Coherence of Argumentative Dialogue Snippets: A Large Scale Evaluation of Inference Anchoring Theory

Anonymous ACL submission

Abstract

This paper describes a large scale experimental study (with 933 dialogue snippets and 87 annotators) addressing the research question 004 Does Inference Anchoring Theory (IAT) model the structure of coherent debate? IAT sets out the relation between dialogue structures (illocutionary acts, turns and their relations) and 800 the inferential relations between the propositions that the interlocutors put forward in their debate with each other. IAT has been used for substantial corpus annotation and practical applications. To validate the structures that the theory assigns to debates, we designed 014 an experiment for systematically comparing the coherence ratings for several variants of short debate snippets. The comparison is between original human-human debate snippets 018 and algorithmically-generated variations that comply to different degrees with the structures mandated by IAT. In particular, we utilise an algorithm for producing alternatives of the original snippets which retain structure but change 023 the content. We found that whereas the original debate snippets and their IAT-compliant variants receive high coherence ratings, snippets that violate IAT-mandated propositional relations received lower ratings (a difference that is statistically highly significant).

1 Introduction

001

011

012

027

034

039

042

The proper modeling of argumentation in dialogue is a long-standing challenge, raising questions about how individual and collective reasoning and argumentation are connected (Yu et al., forthcoming; Ivanova and Gubelmann, 2025). In particular, a significant question is how coherence relations in debate are connected to the propositional relations of logical reasoning, that is conflict/denial, and inference/implication. An important proposal clarifying this relation is Inference Anchoring Theory (IAT) (Reed, 2011; Reed and Budzynska, 2011; Budzynska et al., 2014). This theory aims to account for the coherence of debates and offers the

tools for argument corpus development (Budzynska et al., 2014), finetuning LLMs (Wu et al., 2024), and shedding light on, for example, the role of questions in debates (Hautli-Janisz et al., 2022).

043

045

047

049

051

054

055

057

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

It is notoriously difficult to evaluate discourse analysis theories and annotation schemes, such as IAT – however, we agree with the assessment of Knott (2007, Page 594), who proposes that evaluation of a theory of coherence is 'considerably more compelling as empirical support' when done by means of an application of the theory for text generation that can then be assessed against judgements from 'actual readers'.

We propose such an evaluation of IAT for examining the following research question: Does Inference Anchoring Theory (IAT) model the structure of coherent debate? We have developed an algorithm that can generate dialogue snippet variants whose structures are entirely or partially IAT-compliant. We obtain coherence ratings, by human readers, for these dialogue snippets as well as for naturallyoccurring dialogue snippets. If IAT captures the structure that accounts for coherence, such IATcompliant generated snippets should be at least as coherent as the original dialogues. Additionally, ablated versions of the algorithm, only being partially IAT-compliant, should lead to less coherent dialogue, if IAT is valid.

The current work sits firmly within Computational Linguistics, as it uses an empirical study enabled by computational means, i.e. an algorithm for generating dialogue snippet variants - to validate IAT, a linguistic theory within the remit of discourse analysis and pragmatics.

At the heart of this paper is an experiment in which we ask judges to rate the coherence of both human-human debate snippets and variants that have been algorithmically generated at different levels of compliance with IAT, see Table 1 for two examples of such snippets.

In the next section, we describe IAT in some

Method	Dialogue example
A human-human de-	Witness: Actually, we need to think about, how we have a more pater-
bate snippet from the	nalist system for people like, for example, a young guy that lost his job
Moral Maze corpus on	at Tesco, mainly because he wasn't turning up to work on time, which is
topic of the Welfare	mainly because he was smoking a lot of spliff and he was basically very
state	disorganised. If you took the people around this table this evening, and
	you took away all of our contacts, our qualifications our great jobs and
	so on, we'd still have more internal resources than that guy.
	Panellist: Isn't the reality that in the last ten or twenty years there's
	been a massive transfer of wealth from the poor to the rich, and from the
	young to the old, and that what you're really trying to do is to justify that
	by blaming the poor for their position?
	Witness: No, I'm not trying to justify anything.
A variant of the de-	Witness: I think that all humans should be vegan. In the sense that a
bate snippet above	world of veganism would be a more ethical world: its morals would bring
generated using our	benefits to human society.
Dialogue Propositional	Panellist: Don't you think that killing animals for food is a survival
Content Replacement	instinct, and so not inherently unethical or morally blameworthy?
(DPCR) algorithm,	Witness: I don't think so, I think that instinctive, natural behavior is
with as input the dia-	counterproductive can create problems, both for the individual and so-
logue structure of the	ciety, and both might want it removed. If both deem it immoral and
Moral maze dialogue	unethical, then it is it as such and the unwanted behaviors should be shied
shown above and an	away from and hopefully removed if possible.
argument map on the	
topic of Veganism	

Table 1: Illustration of two instances of our experimental materials. At the top is a dialogue snippet from the Moral Maze radio programme and at the bottom is a dialogue snippet on a different topic but with the same dialogue structure. Both dialogues consist of three turns by Witness, Panellist and Witness. The first turn starts with an assertion and a claim in support of it. In the second turn, the content of the first turn is contradicted and the third turn contradicts the content of the second turn. We discuss the detailed representation of the underlying dialogue structure in Section 2.

100

102

103

104

106

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

detail. The remainder of the paper follows the standard American Psychological Association (APA) format for reporting experimental research.

