human forecasts. Throughout, the forecasts are done on a balanced dataset of derailing and non-derailing conversations paired by topic, following the setup of the CGA dataset (Section 2); thus the overall topic of the discussion plays a minimal role and the random baseline is 50%. To leave room for future work and avoid polluting the available data, we leave the original CGA test set untouched. Using truncated transcripts from the original train split, we construct a new train set (234 conversations), a new dev set (100), and a new test set (100); the new test set includes the 50 conversations for which we also have human-written summaries (Section 2).

4.1 Useful for automated forecasts?

We train classifiers to predict if a conversation will eventually derail based on the various types of summaries of the truncated transcripts. We adopt GPT-3.5-turbo to develop few-shot classifiers for each summary type, using examples from outside the test split. We also train BART (Lewis et al., 2020) and longformer (Beltagy et al., 2020) as finetuned classifiers. To provide more robust estimates, for each summarization method we generate 4 different summaries for each conversation, and average the classifiers’ performance on them (details in Appendix D.3). We find GPT-3.5 few-shot classifier gives the best performance across all types of summaries, so we use it for our main analysis here. While the performance of other classifiers is substantially lower (Appendix D.3), the comparisons discussed below still hold.

Comparison of summaries. As shown in Table 1, the classifier based on the procedural prompt achieves the best accuracy, significantly outperforming the other types of summaries ($p < 0.05$, throughout we use the Wilcoxon signed-rank test significance testing). In particular, the information

Based on...	Accuracy
transcripts (CRAFT classifier)	56.2
transcripts (GPT-16k classifier)	60.0
traditional prompt summaries	58.3 (5.85)
zeroshot prompt summaries	58.8 (6.24)
procedural prompt summaries	67.3* (2.63)

Table 1: Derailment forecasting results for systems based on truncated transcripts and on different types of machine-generated summaries. For summary-based systems, we report standard deviation across 4 summary-generation trials, and indicate with * the highest performance ($p < 0.05$, Wilcoxon signed-rank test). Results for the GPT-3.5-turbo few-shot classifiers are shown unless otherwise noted.

conveyed by the SCDs generated with the procedural prompt appears to be more useful for the automatic derailment forecaster than that included in traditional summaries. Other metrics (Macro-F1, precision, recall) support the same conclusion (Appendix F.5).

The finetuned models—finetuned on the 50 human-written examples and evaluated on the remaining of the test set—perform worse than the procedural prompt on the same set (Appendix D). This could be due to the relatively small collection of human-written summaries, as well as the generic fine-tuning methodology.

Summary vs transcript. For reference, we also include two baselines operating directly on the truncated transcripts. The first baseline, CRAFT, was introduced before the advent of the LLM era and remained the state-of-the-art system for this task (Chang and Danescu-Niculescu-Mizil, 2019).⁶ The second baseline is a few-shot GPT-3.5-turbo-16k classifier, which can take up to 16k tokens to cope with the greater input lengths of the transcripts.⁷

As shown by Table 1, predictions based on procedural prompt outperform those based directly on the transcripts. This suggests that SCDs are effective in distilling from the transcripts information that is useful for the forecasting task. Perhaps more importantly, the feasibility of this ‘summarize-then-

⁶For a fair comparison, we modify the ConvoKit implementation of CRAFT (Chang et al., 2020) to trigger forecasts exactly 3 utterances before the end of the conversation. This setup is harder than the original setup in which the system could make predictions all the way up to right before the attack or the end of the conversation, potentially having access to more explicit signals of upcoming toxic behavior.

⁷Both baseline systems might have an advantage in that they might have accessed the full untruncated transcripts during pre-training.

forecast’ approach points out a promising future direction for improving the interpretability of the user-facing forecasting systems, where the summary could be presented as an easily digestible rationale for the prediction. In fact, users of such systems have identified the lack of explanations as one of their most important drawbacks (Chang et al., 2022).

4.2 Useful for human forecasts?

We now switch to the other main motivation: can SCDs help *humans* quickly grasp the trajectory of a conversation? To answer this question we devise an experiment in which subjects are asked to guess whether a conversation will eventually derail based either on a transcript or its SCD. We compare both their accuracy and efficiency, in terms of the time they spend to make their guess, as well as their confidence in their guess.

To better focus our resources, we use a subset of 20 paired conversations out of those for which we created human summaries. In addition to the transcripts and the human-written summaries, we also consider the corresponding procedural prompt summaries (since those were shown to fare best in the automatic prediction task).

We recruit 20 university students fluent in English as participants. A subset of participants make their guesses based on the transcripts only, while another subset will make guesses based on summaries only. Each participant in the latter subset will see a mix of human-written and machine-generated summaries (without being aware that these are produced differently) such that any observed differences between them cannot be attributed to participant idiosyncrasies. In addition to providing a guess of whether the conversation will derail or not, each participant is asked to rate their confidence in their guess (on a scale from 1 to 5). We also record the time it took for the participants to make their guess (starting from the time they see the transcript or summary until the time they select their guess), and instruct them to work on each question without pausing. The specific instructions and details about how participants are grouped are in Appendix C.

Unlike in the automatic evaluation in Section 4, we adopt a zero-shot prediction setting, in which humans do not have labeled examples of summaries (or transcripts) to assist their guessing. This way, we can better test if the summaries are immediately intuitive to humans rather than testing the participants’ ability to learn patterns that might not

be visible to untrained individuals. This means, however, that the accuracies of the human participants are not directly comparable with those of the automated system.

Summaries vs transcripts. As shown in Table 2, participants can make guesses 3-4 times faster based on SCDs while maintaining similar accuracy. This improvement in efficiency is critical for applications such as proactive online moderation, as earlier work has found that moderators are faced with “too many [potentially at-risk conversations] to proactively monitor (Schluger et al., 2022).”

Human vs generated summary. Participants are significantly more confident when making predictions based on human-written summaries than on machine-generated summaries (and even on the transcript).⁸ This gap is important for applications where summaries are used for decision-making (e.g., moderation) and motivates future work on improving summarization models. Another noticeable difference is that machine-generated summaries provide a better understanding of the topical content of the discussion, perhaps to the detriment of better coverage of aspects of conversation dynamics. In Section 5 we further explore this trade-off via a qualitative analysis of the summaries.⁹

Based on...	Acc	Conf	Topic	Time
transcripts	60	3.5	-	132
gen. summ.	59	3.6	3.9	45*
human summ.	62	4.0*†	3.4†	31*

Table 2: Results on the human forecasting experiment. “gen. summ.” refers to the summaries generated using the procedural prompt. Time is measured in seconds. * indicates a significant difference when compared with transcripts ($p < 0.05$, Wilcoxon signed rank test), † indicates a significant difference when comparing human-written with automated summaries ($p < 0.05$).

5 Qualitative Analysis

To complement our quantitative evaluation and understand what might drive the differences between human and machine-generated summaries, we now

⁸This difference continues to hold when only considering correct guesses. Also, reassuringly, confidence in correct guesses is higher than in incorrect ones throughout.

⁹We also experimented with directly asking participants to report their understanding of the trajectory of the conversation, on a scale from 1 to 5. There was no significant difference between human and machine-written summaries (4 and 3.9 respectively), perhaps due to the difficulty of briefly explaining what a trajectory is and how it differs from the guess, a confusion that surfaced during debriefing.