2 **Inference Anchoring Theory (IAT)**

For readers unfamiliar with IAT we provide a short introduction by discussing a debate snippet from the Moral Maze corpus (Janier, 2017) in terms of IAT. Figure 1 shows a snippet of Moral Maze dialogue on the Welfare state annotated with IAT. On the right-hand side, we can see four locutions. For example, the first one states 'Neil: Actually, we need to think about, how we have a more paternalist system for people like that' (with speaker Neil). As the example shows, a locution consists of a speaker designation and an utterance. Each locution is anchored to a *proposition* (shown in the four left-most blue boxes) via an Illocutionary Connection (IC) (in the middle yellow boxes). In this case, the two first locutions' Illocutionary Connection, i.e. IC, is 'Asserting', the third is 'Assertive Questioning' and the last one is 'Asserting'. Propositions and ICs represent the propositional content and Illocutionary force, i.e. the speaker's communicative intention, from Speech Act theory (Searle, 1969).

The locutions are linked by a transition box, which signifies that a locution is a reply or response to its predecessor. Each transition between locutions is anchored, via an Illocutionary Connection for Transition (ICTA), to a Propositional relation. In our example, two Transitions are anchored via the 'Arguing' Illocutionary connection to the 'Inference' Propositional relation and one Transition is anchored via the 'Disagreeing' Illocutionary connection to the Conflict Propositional relation. As the example shows, 'Inference' is used when one proposition provides a reason to accept the other proposition. In contrast, 'Conflict' is used when one proposition provides a reason to not accept the other proposition.¹

3 Method

3.1 Materials

Coherence Rating Scale For our experiment, we developed, based on pilot studies and previous work on coherence annotation such as Cervone and

Riccardi (2020), a scale from 1 to 7 for rating the 129 coherence of argumentative dialogue snippets (with 130 1 =**incoherent** and 7 =**coherent**).² Annotaters are 131 asked to rate a dialogue snippet as coherent if the 132 following apply: 133 (a) all sentences in the dialogue make sense by themselves and are clear (at the point in 135 dialogue where they occur), and 136 (b) all sentences in the dialogue link together 137 well with each other so that the dialogue is 138 clear and sensible. 139 **Dialogue Snippet Variants** To address our aim 140 of validating IAT, the participants in the current 141 study rated the coherence of dialogue snippets be-142 longing to one of the following categories: 143 1. original naturally-occurring argumentative di-144 alogue snippets (MORAL MAZE_{original}) from 145 the Moral Maze corpus, 146 2. original naturally-occurring argumentative di-147 alogue snippets, but after contextual enhance-148 ment: i.e., where there is anaphora or el-149 lipsis, we manually expand these to make 150 the dialogue more self-contained (MORAL 151 $MAZE_{original}^{+context}$), since this could affect coher-152 ence ratings, 153 3. argumentative dialogue snippets generated by 154 means of the Dialogue Propositional Content 155 Replacement (DPCR) algorithm described be-156 157 4. argumentative dialogue snippets generated by the 'No sentence templates' algorithm 159 (DPCR_{-templ}). This algorithm generates 160 new argumentative dialogues according to the 161 same algorithm as DPCR but without apply-162 ing sentence templates. 163 5. argumentative dialogue snippets generated by 164 the 'Random propositional relations' algo-165 rithm (DPCR $^{-rel}$). This algorithm applies 166 sentence templates corresponding with Illocu-167 tionary Connections (ICs) in locution patterns 168 (LPs), but selects a random propositional rela-169 tion rather than the relation selected according 170 to IAT, and 171

low,

¹IAT singles out further propositional relations, for example, Rephrase (when one proposition is more or less a paraphrase of the other) but we ignore them for the purpose of this paper.

²The guidelines used in the experiment can be found in Appendix E on Page 18. Further detail can be found in the Data Supplementary Materials folder.



Figure 1: Example from Moral Maze Welfare State (Map6273, 2017). In this episode, Neil and Clifford are respectively a Witness and a Panellist.

6. argumentative dialogue snippets generated by the 'No sentence templates and random propositional relations' algorithm (DPCR^{-rel}_{-templ}). This algorithm selects a random proposition and does not apply the sentence templates to generate locutions.

172

173

178Note that for our study we use four algorithms179for snippet generation: the full DPCR algorithm180as well as three ablated versions of this algo-181rithm. We also have human-generated argumen-182tative dialogues snippets (MORAL MAZE_original183and MORAL MAZE_context). As shown in Table 2184these algorithms generated 883 dialogue snippets.185Additionally, there are 50 dialogues from MORAL

MAZE_{original} and MORAL MAZE^{+context}_{original} adding up to total of 933 snippets. The Data part of the Supplementary Materials for this paper includes the full set of dialogue snippets. Additionally, for representative examples, see Appendix F on Page 19.

186

188

189

190

191

192

193

194

195

196

197

198

200

Once generated, the argumentative dialogue snippets were split into batches of 13. In each batch, two snippets were repeated twice each, to be used for annotator quality control. The two repeated snippets that were presented twice at random places in the batch allowed us to assess the annotator's self-consistency. Overall, we have 71 batches of 15 dialogues (13 plus 2 repetitions) and 1 batch of 12 dialogues (10 plus 2 repetitions).

Algorithm	Brexit	Veganism	Vaccination	Total
DPCR	72	75	72	219
DPCR _{-templ}	72	75	72	219
DPCR ^{-rel}	74	72	75	221
$\mathbf{DPCR}^{-rel}_{-templ}$	75	74	75	224
Total	293	296	294	883

Table 2: Number of argumentative dialogue snippets generated per topic and per algorithm for Brexit, Veganism and Vaccination.

Method for Dialogue Propositional Content Replacement (DPCR) For the current work, we made use of an enhanced version of the Moral Maze MM2012c dataset (Janier, 2017): the QTMM2012c+ dataset (Amidei et al., 2021). The latter includes the following additional information: (a) each speaker is labelled with their role (one of Chair, Panellist or Witness), (b) speakers are associated with a stance towards the claim or thesis under discussion (neutral, pro and con) and (c) information on the locutions chronological order is made explicit.

202

207

209

210

213

214

215

216

217

218

219

221

226

229

230

237

The second main resource that the current work draws on are argument maps. An argument map, such as the tree-structured one depicted in Figure 2, starts with a thesis (top claim, blue box). The thesis can be supported or attacked by *pro* (green dashed boxes) and *con* (red boxes) arguments. In turn, both pro-arguments and con-arguments can branch into subsequent arguments that support or attack them. Argument maps and related structures such as argument graphs have been used previously to drive persuasive chatbots, see Chalaguine and Hunter (2020). The DPCR algorithm is not tied to a specific dataset of argument maps, but for our study we will be using maps with claims from Kialo.com.

In a nutshell, with DPCR we take an existing snippet of an argumentative dialogue from an argumentative dialogue corpus and replace its locutions with claims lifted from an argument map on a different topic whilst retaining the IAT dialogue structure, including propositional relations between the locutions' contents. Formally, given an argumentative dialogue snippet D consisting of the sequence of locutions l_1, \ldots, l_n on topic T with:

• *locutions* as the set of possible locutions;

• $type : locutions \longrightarrow dialogue_act_type; ^3$

- $speaker_role : locutions \longrightarrow roles; ^4$ 238
- content : locutions \longrightarrow propositions; ⁵ 239
- $prop_rel$: $propositions \times propositions$ 240 $\longrightarrow propositional relation.$ 6 241

242

243

244

245

246

247

248

249

250

251

253

254

255

257

258

259

260

261

262

Argumentative Dialogue Propositional Content Replacement (DPCR) is defined as obtaining a dialogue snippet $D' = l'_1, \ldots, l'_n$ on topic T' from a dialogue snippet $D = l_1, \ldots, l_n$ on topic T such that:

- for all $1 \le x \le n : type(l_x) = type(l'_x)$
- for all $1 \le x \le n$: $speaker_role(l_x) = speaker \ role(l'_x)$
- for all $1 \le x, y \le n$: $prop_rel(content(l_x), content(l_y))$ = $prop_rel(content(l'_x), content(l'_y))$

This definition stipulates what counts as DPCR, i.e. replacing propositional content on topic T with content on topic T' applied to argumentative dialogue snippet D on topic T, resulting in snippet D': as we replace locutions on one topic for those on another, (a) the dialogue act types and speaker roles belonging with the replaced for locutions should remain the same and (b) where there are propositional relations between the contents of the original locutions, these should also hold between the contents

³For more detail on the dialogue act types used in this paper, we refer to Table 7 in Appendix A.

⁴For this paper we use the roles *chair*, *panellist* and *witness*.

⁵Propositions are represented as paraphrases of the locutions, with context-dependence removed where possible.

⁶For the purpose of this work we only distinguish two propositional relations: *pro* (or *inference*) and *con* (or *conflict*). The labels *pro/con* are used for propositional relations in argument maps to signify support (*pro*) and opposition (*con*) between two propositions. These correspond to the IAT propositional relations *inference* and *conflict*. Note that by restricting our work to these two relations, the current function *prop_rel* is partial. For pairs of propositions where the relation between the propositions is not one of the aforementioned two relations, we assume that it maps to \star .



Figure 2: Example of argument map about veganism, with claims based on (Veganism example, 2022).

263of the replacement locutions (taken from argument264map). Thus we obtain a new snippet that has the265same IAT structure as the original snippet, but deals266with a different topic. DPCR $_templ$ violates (a) by267rendering almost all acts as assertions, DPCR $^{-rel}$ 268violates (b) by selecting propositional relations at269random, and DPCR $^{-rel}_{-templ}$ violates both (a) and (b).270Appendix C on Page 16 contains a full description271of the DPCR algorithm and its ablations, whilst the272code is supplied as Supplementary Material.

3.2 Participants

274

275

276

277

278

281

282

285

For our experiment, we used the Amazon Mechanical Turk platform (Mturk, 2022) with 10 annotators per batch. Each annotator was paid \$4, for a task of 20 minutes.⁷ The annotators were Master annotators⁸ from the UK and USA with as their minimum education a US Bachelor degree. We had 89 annotators who performed the task. Two of them were rejected resulting in data being used from 87 annotators. The two annotators were rejected on the basis of a test-retest setup. In each batch, the test-retest setup was based on two dialogues each being repeated once. We expected the participants to assign identical or close scores to identical (repeated) dialogues. We split the scores into three sets: $\{1, 2, 3\}, \{3, 4, 5\}$ and $\{5, 6, 7\}$. If the scores from a repeated dialogue were different and part of two different sets, then we consider the test failed and rejected the annotator. 287

288

290

291

292

294

295

297

298

299

300

301

302

304

305

306

307

308

309

310

311

312

3.3 Design

We start from the research question *Does Inference Anchoring Theory (IAT) model the structure of coherent debate?* Our overall claim or hypothesis is that it does. We break this down into four testable hypotheses:

- (H1) DCPR-generated snippets are at least as coherent as MORAL MAZE snippets.
- (H2) DCPR-generated snippets are more coherent than $DPCR^{-rel}$ snippets.
- (H3) DCPR-generated snippets are more coherent than DCPR_{-templ} snippets.
- (H4) DCPR-generated snippets are more coherent than $DPCR_{-templ}^{-rel}$ snippets.