493	turn to the actual content of the SCDs. Through	545
494	a close reading of the 20 human-written and 20	546
495	machine-generated summaries used in the experi-	547
496	ment described above, we identify, annotate, and	548
497	compare several aspects that were shown to provide	549
498	clues about the conversation trajectories.	550
499	Tone. Whether ‘polite,’ ‘rude,’ ‘aggressive,’ ‘con-	551
500	descending,’ or ‘sarcastic’ (Tannen, 2005; Brown	552
501	and Levinson, 1987), the tone employed by the	553
502	participants is a prominent feature of the SCDs.	554
503	Tone can be explicitly stated, as in ‘Speaker1 dis-	555
504	agrees [...] in a somewhat passive-aggressive	556
505	tone. ’ Other times, especially in human-written	557
506	summaries, it is expressed as modifying a speech	558
507	act, as in ‘contradicts sarcastically ,’ ‘disagrees	559
508	rudely ,’ and ‘ admantly defends.’ Overall, tone is	560
509	indicated less frequently in the machine-generated	561
510	summaries (75% of them mention tone at least	562
511	once) than in the human-written summaries of the	563
512	same conversations (all mention tone at least once),	564
513	suggesting a potential path for improvement.	565
514	Changes in tone. Tone can evolve throughout a	566
515	conversation, and changes in tone can provide an	567
516	intuition about its trajectory (Niculae et al., 2015).	568
517	When participants use an ‘ increasingly passive	569
518	aggressive tone,’ or when the ‘tension rises ’ the	570
519	conversation seems more likely to be getting out of	571
520	hand than when a ‘slight tension [...] is maintained	572
521	but doesn’t escalate ’ or when the ‘tone remains	573
522	argumentative but civil’. The latter quote is an ex-	574
523	ample of an overall assessment of tone dynamics	575
524	that both humans and (more commonly) automated	576
525	systems sometimes include at the end of the sum-	577
526	mary, even though neither is explicitly instructed	578
527	to do so. Overall, 75% of the human summaries	579
528	feature phrases explicitly mentioning changes in	580
529	tone whereas only 50% of the machine-generated	581
530	counterparts do so.	
531	Patterns of interaction. Beyond the content of the	582
532	messages, the structural properties of the interac-	
533	tions were shown to be indicative of future trajec-	
534	tories (Backstrom et al., 2013; Zhang et al., 2018b).	
535	Two participants can have a ‘ brief exchange ’ or	
536	an extended ‘ back-and-forth ,’ which can be inter-	
537	rupted when another participant ‘ jumps in ’. While	
538	explicit mentions of such patterns are relatively	
539	rare (found in 45% of the human summaries and	
540	31% of the machine-generated summaries), they	
541	can often be inferred by following the sequence of	
542	speakers mentioned in the summaries.	
543	Conversation strategies. Interlocutors employ	
544	strategies that can put the conversation on vari-	
	ous trajectories. For example, ‘pos[ing] a rhetor-	545
	ical question’ or ‘questioning each other’s logic’,	546
	can often lead to personal attacks (Habernal et al.,	547
	2018), whereas expressing uncertainty about one’s	548
	own view (e.g., via hedging), would soften an im-	549
	pending disagreement and prevent the escalation of	550
	tension (Zhang et al., 2018a). ‘Supporting [a] point	551
	with evidence’, ‘justifying objective claims with	552
	personal experiences’, ‘draw[ing] a comparison’	553
	or ‘question[ing] the importance of specific details’	554
	are classic persuasion strategies (Zeng et al., 2020;	555
	Li et al., 2020a). A list of strategies considered in	556
	this analysis is included in Appendix E. Overall,	557
	we find that mentions of strategies are evenly dis-	558
	tributed among human-written (80%) and machine-	559
	generated summaries (85%).	560
	Topical context. Finally, these dynamics can only	561
	exist in the context of the content being discussed.	562
	Though not the primary focus of SCDs, a small	563
	amount of topical context is needed to provide a	564
	scaffolding for the phenomena discussed above.	565
	Both human and machine-generated summaries	566
	generally start with a sentence about the general	567
	topic of the discussion. Beyond that, machine-	568
	generated summaries include substantially more	569
	topical context to the detriment of actual aspects	570
	of conversation dynamics, despite the explicit in-	571
	struction and in-context-learning examples against	572
	this behavior. This echoes the subjective ratings	573
	of the participants in the human forecasting experi-	574
	ment (Table 2). This phenomenon suggests that in-	575
	context learning is not sufficient to ‘untrain’ LLMs	576
	from the traditional summary examples seen in pre-	577
	training. This motivates developing models that	578
	are specifically designed to select and synthesize	579
	aspects of conversation dynamics, perhaps inspired	580
	by the interactive human-writing procedure.	581
	6 Further Related Work	582
	Our work falls at the intersection of three broad	583
	areas of NLP: studies of conversation dynamics,	584
	summarization, and conversation forecasting.	585
	Conversation dynamics. We are primarily in-	586
	spired by extensive computational work on mod-	587
	eling various aspects of conversation dynamics.	588
	Some studies have focused on identifying specific	589
	aspects, such as such as politeness (Burke and	590
	Kraut, 2008; Danescu-Niculescu-Mizil et al., 2013;	591
	Li et al., 2020b), formality (Pavlick and Tetreault,	592
	2016; Krishnan and Eisenstein, 2015), passive-	593
	aggressiveness (Chhaya et al., 2018), condescen-	594

sion (Wang and Potts, 2019) or sarcasm (Oraby et al., 2017). Others have studied changes along these dimensions during the discussion (Wang and Cardie, 2014; Niculae et al., 2015; Niculae and Danescu-Niculescu-Mizil, 2016). A separate but related thrust focused on persuasive strategies interlocutors employ in a conversation, mostly in the context of debates (see Lawrence and Reed (2020) for a survey). Unlike these studies, the goal of SCDs is not to exhaustively identify occurrences of either one of these phenomena, but to convey how such key aspects combine towards an understanding of the conversation’s trajectory.

Dialogue summarization. The vast majority of dialogue summarization systems focus on the content of the utterances, rather than on the more subtle non-topical dynamics. Early approaches to dialogue summarization focused on using external tools to explicitly model dialogue structures, such as topic segmentation and conversation stages (Li et al., 2019; Chen and Yang, 2020), dialogue acts (Goo and Chen, 2018), along with discourse dependency and speaker-action relations (Chen and Yang, 2021), which are processed into features that can help language models. Later, pretraining on dialogue corpora also attracted increasing research interest and achieved state-of-the-art results on many datasets (Zhong et al., 2022). Most recently, extensively pretrained, instruction-tuned, LLMs, such as the GPT-family models, have achieved superior results on various summarization leaderboards (Goyal et al., 2023; Zhang et al., 2023; Yang et al., 2023). In dialogue summarization, these instruction-tuned LLMs possess strong in-context-learning capabilities (Wu et al., 2023), making them strong candidates for solving new summarization tasks that have limited training data.

Conversation forecasting. We motivate and evaluate dynamics summaries with applications requiring an understanding of a conversation’s trajectory. Beyond forecasting derailment (Zhang et al., 2018a; Liu et al., 2018; Chang and Danescu-Niculescu-Mizil, 2019), other tasks include forecasting thread growth (Backstrom et al., 2013), prosocial outcomes (Bao et al., 2021), editorial decisions (Mayfield and Black, 2019), the outcomes of negotiations (Chawla et al., 2020) or team resilience (Whiting et al., 2019). It would be interesting to consider the extent to which SCDs can aid with these other forecasting tasks, and what modifications might be needed to summaries specifically dedicated to these tasks.

To the best of our knowledge, all conversational forecasting work operate directly on conversation transcripts. The early work by Chang and Danescu-Niculescu-Mizil (2019) adopts a recurrent network and applies unsupervised training to learn a representation of conversation dynamics. More recently, Kementchedjheva and Sogaard (2021) explores pretraining and various training paradigms for this task, Altarawneh et al. (2023) applies a graph convolutional network, and Yuan and Singh (2023) uses a hierarchical transformer-based framework to combine utterance-level and conversation-level information. However, since it aims to guess the future, this task remains challenging and none of the mentioned approaches achieved an accuracy beyond 65% on the CGA balanced dataset.

Unlike detecting the toxic language after the fact (Wulczyn et al., 2017; Breiffeller et al., 2019), the signs of future derailment are subtle and require a more thorough understanding of the conversation trajectory. Our results suggest that SCDs can provide this information concisely and effectively, suggesting a new summarize-then-forecast approach to conversational forecasting tasks. This inspires future work that integrates SCDs in real-time forecasting systems, which would require tackling shorter conversations where summaries might not be appropriate, as well as the ‘unknown horizon’ problem: not knowing when to trigger the prediction (Chang and Danescu-Niculescu-Mizil, 2019).

7 Conclusions

In this work, we introduce the task of summarizing the dynamics of interaction between participants in a text-based conversation. By introducing human and automated procedures for writing such summaries, we show that they can capture information that is mostly missing from traditional summaries, such as the tone in which the participants write and how it changes throughout a conversation. Summaries of these dynamics are useful to both humans and automated systems for understanding the overall trajectory of the conversation, as shown through the downstream evaluation task of forecasting whether a conversation will eventually derail or not. Humans can make similarly accurate forecasts three to four times faster by starting from SCDs than by reading the transcripts. Automated systems can surpass the performance of state-of-the-art forecasting systems when generating SCDs as an intermediate step for forecasting.

8 Limitations

This work, however, only takes the first steps towards solving and evaluating this task automatically. In fact, we show that there is a substantial gap remaining between human-written summaries and machine-generated ones. Since in this work we focus on defining the task and demonstrating its feasibility, we only employ simple prompting and standard fine-tuning procedures. This sets the stage for the future development of more specialized models and training regimes. These models could be more tightly integrated with the downstream task, learning to select aspects of the dynamics that are most relevant as well as to determine the right level of abstraction.

To continue improving on dynamics generation models, more diverse automated evaluation methods are required. Given the highly abstractive nature of the task, traditional metrics based on token overlap or semantic similarity are not immediately applicable (Goyal et al., 2023). Our informativeness check provides an avenue for evaluation that could potentially be scaled up through automation. Furthermore, considering other downstream applications, such as forecasting prosocial outcomes (Bao et al., 2021) or how likely it is for participants to change their mind (Tan et al., 2016; Hovy and Yang, 2021), could further help evaluate the usefulness of dynamics summaries.

While the current work is restricted to summaries of text-based conversations, important dynamics can be encoded in vocal features (e.g., intonation, or pitch) or gestures (laughter, body positioning). A multimodal approach could enable applications that go beyond text-based conversations, and could even aid conversational analysis researchers explore the intricate ways in which conversations develop (Sidnell, 2011).

Additionally, while we tested how useful summaries are for humans in a small-scale control setting, further work could test this more comprehensively through user studies, for example by integrating these summaries into conversational assistance tools (Chang et al., 2022) or moderation assistance tools (Schluger et al., 2022). From a technical perspective, a real-time deployment would require iteratively generating summaries in real-time, as the conversation progresses, rather than at a set moment in the conversation as we do in this work for the sake of scalability.

Ethical concerns surrounding fairness and bias should necessarily take center stage in any deployment of summarization systems, especially since SCDs may include mentions of emotions and affect of the people involved in the conversation (Zhou and Tan, 2023). Any broad usage scenario should undergo rigorous scrutiny of potential for unintended consequences (Weidinger et al., 2022). For example, SCDs and automated forecasts relying on them should not be used to make automated censoring or moderation decisions, to avoid propagating biases embedded in the underlying large language models. If future developments will result in summaries that are reliable enough to inform human decisions (e.g., helping moderators decide whether to closely monitor an ongoing conversation), the users should be informed about systematic mistakes the summary is likely to make in that respective setting.

References

- Enas Altarawneh, Ammeta Agrawal, Michael Jenkin, and Manos Papagelis. 2023. *Conversation Derailment Forecasting with Graph Convolutional Networks*. ArXiv:2306.12982 [cs].
- 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, WSDM*.
- Jiajun Bao, Junjie Wu, Yiming Zhang, Eshwar Chandrasekharan, and David Jurgens. 2021. *Conversations Gone Alright: Quantifying and Predicting Prosocial Outcomes in Online Conversations*. In *Proceedings of WWW*, pages 1134–1145.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. *Longformer: The Long-Document Transformer*. ArXiv:2004.05150 [cs].
- Luke Breitfeller, Emily Ahn, David Jurgens, and Yulia Tsvetkov. 2019. *Finding Microaggressions in the Wild: A Case for Locating Elusive Phenomena in Social Media Posts*. 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 1664–1674, Hong Kong, China. Association for Computational Linguistics.
- Penelope Brown. 2015. *Politeness and Language*. In *International Encyclopedia of the Social & Behavioral Sciences*, pages 326–330. Elsevier.
- Penelope Brown and Stephen C. Levinson. 1987. *Politeness: Some Universals in Language Usage*. Cambridge University Press.

800	Moira Burke and Robert Kraut. 2008. Mind Your Ps and Qs: The Impact of Politeness and Rudeness in Online Communities. In <i>Proceedings of CSCW</i> .	<i>and Emotions in Social Media</i> , pages 76–86, New Orleans, Louisiana, USA. Association for Computational Linguistics.	858
801			859
802			860
803	Jean Carletta, Simone Ashby, Sebastien Bourban, Mike Flynn, Mael Guillemot, Thomas Hain, Jaroslav Kadlec, Vasilis Karaiskos, Wessel Kraaij, Melissa Kronenthal, Guillaume Lathoud, Mike Lincoln, Agnes Lisowska, Iain McCowan, Wilfried Post, Dennis Reidsma, and Pierre Wellner. 2006. <i>The AMI Meeting Corpus: A Pre-announcement</i> . In <i>Machine Learning for Multimodal Interaction</i> , Lecture Notes in Computer Science, pages 28–39, Berlin, Heidelberg. Springer.	Cristian Danescu-Niculescu-Mizil, Moritz Sudhof, Dan Jurafsky, Jure Leskovec, and Christopher Potts. 2013. <i>A Computational Approach to Politeness with Application to Social Factors</i> . In <i>Proceedings of ACL</i> .	861
804			862
805			863
806			864
807			
808		Daniel Deutsch, Rotem Dror, and Dan Roth. 2021. <i>A Statistical Analysis of Summarization Evaluation Metrics Using Resampling Methods</i> . <i>Transactions of the Association for Computational Linguistics</i> , 9:1132–1146. Place: Cambridge, MA Publisher: MIT Press.	865
809			866
810			867
811			868
812			869
813	Jonathan P. Chang, Caleb Chiam, Liye Fu, Andrew Wang, Justine Zhang, and Cristian Danescu-Niculescu-Mizil. 2020. <i>ConvoKit: A Toolkit for the Analysis of Conversations</i> . In <i>Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 57–60, 1st virtual meeting. Association for Computational Linguistics.		870
814			
815		Chih-Wen Goo and Yun-Nung Chen. 2018. <i>Abstractive Dialogue Summarization with Sentence-Gated Modeling Optimized by Dialogue Acts</i> . In <i>2018 IEEE Spoken Language Technology Workshop (SLT)</i> , pages 735–742, Athens, Greece. IEEE.	871
816			872
817			873
818			874
819			875
820	Jonathan P. Chang and Cristian Danescu-Niculescu-Mizil. 2019. <i>Trouble on the Horizon: Forecasting the Derailment of Online Conversations as they Develop</i> . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 4743–4754, Hong Kong, China. Association for Computational Linguistics.	Tanya Goyal, Junyi Jessy Li, and Greg Durrett. 2023. <i>News Summarization and Evaluation in the Era of GPT-3</i> . ArXiv:2209.12356 [cs].	876
821			877
822			878
823		Ivan Habernal, Henning Wachsmuth, Iryna Gurevych, and Benno Stein. 2018. <i>Before Name-Calling: Dynamics and Triggers of Ad Hominem Fallacies in Web Argumentation</i> . In <i>Proceedings of NAACL</i> , pages 386–396.	879
824			880
825			881
826			882
827			883
828			
829	Jonathan P. Chang, Charlotte Schluger, and Cristian Danescu-Niculescu-Mizil. 2022. <i>Thread With Caution: Proactively Helping Users Assess and Deescalate Tension in Their Online Discussions</i> . <i>Proceedings of the ACM on Human-Computer Interaction</i> , 6(CSCW2):545:1–545:37.	Chung-hye Han. 2002. <i>Interpreting interrogatives as rhetorical questions</i> . <i>Lingua</i> , 112(3):201–229.	884
830			885
831		Dirk Hovy and Diyi Yang. 2021. <i>The Importance of Modeling Social Factors of Language: Theory and Practice</i> . In <i>Proceedings of NAACL</i> , page 15.	886
832			887
833			888
834			
835	Kushal Chawla, Gale Lucas, Jonathan Gratch, and Jonathan May. 2020. <i>BERT in Negotiations: Early Prediction of Buyer-Seller Negotiation Outcomes</i> . <i>arXiv:2004.02363 [cs]</i> . ArXiv: 2004.02363.	T. Huckin. 2002. <i>Critical Discourse Analysis and the Discourse of Condescension</i> .	889
836			890
837		A. Janin, D. Baron, J. Edwards, D. Ellis, D. Gelbart, N. Morgan, B. Peskin, T. Pfau, E. Shriberg, A. Stolcke, and C. Wooters. 2003. <i>The ICSI Meeting Corpus</i> . In <i>2003 IEEE International Conference on Acoustics, Speech, and Signal Processing, 2003. Proceedings. (ICASSP '03)</i> , volume 1, pages I–I. ISSN: 1520-6149.	891
838			892
839	Jiaao Chen and Diyi Yang. 2020. <i>Multi-View Sequence-to-Sequence Models with Conversational Structure for Abstractive Dialogue Summarization</i> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 4106–4118, Online. Association for Computational Linguistics.		893
840			894
841			895
842			896
843			897
844			
845		Julia Jorgensen. 1996. <i>The functions of sarcastic irony in speech</i> . <i>Journal of Pragmatics</i> , 26(5):613–634.	898
846	Jiaao Chen and Diyi Yang. 2021. <i>Structure-Aware Abstractive Conversation Summarization via Discourse and Action Graphs</i> . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1380–1391, Online. Association for Computational Linguistics.		899
847			
848		Yova Kementchedjheva and Anders Søgaard. 2021. <i>Dynamic Forecasting of Conversation Derailment</i> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 7915–7919, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	900
849			901
850			902
851			903
852			904
853	Niyati Chhaya, Kushal Chawla, Tanya Goyal, Projjal Chanda, and Jaya Singh. 2018. <i>Frustrated, Polite, or Formal: Quantifying Feelings and Tone in Email</i> . In <i>Proceedings of the Second Workshop on Computational Modeling of People’s Opinions, Personality, and Emotions in Social Media</i> , pages 76–86, New Orleans, Louisiana, USA. Association for Computational Linguistics.		905
854			
855		Huan Yee Koh, Jiaxin Ju, Ming Liu, and Shirui Pan. 2022. <i>An Empirical Survey on Long Document Summarization: Datasets, Models, and Metrics</i> . <i>ACM Computing Surveys</i> , 55(8):154:1–154:35.	906
856			907
857			908
			909