The first hypothesis checks that the level of coherence of IAT-structured generated dialogue snippets is at least at the same level as that of natural dialogues. The remaining hypotheses compare snippets that are fully IAT-compliant with those that are only partially or not at all compliant. Together, these hypotheses tests to what extent the structures

⁷With the task taking up to 20 minutes at £3.38 (based on exchange rate at the time), this amounts to remuneration at £10.14/hour. When we carried out the experiment, in July 2022, the minimum wage in the UK was £9.50/hour.

⁸Master Workers are a top Worker of the MTurk marketplace. For more details see Mturk FAQs (2022).

posited by IAT allow us to create dialogue snippets 313 that, on the one hand, are comparable in coherence 314 with snippets from naturally-occuring dialogue and, 315 on the other hand, are superior in coherence when compared with dialogue snippets that at best only partially conform with with IAT dialogue structure.

3.4 Procedure 319

321

322

323

327

328

331

333

334

335

339

341

342

345

Annotators judged one debate snippet at a time. Snippets were grouped into batches, as described above, where the order was randomised per participant (to avoid ordering effects). Annotators could annotate more than one batch, but never the same batch twice. 325

4 **Results**

Annotator reliability Table 3 reports the Inter Annotator Agreement (IAA) value measured with two different metrics, to provide a good overview of the data reliability.⁹ The values in Table 3 are based on the IAA for each of 72 batches.

Value	%	AC2
Mean	0.82	0.39
Max	0.92	0.80
Min	0.77	0.21
Median	0.84	0.49
Variance	0.0008	0.01

Table 3: Value of Inter Annotator Agreement measured among batches. Where % is the Percent Agreement and AC2 is the Gwet AC2 coefficient.

In our experiment, the lower categories 1–4 are used much less than the higher categories 5-7. This makes our annotation unbalanced towards the lower categories. Under such conditions, chancecorrected coefficients such Krippendorff's α (Krippendorff, 1980), Fleiss's κ (Fleiss, 1971) and Cohen's κ (Cohen, 1960) are subject to the prevalence paradox (Artstein and Poesio, 2008) and suboptimal. For this reason, we decided to report IAA based on the Gwet AC2 coefficient (Gwet, 2014a) which is deemed to be more robust.

To interpret the IAA values we used the Landis and Koch (1977) benchmark scale as revised and adjusted by Gwet (2014a).¹⁰ Based on this analysis we got a level of agreement equal to or higher than *fair* for 83% of the batches. More precisely, 3%of the bathes reached a substantial level of agreement, 30% of the bathes reached a *moderate* level of agreement and 50% of the batches reached a fair level of agreement. Finally, a slight level of agreement was reached for 17% of the batches. Judging the coherence of a dialogue is not straightforward. Many factors can impact dialogue coherence, and make a dialogue more or less coherent. This made the task of judging dialogue coherence a subjective one. Accordingly, we consider the agreement reached in our study a satisfactory level of agreement.

346

347

348

351

352

353

355

356

357

358

359

360

361

363

364

365

366

367

369

370

371

372

373

374

375

376

377

378

Hypotheses Table 4 shows the results of our empirical evaluation. We discuss these in terms of the hypotheses from Section 3.3.

Algorithm	Med. Coh.	Mean Coh.
DPCR _{-templ}	7	5.99
DPCR	6	5.88
MORAL M $^{+context}_{original}$	5	5.16***
MORAL Moriginal	5	4.9***
$DPCR^{-rel}_{-templ}$	5	4.66***
DPCR ^{-rel}	5	4.46***

Table 4: Experiment results. Med/Mean Coh. Score is the median/mean of the coherence scores given to an algorithm. *** indicates that the difference between the algorithm in this row and DPCR, measured by the Student's t-test, is highly significant (at $P \le 0.001$).

(H1) DCPR-generated snippets are at least as coherent as MORAL MAZE snippets. This hypothesis is confirmed: DCPR-generated snippets are not just as coherent as MORAL MAZE snippets but even more coherent (according to the raters): coherence of DCPR is higher than MORAL MAZE original and MORAL MAZE $_{original}^{+context}$ (5.88 versus 4.9 and 5.16, $P \le 0.001$).

(H2) DCPR-generated snippets are more coherent than $DPCR^{-rel}$ snippets. This hypothesis is also confirmed (5.88 versus 4.46, $P \leq 0.001$).

(H3) DCPR-generated snippets are more coherent than DCPR-templ snippets. This hypothesis could not be confirmed. There is no statistically significant difference between DCPR-generated and DCPR_{-templ} snippets (5.88 versus 5.99).

⁹All the criteria were measured by the use of *irrCAC* library provided by the R software (Gwet, 2014b). More specifically we used the functions pa.coeff.raw() and gwet.ac1.raw(), all with ordinal weight.

¹⁰Also in this case, to interpret the IAA values, we used

the *irrCAC* library provided by the *R* software (Gwet, 2014b). More specifically we used the functions landis.koch.bf().

382

384

391

394

397

400 401

402

403

404

405

406

(H4) DCPR-generated snippets are more coher-
ent than $DPCR_{-templ}^{-rel}$ snippets. This hypothesis is
also confirmed (5.88 versus 4.46, $P \le 0.001$).