910	Vinodh Krishnan and Jacob Eisenstein. 2015. “You’re Mr. Lebowksi, I’m the Dude”: Inducing Address Term Formality in Signed Social Networks. In <i>Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1616–1626, Denver, Colorado. Association for Computational Linguistics.	967
911		968
912		969
913		
914		970
915		971
916		972
917		973
918	Robin T. Lakoff. 1973. <i>The Logic of Politeness: Mind-ing Your P’s and Q’s</i> . Chicago Linguistic Society. Google-Books-ID: DfWfNAAACAAJ.	974
919		975
920		976
921	John Lawrence and Chris Reed. 2020. Argument Mining: A Survey . <i>Computational Linguistics</i> , 45(4):765–818.	977
922		978
923		
924	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880, Online. Association for Computational Linguistics.	979
925		980
926		981
927		
928		982
929		983
930		984
931		985
932		986
933	Jialu Li, Esin Durmus, and Claire Cardie. 2020a. Exploring the Role of Argument Structure in Online Debate Persuasion . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 8905–8912, Online. Association for Computational Linguistics.	987
934		988
935		989
936		990
937		991
938		992
939	Manling Li, Lingyu Zhang, Heng Ji, and Richard J. Radke. 2019. Keep Meeting Summaries on Topic: Abstractive Multi-Modal Meeting Summarization . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2190–2196, Florence, Italy. Association for Computational Linguistics.	993
940		994
941		
942		995
943		996
944		997
945		998
946	Mingyang Li, Louis Hickman, Louis Tay, Lyle Ungar, and Sharath Chandra Guntuku. 2020b. Studying Politeness across Cultures Using English Twitter and Mandarin Weibo . <i>arXiv:2008.02449 [cs]</i> . ArXiv: 2008.02449.	999
947		1000
948		
949		995
950		996
951	Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.	997
952		998
953		999
954		1000
955	Ping Liu, Joshua Guberman, Libby Hemphill, and Aron Culotta. 2018. Forecasting the Presence and Intensity of Hostility on Instagram Using Linguistic and Social Features . In <i>Proceedings of ICWSM</i> .	1001
956		1002
957		1003
958		1004
959	Elijah Mayfield and Alan W. Black. 2019. Analyzing Wikipedia Deletion Debates with a Group Decision-Making Forecast Model . <i>Proceedings of the ACM on Human-Computer Interaction</i> , 3(CSCW):1–26.	1005
960		1006
961		1007
962		1008
963	Benjamin Newman, Reuben Cohn-Gordon, and Christopher Potts. 2019. Communication-based Evaluation for Natural Language Generation . ArXiv:1909.07290 [cs].	1009
964		1010
965		1011
966		
	Vlad Niculae and Cristian Danescu-Niculescu-Mizil. 2016. Conversational Markers of Constructive Discussions . In <i>Proceedings of NAACL</i> .	970
		971
		972
		973
	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 <i>Proceedings of ACL</i> .	974
		975
		976
		977
		978
	Shereen Oraby, Vrindavan Harrison, Amita Misra, Ellen Riloff, and Marilyn Walker. 2017. Are you serious?: Rhetorical Questions and Sarcasm in Social Media Dialog . pages 310–319, Saarbrücken, Germany. Association for Computational Linguistics.	979
		980
		981
	Ellie Pavlick and Joel Tetreault. 2016. An Empirical Analysis of Formality in Online Communication . <i>TACL</i> .	982
		983
		984
		985
		986
	Harvey Sacks, Emanuel A. Schegloff, and Gail Jefferson. 1974. A Simplest Systematics for the Organization of Turn-Taking for Conversation . <i>Language</i> , 50(4):696–735. Publisher: Linguistic Society of America.	987
		988
		989
		990
		991
		992
	Charlotte Schluger, Jonathan P. Chang, Cristian Danescu-Niculescu-Mizil, and Karen Levy. 2022. Proactive Moderation of Online Discussions: Existing Practices and the Potential for Algorithmic Support . <i>Proceedings of the ACM on Human-Computer Interaction</i> , 6(CSCW2):370:1–370:27.	993
		994
	Jack Sidnell. 2011. <i>Conversation Analysis: An Introduction</i> . John Wiley & Sons.	995
		996
		997
		998
		999
		1000
	Michael Silverstein. 1984. On the pragmatic ‘poetry’ of prose: Parallelism, repetition, and cohesive structure in the time course of dyadic conversation. <i>Meaning, form, and use in context: Linguistic applications</i> , pages 181–99. Publisher: Georgetown University Press Washington, DC.	1001
		1002
		1003
		1004
		1005
	Chenhao Tan, Vlad Niculae, Cristian Danescu-Niculescu, and Lillian Lee. 2016. Winning Arguments: Interaction Dynamics and Persuasion Strategies in Good-faith Online Discussions . In <i>Proceedings of WWW</i> .	1006
		1007
		1008
	Deborah Tannen. 2005. <i>Conversational style : analyzing talk among friends</i> . Oxford University Press, New York.	1009
		1010
		1011
	Lu Wang and Claire Cardie. 2014. A Piece of My Mind: A Sentiment Analysis Approach for Online Dispute Detection . pages 693–699.	1012
		1013
		1014
		1015
	Zijian Wang and Christopher Potts. 2019. TalkDown: A Corpus for Condescension Detection in Context . <i>arXiv:1909.11272 [cs]</i> . ArXiv: 1909.11272 version: 1.	1016
		1017
		1018
		1019
		1020
	Laura Weidinger, Jonathan Uesato, Maribeth Rauh, Conor Griffin, Po-Sen Huang, John Mellor, Amelia Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh, Courtney Biles, Sasha Brown, Zac Kenton, Will Hawkins, Tom Stepleton, Abeba Birhane, Lisa Anne	

1021	Hendricks, Laura Rimell, William Isaac, Julia Haas,	Yuhao Zhang, Daisy Yi Ding, Tianpei Qian, Christo-	1076
1022	Sean Legassick, Geoffrey Irving, and Iason Gabriel.	pher D. Manning, and Curtis P. Langlotz. 2018c.	1077
1023	2022. Taxonomy of Risks posed by Language Mod-	Learning to Summarize Radiology Findings . In <i>Pro-</i>	1078
1024	els . In <i>Proceedings of the 2022 ACM Conference on</i>	<i>ceedings of the Ninth International Workshop on</i>	1079
1025	<i>Fairness, Accountability, and Transparency</i> , FAccT	<i>Health Text Mining and Information Analysis</i> , pages	1080
1026	'22, pages 214–229, New York, NY, USA. Associa-	204–213, Brussels, Belgium. Association for Com-	1081
1027	tion for Computing Machinery.	putational Linguistics.	1082
1028	Mark E. Whiting, Allie Blaising, Chloe Barreau, Laura	Ming Zhong, Yang Liu, Yichong Xu, Chenguang Zhu,	1083
1029	Fiuza, Nik Marda, Melissa Valentine, and Michael S.	and Michael Zeng. 2022. DialogLM: Pre-trained	1084
1030	Bernstein. 2019. Did It Have To End This Way?:	Model for Long Dialogue Understanding and Sum-	1085
1031	Understanding The Consistency of Team Fracture .	marization . ArXiv:2109.02492 [cs].	1086
1032	In <i>Proceedings of CHI</i> , volume 3, pages 1–23.		
1033	Tianyu Wu, Shizhu He, Jingping Liu, Siqi Sun, Kang	Karen Zhou and Chenhao Tan. 2023. Characterizing Po-	1087
1034	Liu, Qing-Long Han, and Yang Tang. 2023. A Brief	litical Bias in Automatic Summaries: A Case Study	1088
1035	Overview of ChatGPT: The History, Status Quo and	of Trump and Biden . ArXiv:2305.02321 [cs].	1089
1036	Potential Future Development . <i>IEEE/CAA Journal</i>		
1037	of Automatica Sinica , 10(5):1122–1136.		
1038	Ellery Wulczyn, Nithum Thain, and Lucas Dixon. 2017.		
1039	Ex Machina: Personal Attacks Seen at Scale . In		
1040	<i>Proceedings of WWW</i> .		
1041	Diyi Yang and Chenguang Zhu. 2023. Summarization		
1042	of Dialogues and Conversations At Scale . In <i>Pro-</i>		
1043	<i>ceedings of the 17th Conference of the European</i>		
1044	<i>Chapter of the Association for Computational Lin-</i>		
1045	<i>guistics: Tutorial Abstracts</i> , pages 13–18, Dubrovnik,		
1046	Croatia. Association for Computational Linguistics.		
1047	Xianjun Yang, Yan Li, Xinlu Zhang, Haifeng Chen,		
1048	and Wei Cheng. 2023. Exploring the Limits of Chat-		
1049	GPT for Query or Aspect-based Text Summarization .		
1050	ArXiv:2302.08081 [cs].		
1051	Jiaqing Yuan and Munindar P. Singh. 2023.		
1052	Conversation Modeling to Predict Derailment .		
1053	ArXiv:2303.11184 [cs].		
1054	Jichuan Zeng, Jing Li, Yulan He, Cuiyun Gao,		
1055	Michael R. Lyu, and Irwin King. 2020. What		
1056	Changed Your Mind: The Roles of Dynamic		
1057	Topics and Discourse in Argumentation Process .		
1058	<i>arXiv:2002.03536 [cs]</i> . ArXiv: 2002.03536.		
1059	Justine Zhang, Jonathan P. Chang, Cristian Danescu-		
1060	Niculescu-Mizil, Lucas Dixon, Nithum Thain,		
1061	Yiqing Hua, and Dario Taraborelli. 2018a. Conversa-		
1062	tions Gone Awry: Detecting Early Signs of Conversa-		
1063	tional Failure . In <i>Proceedings of ACL</i> .		
1064	Justine Zhang, Cristian Danescu-Niculescu-Mizil,		
1065	Christina Sauper, and Sean J. Taylor. 2018b. Charac-		
1066	terizing Online Public Discussions Through Patterns		
1067	of Participant Interactions . In <i>Proceedings of CSCW</i> .		
1068	Tianyi Zhang, Varsha Kishore, Felix Wu, Kil-		
1069	ian Q. Weinberger, and Yoav Artzi. 2020.		
1070	BERTScore: Evaluating Text Generation with		
1071	BERT . ArXiv:1904.09675 [cs].		
1072	Tianyi Zhang, Faisal Ladhak, Esin Durmus, Percy Liang,		
1073	Kathleen McKeown, and Tatsunori B. Hashimoto.		
1074	2023. Benchmarking Large Language Models for		
1075	News Summarization . ArXiv:2301.13848 [cs].		

1090 A Instructions for Writing Summaries

1091 In this section we explain our annotation procedure
1092 and provide definitions for the terminologies in our
1093 instructions along the way. As described in Sec-
1094 tion 2, the procedure is divided into two parts: one
1095 in which an annotator works individually and the
1096 other in which they interact with another annotator.

1097 A.1 Individual Work

1098 Instructions for an individual annotator:

1099 1. Depending on the *complexity* of the conver-
1100 sations, either 1) thoroughly read the whole
1101 conversation or 2) skim through the conversa-
1102 tion to understand the general idea

1103 • *Complexity*: number of speakers, fa-
1104 miliarity of the topic to the annotators,
1105 length. For shorter conversations, it is
1106 easier to read through the whole conver-
1107 sation before moving on to summarizing,
1108 while for really longer ones, annotators
1109 would read a few comments at a time,
1110 summarize, read the next few, etc.

1111 2. Go through the conversation *comment-by-*
1112 *comment* and write a comprehensive summary
1113 that captures the content of each *comment* and
1114 any *key points*.

1115 • *comment*: all speakers' utterances are in
1116 the form of reddit comments.
1117 • *key points/moments*: also referred to as
1118 "turning points" are where the tension of
1119 the conversation or the speakers' opin-
1120 ions notably change. Annotators should
1121 highlight them in both the original tran-
1122 script and the summary in the following
1123 way: increase in tension (red), decrease
1124 in tension (blue), change in opinions to-
1125 wards disagreement (yellow), change in
1126 opinions towards agreement (green)

1127 3. Then revise the comprehensive summary to

1128 (a) change any wording that's confusing (not
1129 accurately describing the original com-
1130 ment)

1131 (b) review if the summary reflects the con-
1132 versation accurately (specifically the con-
1133 versation dynamics and tension) and add
1134 any tone indicators that might be missing

i. Indicate changing tension (e.g. curse 1135
words, all-caps, rhetorical questions, 1136
polite words) and indicate senti- 1137
ments with phrases like "sarcasti- 1138
cally," "passive-aggressively," "po- 1139
litely," etc. Use direct quotes (no 1140
need to explicitly describe the emo- 1141
tion) if they are concise and hard to 1142
capture in a summary. Focus on the 1143
highlighted elements of the conversa- 1144
tion when adding indicators in order 1145
to capture changes in tension. 1146

(c) Condense the summary to 150 words 1147
while trying to preserve the turning 1148
points from step 2 and tone indicators 1149
indicated during revision. Omit the parts 1150
of the conversation that didn't contribute 1151
much to the overall trajectory and other- 1152
wise reword for brevity. For example, 1153

i. Condense lengthy or redundant back- 1154
and-forth conversation that doesn't 1155
introduce new points (but may im- 1156
pact tension) into fewer sentences 1157
summarizing the main developments 1158

ii. Omit irrelevant comments (e.g. brief 1159
interjections by a new user that did 1160
not have any substantial follow-ups) 1161

iii. Change a few direct quotes/details to 1162
more concise sentiment words (ex. 1163
"calling this blatant racism" → "... 1164
with condemnation") 1165

iv. Other editorial changes 1166

4. After comprehensive summary, write the 1167
speaker summary by 1168

(a) First identifying the key speakers based 1169
on the comprehensive summary. 1170

• Usually whichever speakers spoke 1171
the most, but also considering those 1172
contributed to the *key moments* 1173

(b) Then for each key speaker, rereading 1174
only their comments in the original con- 1175
versation, in order to describe their spe- 1176
cific changes in tone/stance/conversation 1177
strategies and interactions/response to 1178
other key speakers 1179

1180 A.2 Interactive Work

1181 Annotators start the interaction from the following 1181
1182 setup: 1182

1183	• Annotator A: having completed the individual work for the conversation, i.e., read the original conversation and wrote the comprehensive summary and the speaker summary	providing sources, requesting sources, insulting, defending, acknowledging, conceding, rhetorical questions, invalidating, repetition, using long comment	1228
1184			1229
1185			1230
1186			1231
1187	• Annotator B: didn't read the original conversation, now writes the summary of conversation dynamics.	• tonal elements example: sarcasm, passive-aggressiveness, bluntness, rudeness, civility, neutrality, passion, harshness, strength, assertiveness, politeness, friendliness, objectivity, annoyance, frustration, tension, provocation, skepticism, demanding	1232
1188			1233
1189			1234
1190	Collaboratively, they follow these steps, which we describe from a third-person perspective for better clarity.		1235
1191			1236
1192			1237
1193	1. Annotator B reads the comprehensive summary and speaker summary out loud. They ask initial questions to Annotator A confirming the order of speaker comments and key speakers ("Speaker1 then Speaker2 then Speaker1 again?", "Speaker1 spoke the most?"), the overall stance/speaker relationship of the argument ("Speaker1 and 3 agreed, and both disagreed with Speaker2?")	• If the indicator of tone is missing or not clear, Annotator B asks Annotator A questions such as the ones below, and Annotator A often goes back to the original conversation to reread comments and provide accurate answers to the questions or even read aloud whole phrases of a comment if needed to give proper context	1238
1194			1239
1195			1240
1196			1241
1197			1242
1198			1243
1199			1244
1200			1245
1201			1246
1202	2. Annotator B begins writing the SCD by first copying the first sentence of the comprehensive summary, which often describes the overall topic of the conversation in a few words.	– B: "Was this said neutrally, or is there something about the tone that I should note?"	1247
1203			1248
1204			1249
1205			1250
1206	3. Annotator B identifies the first <i>section</i> of the comprehensive summary, highlighting the summary sentences on the document so that Annotator A can also reference.	– B: "Is the comment overtly rude, or is it just passive-aggressive or blunt?"	1251
1207			1252
1208			1253
1209			1254
1210	• <i>section</i> – usually 1-3 <i>comments</i> that fall before/in between any <i>key moments</i> . These <i>comments</i> should have a similar impact on the overall conversation dynamics, so that it makes sense to condense them into one sentence in the SCD	• Annotator A reviews the work done on this <i>section</i> and makes corrections or suggestions if they think the conversation dynamics summary isn't an accurate representation of the conversation. And, Annotator A and B would revise the sentences together.	1255
1211			1256
1212			1257
1213			1258
1214			1259
1215			1260
1216	• Annotator A may disagree with condensing the section if they think important information from within the section would be lost (e.g. different tone/rhetorical elements, argumentative stances)	• They repeat this process for each <i>section</i> .	1261
1217			1262
1218			1263
1219			1264
1220			1265
1221	4. For each <i>section</i> , Annotator B writes a corresponding summary capturing the dialogue acts, conversation strategies, and tonal elements, without any topical details.	• Annotator A rereads the whole conversation dynamics summary, noting if any part does not seem to accurately reflect the original conversation/comprehensive summary. Both people work together to correct any such cases with the question-asking method above.	1266
1222			1267
1223			1268
1224			1269
1225	• dialogue acts and conversation strategies examples: disagreement, agreement, counterargument, criticism, accusation,	– If needed, annotators would condense the summaries to be under 80 words, but usually they were already within range.	1270
1226			1271
1227			1272
		B Informativeness Check	1272
		Conversations covered in the check. We first sample 10 conversations on 10 different topics. 5 of the conversations are 'derailing' and 5 are 'non-derailing'. Each of these conversations makes one	1273
			1274
			1275
			1276