Algorithm	Med. Coh.	Mean Coh.
Brexit		
DPCR _{-templ}	6	5.91
DPCR	6	5.86
$DPCR^{-rel}_{-templ}$	5	4.64***
DPCR ^{-rel}	5	4.50***
Veganism		
DPCR _{-templ}	7	6.01
DPCR	6	5.82
$DPCR^{-rel}_{-templ}$	5	4.64***
DPCR ^{-rel}	5	4.31***
Vaccination		
DPCR _{-templ}	7	6.06
DPCR	6	5.95
$DPCR^{-rel}_{-templ}$	5	4.71***
DPCR ^{-rel}	5	4.57***

Table 5: Median and mean coherence scores by topic. *** indicates highly significant differences with DCPR (Student's t-test, $P \le 0.001$).

Table 5 compares the DPCR-generated variants per topic. The overall results from Table 4 are reproduced: DPCR_{-templ} and DPCR are tied in first place, with no statistically significant difference between them, whereas both outperform DPCR^{-rel} and DPCR^{-rel}.

Table 6 shows the mean scores depending on the number of turns per dialogue. The table suggests that the perceived dialogue coherence is impacted by the number of turns. For the DPCR and DPCR_{-templ}-generated snippets, as there are more turns, the score decreases gradually. In contrast, for MORAL MAZE $^{+context}_{original}$ and MORAL MAZE $_{original}$ the score increases as the number of turns increases. Overall, the trend is that as turn number increases, the diffence between coherence levels of, on the one hand, the DPCR and DPCR-templgenerated snippets and, on the other hand, the MORAL MAZE *original* and MORAL MAZE *original* disappears. In contrast, for the ablated versions $DPCR^{-rel}$ and $DPCR^{-rel}_{-templ}$, as dialogue length increases coherence decreases.

5 Discussion

Inference Anchoring Theory (IAT) is a widely used theory that presents an appealing perspective on

Algorithms	2 Tns	3 Tns	4 Tns
DPCR _{-templ}	6.27	5.84	5.5
DPCR	6.13	5.63	5.41
Moral maze $_{original}^{+context}$	5.42	4.56	5.35
Moral maze _{original}	4.89	4.34	5.53
$\mathrm{DPCR}^{-rel}_{-templ}$	4.78	4.56	4.44
$DPCR^{-rel}$	4.64	4.28	4.28

Table 6: Average score per number of turns (Tns).

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

how the illocutionary and transitional dimensions of argumentative dialogue are intertwined with logical relations of conflict and inference between propositional contents. It has proven effective for, among other things, discourse annotation. The current paper provides the, to our knowledge, first empirical evidence that the underlying structures IAT assigns to debates account, at least partially, for the coherence of those debates.

We tested four hypotheses derived from the claim that *IAT models the structure of coherent dialogue snippets*. One hypothesis (H3) regarding DPCR_{-templ} was not confirmed and we discuss this result in some detail in the limitations section. The other three hypotheses – (H1), (H2) and (H4) – were confirmed: We saw that IAT-generated snippets were at least as coherent as naturally-occurring dialogue snippets from the Moral Maze corpus. We also saw that dialogue snippets whose underlying propositional relations were selected at random (which is the case for both DPCR^{-rel} and DPCR^{-rel} and DPCR^{-rel}, rather than driven by IAT, are inferior in coherence to dialogue snippets that conform with the propositional relations mandated by IAT.

Apart from the important partial validation of IAT, the current work contributes to the computational study of argumentative dialogue by offering the DPCR algorithm (full and ablated versions) and its implementation for use by the research community as well as the corpus of 933 generated and natural-occurring dialogue snippets together with their coherence ratings, with each snippet rated by 10 annotators out of group of 87 annotators.

6 Limitations

The current study has a number of limitations:

 Only short, up to four-turn, snippets were in scope, given that longer snippets typically do not have a fully connected IAT structure underpinning them. IAT works well where relations

between locutions can be mapped to under-446 lying propositional relations. Where these 447 are absent, other factors may influence dia-448 logue coherence. This requires further study 449 and potentially use of constructs from other 450 discourse and dialogue structure theories that 451 go beyond argumentative/propositional rela-452 tions – e.g. RST (Mann and Thompson, 1988), 453 SDRT (Asher and Lascarides, 2003), or QUD 454 (Ginzburg, 2012). 455

• The DPCR and DPCR-templ algorithm-456 generated snippets were judged as more co-457 herent than the original dialogue snippets. 458 This may be because the claims that the al-459 460 gorithm takes from the Kialo maps are generally well-written and self-contained. In con-461 trast, some of the original spoken language 462 locutions in the Moral Maze snippets are less 463 self-contained and context-dependent. We 464 tried to compensate for this by creating ver-465 sions of the Moral Maze snippets with added 466 context. Adding context did help somewhat, 467 with MORAL MAZE^{+context} dialogues rated 468 slightly higher (regarding coherence) than 469 MORAL MAZE_{original}, but still not at the level 470 of DPCR and DPCR-templ. 471

• Contra our hypothesis (H3), the DPCR_{-templ} algorithm-generated snippets were rated as highly in terms of coherence as DPCR-generated snippets. We had expected that by switching off the template generation, the coherence would detioriate. The idea behind the template generation was to convert propositions (stated as assertions) into the correct dialogue acts, i.e. the act type observed in the original naturally-occuring snippet from which the generated snipped was derived via DPCR.