D.2 Finetuned Summarization Systems

For finetuned summarization systems, we use 40 transcript-summary pairs from our human summary dataset for finetuning, 10 pairs for development, and generate summaries for the remaining 50 test set conversations that do not have human summaries. The generated summaries are then evaluated with our downstream task in Section 4.

We first experimented with the SOTA conversation summarization systems, BART-large and DialogLED (Lewis et al., 2020; Zhong et al., 2022). Both systems previously showed strong performance on long dialogue summarization datasets with small train sets, such as AMI (train size 97) (Carletta et al., 2006) and ICSI (train size 43) (Janin et al., 2003), as reported in Zhong et al. (2022). Table 3 reports the performance brought by summaries from finetuned BART and DialogLED in our downstream task. We find that these models finetuned on the 40 human written summaries, do not produce summaries that lead to better forecasting results than procedural prompt summaries.

Additionally, we attempted to finetune GPT-3.5-turbo using OpenAI’s API. Due to the high cost OpenAI charges for finetuning and inferencing on finetuned checkpoints, we find adequate hyperparameter search unfeasible and stopped after obtaining one checkpoint with reasonable summary quality. The summaries by this checkpoint led to an accuracy of 61.9% in the downstream task, substantially lower than the accuracy brought by the procedural prompt summaries (Table 3).

Based on...	Accuracy
transcripts	56.0
procedural prompt summ.	71.5 (2.52)
finetuned BART summ.	57.5 (2.52)
finetuned DialogLED summ.	54.5 (5.26)

Table 3: Few-shot GPT derailment forecaster performance based on finetuned models summaries (for the 50 test set conversations that do not have human summaries). We include results on transcripts and procedural prompt summaries of the same 50 conversations for reference.

D.3 Other Forecasting Systems

For using GPT-3.5-turbo as few-shot classifiers, we set the sampling temperature to 0 for deterministic behaviors (because we could not set a random seed with OpenAI’s API).

Additionally, we also experimented with other classifiers using supervised training to forecast conversation derailment. We use the transcripts or the generated summaries of the train (234 conversations) and dev (100 conversations) splits of our dataset to obtain trained classifiers and run inference on the transcripts or generated summaries of the test split (100 conversations). We examine two strong baseline models for text classification for this supervised setting, namely BART and Longformer. Although these supervised models are consistently outperformed by the GPT few-shot classifier (Table 4), when comparing their performances on the generated summaries, we still find that procedural prompt summaries best help the downstream forecasting of conversation derailment, indicating that our conversation dynamics summary task indeed helps automatic systems to forecast conversation derailment.

Based on...	Acc. by Forecasting Model		
	GPT	BART	Longformer
transcripts	60	46	50
convent. summ.	58	57	58
proced. summ.	67	63	62
BART summ.	58	54	54
LED summ.	55	52	58

Table 4: Comparing different classifier architecture for derailment forecasting. Similar to Table 1, “convent. summ.” and “procedu. summ.” refer to the GPT-3.5 generated traditional prompt and procedural prompt summaries; “BART summ.” and “LED summ.” refer to summaries by finetuning BART and DialogLED. Additionally, the column name “GPT” refers to GPT-3.5 few-shot classifier.

D.4 Prompt Engineering

When developing our zeroshot and procedural prompts for dynamics summaries, we tried different synonyms for conversation dynamics and specific dynamics elements, as well as changing the phrasing of their definitions and examples. For example, instead of simply prompting the model to summarize ‘conversation dynamics’, which might appear as a novel jargon to the model’s parametric knowledge, we instruct the model to write a summary that captures the trajectory of the conversation, especially focusing on how elements like tone, sentiment, conversation strategies may change or remain the same throughout the conver-

1445
1446
1447
1448
1449
1450
1451
1452
1453
1454

1455

1456
1457
1458
1459
1460
1461

1462

1463
1464
1465
1466
1467

1468

1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489

sation. We then manually examine the quality of generated summaries for a small prompt engineering dataset (size 10) that’s disjoint with our dev and test splits. For the procedural prompts, in particular, we manually wrote example summary segments to contrast different aspects of traditional summaries with those of SCDs, and included these examples in the procedural prompt. Figure 3 shows the two prompts we eventually chose as the zeroshot and procedural prompts for SCDs.

E Qualitative Analysis

Inspired by prior literature (Habernal et al., 2018; Zeng et al., 2020; Zhang et al., 2018a; Li et al., 2020a), we focus on a set of conversation strategies related to conversation trajectories for our qualitative analysis in Section 5. Here, we present the list in Table 11.

F Miscellaneous

F.1 Transcript of the Introductory Example

We provide the transcript of our introductory example in Figure 4 to 7. The last 3 utterances of the transcript are omitted as how it appears in our dataset.

F.2 Data Collection

Anonymization. We collect human summaries for conversation transcripts from the published dataset CGA, which we accessed through ConvoKit 2.5.3. The dataset contains the usernames of the conversation participants, which we replace with ‘Speaker1’, ‘Speaker2’, and etc. to respect the users’ identity.

Annotators. All annotators for our evaluations are recruited as volunteers from university students in the US. The two annotators who wrote the summaries of conversation dynamics are co-authors of this paper. The data collection was approved by an Institutional Review Board at the authors’ institution. All annotators were informed that their data would be used for an NLP research and eventually a published paper before they gave consent.

Disclaimer of Risks. All annotators are informed that “some of the conversations presented in the annotation task can be extremely biased and offensive and speak of sensitive topics.” All annotators gave their consent to participate.

F.3 Implementation Details

For our finetuned models, we conducted hyperparameter search over learning rates [3e-6, 5e-6, 1e-5, 2e-5, 3e-5, 5e-5, 1e-4] and warmup steps ([40, 80] for summarizers and [234, 468] for classifiers), and used the default values from their original implementation for other hyperparameters. For the DialogLED and BART summarizers, we eventually used a learning rate of 3e-5 and 80 warmup steps. For the BART classifier, we used a learning rate of 3e-6 and 468 warmup steps. For the Longformer classifier, we used a learning rate of 5e-6 and 468 warmup steps. The finetuning experiments in total took about 150 GPU hrs on an Nvidia A40 GPU.

F.4 Used Artifacts

We include a list of existing artifacts we used. Some of them have been cited in the main sections of this paper above. We have closely followed their intended use.