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

A qualitative analysis of the dialogue snippets generated with DPCR_{-templ} and DPCR revealed coherence score differences where the snippets begin with a speaker that uttered a questioning followed by an asserting. This suggests that there is scope to improve the relevant sentence templates for this situation. For example:

> Assertive Questioning: Do you believe that the UK should remain in the EU if a hard Brexit is the only alternative option?

followed by:

Asserting: In other words I think	
that by remaining in the EU, the UK	
would be able to operate broadly	
as before but with clear caveats re-	
garding some issues that concern	
its citizens.	

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

In cases like this, a dialogue generated with the DPCR $_{-templ}$ algorithm can result in dialogue that is perceived as more coherent. For example:

The UK should remain in the EU if	
a hard Brexit is the only alternative	
option.	

followed by:

By remaining in the EU, the UK would be able to operate broadly as before but with clear caveats regarding some issues that concern its citizens.

As illustrated above, the DPCR-templ algorithm does not make use of the Illocutionary Connections (ICs) to rephrase the argument map claims that are used. It will not convert the argument map claim to the sentence type associated with the IC and instead always use the declarative sentence from the argument map verbatim. However, for those locutions where there is no argument map claim involved, such as various forms of challenging, DPCR-templ does use the relevant canned text: for example a Pure challenging IC is realised as one of the following 'Why is that?' or 'Why?'. Similarly, there is canned text for Assertive and Rhetorical Challenging. This means that $DPCR_{-templ}$ will not just yield a sequence of assertions: for all the aforementioned ICs involving challenging, variety is introduced through the canned text associated with these ICs. All in all this means that $DPCR_{-templ}$ is at least partially IAT-compliant after all.

Additionally, though DPCR $_{-templ}$ does not convert argument map claims into questions (i.e. interrogative sentences) where the IC requires this, in dialogue, whether a locution with a declarative sentence type is intended as a question can usually be inferred from the

communicative context (Beun, 1990). In our 544 argumentative dialogue set-up, involving dis-545 cussions between a panellists and witnesses 546 about contentious topics, a natural interpretation of locutions consisting of a declarative sentence is as raising questions for discussion 549 by the other party – this is also in line with the 550 more general idea that assertions can can initiate issues, i.e. introduce questions for discussion, which then become part of the QUD, i.e. 553 questions under discussion (Ginzburg, 2012). 554 As part of our qualitative analysis we also ob-555 served that the dialogue generated with DPCR does look more like natural dialogue, than the ones generated with the DPCR-templ algorithm. For an example of the contrast between the two types of dialogue that are generated, 560 see Appendix F starting on Page 19: Tables 20 and 21. Note the difference between the DPCR dialogue, which involves for example questions and hedges and DPCR-templ 564 dialogue, which is a simple sequence of as-565 sertions. It may be that naturalness needs 566 to be considered separately from coherence, which was the focus of this paper. Whereas we found evidence for the relation between 569 propositional relation choice and coherence, 570 this relation does not seem to be as strong 571 or existent between dialogue act type choice 572 and coherence. Further research is needed to establish whether the latter relation is more 574 closely associated with dialogue naturalness. • We consider the current study large scale, with almost 1000 dialogue snippets judged by 87

• We consider the current study large scale, with almost 1000 dialogue snippets judged by 87 annotators. Of course, size is relative and compared to datasets and annotations undertaken by commercially-driven labs (e.g. to train LLMs), our data set is comparatively small. And yet, for theory-driven empirical work, this study is of a significant size. It is also worth noting that with this paper we are making *all* data and code available.

Acknowledgements

587

588

589

590

591

578

579

581

583

References

. . .

Jacopo Amidei, Paul Piwek, and Svetlana Stoyanchev. 2021. QTMM2012c+: A Queryable Empirically-Grounded Resource of Dialogue with Argumentation. In 5th Workshop on Advances in Argumentation in Artificial Intelligence.

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

- Ron Artstein and Massimo Poesio. 2008. Inter-coder agreement for computational linguistics. *Computational linguistics*, 34(4):555–596.
- N. Asher and A. Lascarides. 2003. *Logics of Conversation.* Cambridge University Press, United States.
- Robbert-Jan Beun. 1990. The recognition of dutch declarative questions. *Journal of Pragmatics*, 14(1):39–56.
- Kasia Budzynska, Mathilde Janier, Chris Reed, Patrick Saint-Dizier, Manfred Stede, and Olena Yakorska. 2014. A model for processing illocutionary structures and argumentation in debates. In *Proceedings* of the Ninth International Conference on Language Resources and Evaluation (LREC'14), pages 917– 924, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Centre for Argument Technology. 2023 manuscript. A Quick Start Guide to Inference Anchoring Theory (IAT).
- Alessandra Cervone and Giuseppe Riccardi. 2020. Is this dialogue coherent? learning from dialogue acts and entities. In *Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, pages 162–174, 1st virtual meeting. Association for Computational Linguistics.
- Lisa A Chalaguine and Anthony Hunter. 2020. A persuasive chatbot using a crowd-sourced argument graph and concerns. In *Computational Models of Argument*, pages 9–20. IOS Press.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.
- Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. *Psychological bulletin*, 76(5):378.
- Jonathan Ginzburg. 2012. *The interactive stance: Meaning for conversation*. Oxford University Press.
- Kilem L. Gwet. 2014a. *Handbook of inter-rater reliability: The definitive guide to measuring the extent of agreement among raters*. Advanced Analytics, LLC.
- Kilem L. Gwet. 2014b. irrCAC library home page. https://rdrr.io/cran/irrCAC/. [Online; accessed 2022].
- Annette Hautli-Janisz, Katarzyna Budzynska, Conor McKillop, Brian Plüss, Valentin Gold, and Chris Reed. 2022. Questions in argumentative dialogue. *Journal of Pragmatics*, 188:56–79.
- Rositsa V Ivanova and Reto Gubelmann. 2025. The shift from logic to dialectic in argumentation theory: Implications for computational argument quality assessment. In *Proceedings of the 31st International Conference on Computational Linguistics*,