- GPT-3.5-turbo-0613: 1509
a snapshot of GPT-3.5-turbo from June 13th, 2023. Closed-source but accessible at a low cost via OpenAI’s API, <https://platform.openai.com/docs/> 1510
- ConvoKit 2.5.3: 1514
<https://convokit.cornell.edu/>, MIT License 1515
- PyTorch 1.8: 1516
<https://pytorch.org>, BSD-3 License 1517
- Transformers 4.25: 1518
<https://github.com/huggingface/transformers>, Apache License 2.0 1519
- Scikit-learn 1.3.2: 1521
<https://scikit-learn.org>, BSD-3 License 1522

F.5 Additional Evaluation Metrics

Here, we provide additional performance metrics (precision, recall, macro-averaged F1) for different summary types, when they are evaluated with our derailment forecasting task. Each summary type is evaluated with its respective GPT-3.5 few-shot derailment forecasting model as described in Section 4.1. 1524

Derailing?	prec.	rec.	F1
False	72.7	32.0	44.4
True	56.4	88.0	68.8
macro avg	64.6	60.0	56.6

Table 5: Additional metrics for derailment forecasting on **transcripts**

Derailing?	prec.	rec.	F1
False	55.3	47.0	50.8
True	53.9	62.0	57.7
macro avg	54.6	54.5	54.2

Table 10: Additional metrics for derailment forecasting on **finetuned DialogLED summaries**

Derailing?	prec.	rec.	F1
False	57.1	66.5	61.4
True	59.9	50.0	54.5
macro avg	58.5	58.3	58.0

Table 6: Additional metrics for derailment forecasting on **traditional summaries**

Derailing?	prec.	rec.	F1
False	56.1	80.0	66.0
True	65.2	37.5	47.6
macro avg	60.7	58.8	56.8

Table 7: Additional metrics for derailment forecasting on **zeroshot prompt summaries**

Derailing?	prec.	rec.	F1
False	62.9	84.0	72.0
True	75.9	50.5	60.7
macro avg	69.4	67.3	66.3

Table 8: Additional metrics for derailment forecasting on **procedural prompt summaries**

Derailing?	prec.	rec.	F1
False	56.1	69.0	61.9
True	59.7	46.0	52.0
macro avg	57.9	57.5	56.3

Table 9: Additional metrics for derailment forecasting on **finetuned BART summaries**

[Conversation Summary]

Speakers discuss the responsibilities of caregivers of autistic children. One Speaker opens up the discussion using strong language. Speaker3 and Speaker4 begin to argue in a passive-aggressive manner, which then transitions into sarcasm, accusations, and questioning each other's logic. Speaker4 supports their point with a personal experience, which Speaker3 refutes rudely.

Will the conversation go awry (derail)?

- Yes
- No

Confidence of your answer (1 for least confident and 5 for most confident)

- 1
- 2
- 3
- 4
- 5

To what extent did the summary help you understand the topic of the conversation (on a scale of 1 to 5)?

- 1: I don't even know the general topic.
- 2
- 3: I know the general topic of the discussion.
- 4
- 5: I know how each Speaker is related to the topic.

To what extent did the summary help you understand the conversation trajectory (on a scale of 1 to 5)?

- 1: I don't have any idea of the trajectory of the conversation.
- 2
- 3: I have a general understanding of the trajectory.
- 4
- 5: I have a thorough understanding of how each Speaker contributed to the trajectory.

Figure 2: Example question for derailment forecasting based on summaries

Zeroshot Prompt:

Write a short summary capturing the trajectory of an online conversation. Do not include specific topics, claims, or arguments from the conversation. Instead, try to capture how the speakers' sentiments, intentions, and conversational/persuasive strategies change or persist throughout the conversation. Limit the trajectory summary to 80 words.

Procedural Prompt:

Write a short summary capturing the trajectory of an online conversation. Do not include specific topics, claims, or arguments from the conversation. The style you should avoid:

Example Sentence 1: "Speaker1, who is Asian, defended Asians and pointed out that a study found that whites, Hispanics, and blacks were accepted into universities in that order, with Asians being accepted the least. Speaker2 acknowledged that Asians have high household income, but argued that this could be a plausible explanation for the study's findings. Speaker1 disagreed and stated that the study did not take wealth into consideration." This style mentions specific claims and topics, which are not needed.

Instead, do include indicators of sentiments (e.g., sarcasm, passive-aggressive, polite, frustration, attack, blame), individual intentions (e.g., agreement, disagreement, persistent-agreement, persistent-disagreement, rebuttal, defense, concession, confusion, clarification, neutral, accusation) and conversational strategies (if any) such as 'rhetorical questions', 'straw man fallacy', 'identify fallacies', and 'appealing to emotions.' The following sentences demonstrate the style you should follow:

Example Sentence 2: "Both speakers have differing opinions and appeared defensive. Speaker1 attacks Speaker2 by diminishing the importance of his argument and Speaker2 blames Speaker1 for using profane words. Both speakers accuse each other of being overly judgemental of their personal qualities rather than arguments."

Example Sentence 3: "The two speakers refuted each other with back and forth accusations. Throughout the conversation, they kept harshly fault-finding with overly critical viewpoints, creating an intense and inefficient discussion."

Example Sentence 4: "Speaker1 attacks Speaker2 by questioning the relevance of his premise and Speaker2 blames Speaker1 for using profane words. Both speakers accuse each other of being overly judgemental of their personal qualities rather than arguments."

Overall, the trajectory summary should capture the key moments where the tension of the conversation notably changes. Here is an example of a complete trajectory summary.

Trajectory summary:

Multiple users discuss minimum wage. Four speakers express their different points of view subsequently, building off of each other's arguments. Speaker1 disagrees with a specific point from Speaker2's argument, triggering Speaker2 to contradict Speaker1 in response. Then, Speaker3 jumps into the conversation to support Speaker1's argument, which leads Speaker2 to adamantly defend their argument. Speaker2 then quotes a deleted comment, giving an extensive counterargument. The overall tone remains civil.

Now, provide the trajectory summary for the following conversation.

Conversation Transcript: [...]

Figure 3: Zero-shot prompt and procedural prompt for SCDs

Strategies	How they can be mentioned in dynamics summaries
Rhetorical questions	“poses a rhetorical question”, “rhetorically asks”
Attacking logic	“point out flaws in [the other speaker]’s arguments”, “accuses [the other speaker] of their logical fallacy”
Anecdotal experience	“shares a personal story”, “uses an anecdotal example”
Evidence	“cites statistics and data to support their viewpoint”, “uses external sources to support”
Juxtaposition	“makes a comparison between”, “provides a detailed explanation of the differences between”
Analogy	“uses an analogy to support”
Pointing at missing or unsupported evidence	“asks for evidence”, “criticizes the lack of evidence”
Accusing of not correctly treating their argument	“accuses [the other speaker] of not reading their arguments”, “accuses [the other speaker] of reinterpreting their positions”
Questioning one’s knowledge or attacking one’s lack of knowledge	“insulting [the other speaker]’s knowledge of [the subject]”, “accusing [the other speaker] of lacking the knowledge of [the subject]”
Hypothetical example	“proposing another hypothetical scenario”
Counterexample	“presents counterexamples”

Table 11: List of conversation strategies a speaker may use. For our qualitative analysis, we consider an SCD ‘mentions’ a conversation strategy only if it explicitly identifies what the strategy is, as shown in the examples listed in the second column of this table. For example, if a summary simply paraphrases the exact rhetorical question or the cited evidence, then we do not count it as ‘mentioning’ a conversation strategy.

Transcript:

Speaker1: Businesses aren't charities. They exist to make a profit. "Morals" and "ethics" are always trumped by profit in the business world.

Speaker2: That's.... Kind of the inherent problem.

Speaker3: For you. Whenever there is a "problem", it is usually some party wanting to further their interests. Remember that morals do not exist out there, they are a construct of society. If a majority is disadvantaged, they may use "morals" to push for their interests.

Speaker2: Unchecked capitalism is economically unsustainable. So yes, it is a problem for me, and for everyone participating in the economy. Business involves more than just profits. It involves human beings investing their labor. These are not machines. The idea that profit alone should drive our economic decisions is morally bankrupt. Human beings deserve a modicum of dignity. If you can't agree to that, and you think that slavery is OK and justified as long as the business is profitable, then I would posit that you too are morally bankrupt.

Speaker4: "Unchecked capitalism is economically unsustainable." what do you mean "unsustainable"? what happens?

edit: [MFW I get my daily reddit downvotes for questioning a socialist](<http://i.imgur.com/QoGM3.gif>)

Speaker2: Well, let's have an economic thought experiment.

Rule 1: Businesses only care about profit (and typically the short term profit at that, not longer term returns).

Rule 2: Businesses have nothing to check them from abusive practices.

Let's think of some common things people believe are good about this scenario:

1) Everyone seeks to further their own means so with supply and demand everything works out in the end!

2) If a company has poor practices, the consumer will migrate to other options.

3) If a company has poor practices, another company will be developed to compete against it, forcing it to rid itself of those poor practices.

Now let's think about whether or not those thoughts can sustain themselves in the thought experiment.

1) Why wouldn't companies seek to form cartels? We already have evidence that they do this, even among competitors. You might say that they'll be seeking to further their own means, so this is only temporary. But I would ask you to consider that if Company A and Company B can make larger profit margins colluding than they can competing, why wouldn't they? Sure, Company A might be able to make slightly higher profits down the road if they were able to beat Company B outright, but that takes short term investment (and therefore cost) and effort. With collusion, we can maximize profits with minimal effort and cost.

2) But surely if A&B collude with one another, C will come out of the woodwork and offer a better product/service at a better price, right? But will they? If A&B catch wind of C, what's to stop them from using potentially coercive means from stopping them from competing. For example: hostile takeovers. A&B are a cartel now, and they are established, so they certainly have the capital necessary to absorb C as a fledgling company.

3) If A&B are colluding, what's to stop them from slowly degrading the conditions of their workforce? Like the employees of office space only working hard enough not to get fired, A&B are only going to create conditions good enough as to not lose profit margins.

4) If A&B are able to block competitors from the industry, what motivation do they have to innovate and improve their products?

5) What's to stop A&B from doing environmentally disastrous things (pollution for example)? If short term profits are all that matter, why would they bother to care?

[continues on the next page..]

Figure 4: Example Transcript, Part 1 of 4. Its conversation ID is cz2r8ig in the CGA corpus.

[continues from the previous page]

All of this isn't to say that Government must be the check on businesses. Government can be the tool of business to enshrine their power as well. But capitalism by itself can have some pretty negative consequences when it is extrapolated out. The problem is the fundamentals of capitalism are built on a foundation of assumption. That assumption is that when people seek to fulfill their individual desires, with supply and demand and perfect competition at play, that everything reaches an equilibrium. The problem here is that perfect competition is never really at play. Like most things "perfect" it's an idea, not a reality. You can't enforce perfect competition. So the next closest thing you can do is create an environment where competition is maintained and our perpetually imperfect competition doesn't get too out of control.

Personally, I'm a believer in social democracy - I believe that capitalism can be good when it is juxtaposed by the idea of solidarity with all of mankind (as opposed to purely selfish means). That's not communism, mind you, it just means that we shouldn't lose sight of humanity in the process. When we think about that as a counter to selfish greed, we start thinking about longer term returns, we start thinking about societal benefit, we start humanizing the capitalistic process (they aren't "capital" or "human resources" or "employees" they're people investing their efforts in making these companies perform).

I hope that addresses some of your question. It was a bit open ended, so I tried not to ramble on too much.

Speaker4: that addressed ZERO of my question. you completely failed to address what happens when a capitalist system reaches its alleged "sustainability" threshold.

my question was not open-ended. what occurs when it is no longer sustainable? you said it was unsustainable. your hypothetical "what ifs" don't support your assertion.

Speaker2: Ok. Thanks for the attitude. Now I'm just going to give you a curt response.

Economic collapse. Mass unemployment. Overly polluted and toxic environments. Did you just want a parade of horrors?

Speaker4: attitude? i just explained that you didn't answer my question. you gave me an opinionated rant instead of an explanation.

what am i supposed to do, just say "oh thanks"?

"economic collapse" yeah, you said that already. why? how?

"mass unemployment" how? everybody just gradually becomes unemployed "because capitalism"?

"toxic environments" why? is there some aspect of socialism that prevents toxicity in products? does socialism provide some sort of waste-disposal service unavailable in a capitalist system? you just keep throwing out matter-of-fact assertions, but i don't see how you are arriving at your conclusions.

and apparently you interpret scrutiny as "attitude" that you take offense to... i dunno. I'm not convinced.

edit: annd im downvoted instead of having any of those valid questions answered.

Speaker2: "attitude? i just explained that you didn't answer my question. you gave me an opinionated rant instead of an explanation." Sorry, perhaps I'm just reading into the emphasis from "ZERO" and "completely". The tone of your response and the one that followed seemed to be ... dickish, for lack of a better word. If I'm reading into it, my apologies.

Economic collapse. Increasingly volatile market behaviors as a result of increasingly risky investments. Essentially: short term profit at any cost. The system fails to account for long term sustainability.

[continues on the next page..]

Figure 5: Example Transcript, Part 2 of 4.

[continues from the previous page]

Mass unemployment. Actually, we're already on our way to mass unemployment due to automation. I just read an article about us losing 5 million (net) jobs by 2020 to automation and AI (that includes the 2mil that will pick up new jobs from the new tech). You also have an increasing population who can be provided for but not necessarily enough jobs to have them making a worthy contribution.

Toxic environments. I actually didn't mention socialism, you did. I said social democracy. The primary economic driver of social democracy is capitalism. It's just checked. But even still: yes. There is something inherent in those systems that greatly reduces the risk of toxicity. The people who make the products and invest the labor are the ones who realize the gains. They are also the ones to decide whether or not those gains are worthwhile given the risks.

In a raw capitalism system: the capitalist realizes the gains and, due to the increased capital, can largely insulate themselves from the associated risks. Drinking water pollution is an example of this. The rich typically don't have this problem as they can afford to purchase potable water. The poor and working class may not always have their luxury (Flint, MI for example). So the people who have to live with the realities and consequences of their production are the ones making the decisions as opposed to someone who doesn't have to live with them or who can use their financially superior position to avoid them.

Other examples include the short term profits associated with taking unnecessary risks: see countless oil spills and deepwater horizon. Also see: covering up of global warming by oil companies or covering up of health hazards by cigarette companies.

The capitalist system *encourages* this behavior. It is financially beneficial (short term profits are encouraged over long term profits and investment - this is legally supported through cases dating back to Dodge v. Ford). A checked capitalist economy disincentivizes that behavior through regulation and social welfare. A socialist economic system would probably do the same through the reality that people who would be causing the harm are the people who have to live with the harm and are less likely to remove themselves from it.

If we want to debate the core tenets of a capitalist economy: that's fine. But I operate on the presumption that we both agree the goal of a capitalist economy is to generate profit (with a weight towards short term profit especially with larger firms having to provide quarterly earnings and financial benefit for shareholders).

To get back to OP - the minimum wage is *one* check on the default capitalist economic framework (to drive costs as low as possible). OSHA is another. The FDA is another. The SEC is another. And so on and so forth.

Speaker4: thank you for that thoughtful response. i do type in a very tonal, speech-based style, so people actually accuse me of being "a dick" or "irate" pretty frequently. i guess i will just have to get used to it.

i just mentioned socialism as an alternative to compare to in the pollution argument, i didn't mean to presume to say you were, necessarily. my apologies if it seemed that way.

So anyway, yes. I agree that capitalism creates volatile markets. However, I don't consider it to be as threatening as you seem to. I can't understand any situation that would cause a cataclysmic death or cessation of those markets. There will always simply be peaks and valleys. "The system fails to account for long term sustainability."

[continues on the next page..]

Figure 6: Example Transcript, Part 3 of 4.

[continues from the previous page]

This is really vague. What happens when it no longer becomes sustainable? Anarchy? Does the wealth stop existing? Where does it go? Does it become impossible to make purchases or be paid for your work? Does everyone die? I am really trying to pin down what exactly you mean by a system "sustaining" itself. It seems more like a system that would have to be electively given up, like a language that dies out, instead of something that can break and has to be discarded, like a broken dish or a burnt-out light bulb.

In this sense, I don't see how any other ideology or economic system would succeed or fail in any of these areas, without the "unchecked" qualifier I see so commonly applied to critique of capitalism.

"unchecked" *is* the problem. greed is the problem. overpopulation is the problem. creation of toxic byproducts and harmful processes is the problem. a sheltered ruling-class is the problem. none of these things will disappear by switching ideologies. no other economic system is better, or worse, -equipped to eradicate these things. similarly, blaming "unchecked" application of any other system for the existence of these problems would be just as silly as blaming capitalism.

Speaker2: That was my point in the beginning though. That unchecked capitalism is bad. I am merely arguing for a check on it. Something to keep it from becoming too volatile, too disruptive, too detrimental. Something to keep it contained. I view capitalism like a nuclear reactor. If you keep the reaction going and have the right containment - you've got a nice source of clean energy. But if you fail to keep it contained - you've got bad news.

To your question of what happens. It can really mean any number of things. It could mean a very slow dystopian degradation of society wherein the rich get progressively richer and the poor become bottom feeders with short lifespans. It could mean everyone dies (perhaps through experimentation gone awry - particularly in the energy or medical fields). I'm not sure about wealth ceasing to exist outright, I don't know what that would look like or how that would come about. Anarchy is certainly possible, so is revolution (as the rich get richer, the poor get angrier - extrapolate that and we have a recipe for repeating history).

Like I said, some of those problems can be avoided when the people profiting from the production are the people doing the production. It connects them with the consequences. It functions as a systematic check. Is it perfect? surely not.

I agree that greed is the problem. But greed isn't a bug in capitalism, it's a feature. Like the nuclear reaction, we're trying to harness the human motivation behind greed and seize it. But, we don't want it to get too far out of control. It's a balancing act.

Figure 7: Example Transcript, Part 4 of 4.