Computational Linguistics. Mathilde Janier. 2017. Dialogical dynamics and argumentative structures in dispute mediation discourse. Ph.D. thesis, University of Dundee. Alistair Knott. 2007. Book Reviews: Coherence in Natural Language: Data Stuctures and Applications, by Florian Wolf and Edward Gibson. Computational Linguistics, 33(4):591-595. Klaus Krippendorff. 1980. Content analysis; an introduction to its methodology. A Sage Publications, Beverly Hills, CA. J. Richard Landis and Gary G. Koch. 1977. The measurement of observer agreement for categorical data. biometrics, pages 159–174. William C. Mann and Sandra A. Thompson. 1988. Rhetorical structure theory: Toward a functional theory of text organization. Text, 8(3):243-281. MM2012c Map6273. 2017. Moral Maze map 6273. http://ova.arg-tech.org/analyse. php?url=local&plus=true&aifdb=6273&akey= 49adc508e74cadf6633d666f9644000e. [Online; accessed 2022]. Mturk. 2022. Mturk home page. https://www.mturk. com/. [Online; accessed 2022]. Mturk FAQs home page. Mturk FAQs. 2022. https://www.mturk.com/worker/help. [Online; accessed 2022]. Chris Reed. 2011. Implicit speech acts are ubiquitous. Why? They join the dots. In Argument cultures: Proceedings of the 8th international conference of the Ontario Society for the Study of Argumentation. Chris Reed and Katarzyna Budzynska. 2011. How dialogues create arguments. In Proceedings of the 7th conference on argumentation of the International Society for the Study of Argumentation, pages 1633-1645. John R. Searle. 1969. Speech acts: An essay in the philosophy of language. Cambridge University Press. Kialo Veganism example. 2022. All humans should be vegan. https://www.kialo.com/ all-humans-should-be-vegan-2762?path= 2762.0~2762.1. [Online; accessed 2022]. Yuetong Wu, Yukai Zhou, Baixuan Xu, Weiqi Wang, and Yangqiu Song. 2024. KnowComp at DialAM-2024: Fine-tuning Pre-trained Language Models for Dialogical Argument Mining with Inference Anchoring Theory. In Proceedings of the 11th Workshop on Argument Mining (ArgMining 2024), pages 103-109, Bangkok, Thailand. Association for Computational Linguistics.

pages 4789-4802, Abu Dhabi, UAE. Association for

647

663

671

672

673

674 675

677

679

681

683

684

Liuwen Yu, Réka Markovic, and Leendert Van der Torre. forthcoming. Thirteen Challenges in Formal and Computational Argumentation. In *Handbbook of Formal Argumentation Vol. 3. / Journal of Applied Logics.*

697

698

699

700

701

702 APPENDIX

703 A IAT Dialogue Act Types

Illoctionary Connection (IC)
Questioning
Rhetorical Questioning
Assertive Questioning
Pure Questioning
Challenging
Rhetorical Challenging
Assertive Challenging
Pure Challenging
Others
Asserting
Popular Conceding
Yes
No

Table 7: Types of Dialogue Acts, referred to as Illocutionary Connections (ICs) in Inference Anchoring Theory (IAT). Detailed descriptions of these acts can be found in Centre for Argument Technology (2023 manuscript)

B Student's t-test statistics

The Student's t-test was used for comparing the DPCR algorithm with the other algorithms.¹¹ The measure was performed on coherence scores associated with the dialogue generated by each algorithm. In the same fashion, we performed the Student's t-test per topic. In this case, we focus on the coherence scores associated with dialogue snippets grouped by topic.

Table 8 reports statistics related to the Student's t-test. Similarly, Tables 9, 11, 10 report statistics related to the Student's t-test respectively for the case of Brexit, vaccination and veganism. Table 12 and Table 13 report respectively the standard deviation of the coherence scores measured for each algorithm and the standard deviation of the coherence scores measured for each algorithm per topic. Table 14 and Table 15 report respectively the variance of the mean coherence scores and the variance of the mean coherence scores per topic.

Algorithms	t-test score	p-value	Degree of freedom
DPCR / MORAL MAZE _{original}	5.77	1.442593054380911e-06	35.58
DPCR / MORAL MAZE ^{+context}	4.53	8.436201745497882e-05	30.48
DPCR / DPCR_templ	-1.88	0.06	504.56
DPCR / DPCR ^{-rel}	19.50	2.22629836154952e-63	464.35
$DPCR / DPCR^{-rel}_{-templ}$	15.99	3.2683704258661824e-46	464.55

Table 8: Student's t-test statistics.

Algorithms	t-test score	p-value	Degree of freedom
DPCR / DPCR_templ	-0.44	0.65	159.47
DPCR / DPCR ^{-rel}	10.69	1.2646138185639769e-20	154.46
DPCR / DPCR $^{-rel}_{-templ}$	9.49	3.4813207445337696e-17	157.86

Table 9: Student's t-test statistics for the topic Brexit.

Algorithms	t-test score	p-value	Degree of freedom
DPCR / DPCR_templ	-1.68	0.09	169
DPCR / DPCR ^{-rel}	11.54	5.687575374991213e-23	152.14
DPCR / DPCR $^{-rel}_{-templ}$	8.88	1.3842848206119827e-15	156.84

Table 10: Student's t-test statistics for the topic veganism.

¹¹We used the function $stats.ttest_ind()$ provided by the python library SciPy for the Student's t-test, setting the variable $equal_var = False$.

Algorithms	t-test score	p-value	Degree of freedom
DPCR / DPCR _{-templ}	-1.14	0.25	165.53
DPCR / DPCR ^{-rel}	11.71	2.0820200193845604e-23	154.47
DPCR / DPCR $^{-rel}_{-templ}$	9.29	2.15709639836805e-16	144.41

Table 11: Student's t-test statistics for the topic vaccination.

Algorithms	Standard deviation
DPCR _{-templ}	1.38
DPCR	1.42
MORAL MAZE $^{+context}_{original}$	1.79
MORAL MAZE _{original}	1.83
$DPCR^{-rel}_{-templ}$	1.93
DPCR ^{-rel}	1.97

Table 12: Standard deviation of the coherence scores.

Algorithms	Brexit	Veganism	Vaccination
DPCR _{-templ}	1.42	1.35	1.34
DPCR	1.44	1.44	1.36
$\mathrm{DPCR}^{-rel}_{-templ}$	1.95	1.91	1.92
$DPCR^{-rel}$	1.97	1.98	1.96

Table 13: Standard deviation of the coherence scores per topic.

Algorithms	Variance
DPCR _{-templ}	0.53
DPCR	0.47
MORAL MAZE ^{+context}	0.62
MORAL MAZE _{original}	0.85
$DPCR^{-rel}_{-templ}$	0.99
DPCR ^{-rel}	0.83

Table 14: Variance of the mean coherence scores.

Algorithms	Brexit	Veganism	Vaccination
DPCR _{-templ}	0.70	0.54	0.37
DPCR	0.45	0.55	0.43
$DPCR^{-rel}_{-templ}$	0.97	0.95	1.07
DPCR ^{-rel}	0.87	0.89	0.71

Table 15: Variance of the mean coherence scores per topic.

Algorithms	Mann-Whitney U score	p-value
DPCR / MORAL MAZE _{original}	1424.5	1.4637139234284339e-09
DPCR / MORAL MAZE ^{+context}	1445.5	5.196092755480786e-07
DPCR / DPCR _{-templ}	27898	0.0051
DPCR / DPCR ^{-rel}	7118	8.28841006753334e-51
$DPCR / DPCR^{-rel}_{-templ}$	10150	1.0811689499536537e-41

Table 16: Mann-Whitney U test statistics. For measuring the Mann-Whitney U test we used the function stats.mannwhitneyu() provide by the python library SciPy.

Algorithms	Mann-Whitney U score	p-value
DPCR / DPCR_templ	3100	0.13
DPCR / DPCR ^{-rel}	832.5	1.0340781888834896e-17
DPCR / DPCR $^{-rel}_{-templ}$	1117.0	2.578864793608788e-15

Table 17: Mann-Whitney U test statistics for the topic Brexit.

Algorithms	Mann-Whitney U score	p-value
DPCR / DPCR_templ	2990.5	0.019
DPCR / DPCR ^{-rel}	736.5	6.949207593279784e-19
DPCR / DPCR $^{-rel}_{-templ}$	1145.5	4.394640292547241e-15

Table 18: Mann-Whitney U test statistics for the topic veganism.

Algorithms	Mann-Whitney U score	p-value
DPCR / DPCR_templ	3180	0.08
DPCR / DPCR ^{-rel}	766.5	2.4894342072291853e-18
DPCR / DPCR $^{-rel}_{-templ}$	1104.5	4.427842296383658e-15

Table 19: Mann-Whitney U test statistics for the topic vaccination.

716	To describe the DPCR Algorithm, we need to first define a number of lists, sets and functions:
717	• Lists and sets:
718	- Arg _{man} is the list of all the claims that make up an argument map.
719	- Sentence _{templates} = $\{IC_1, \ldots, IC_m\}$, where each IC_i (for $i = 1, \ldots, m$) is a set of sentence
720	templates for an Illocutionary Connection (IC).
721	- Final _{dialogue} is a list of locutions that make up the generated dialogue.
722	• Functions:
723	- Random is a function that takes a set as an input and returns a random element of the input set
724	as an output.
725	- GenerateLocution is a function that takes a speaker role, a claim (from an argument map) and
726	a sentence template as input and combines them into a locution that is returned as the output.
727	The function removes (if present) any word repetition between the argument map claim and the
728	Sentence template. Finally, the function adds a question mark to the output sentence when the
729	<i>ChildClaim</i> is a function that takes an argument man a propositional relation (pro or
730	- $China China projector is a function that takes an argument map, a propositional relation (proofcon) and a parent claim as input and gives as output a claim that stands in the proor con relation$
732	(of the input) to the parent claim in the argument map. If such a proposition does not exist, the
733	function gives the string 'FinishedBranch' as an output.
734	- <i>Remove</i> is a function that takes an argument map and a claim and removes that claim from the
735	argument map.
736	- Add is a function that takes a list (L) a locution (Loc) and an index (i) and adds the locution
737	Loc into L at the index i .
	Furthermore, an argumentative dialogue pattern (ADP) is a sequence of locution patterns (LP). An LP is defined in terms of the following components:
	(LP) [Speaker Role, Stance, Prop. Relations List, Illocutionary Connection, L_{ID}]
738	where:
739	• Speaker Role is one of the following: Chair, Witness, Panellist. It represents the role of the speaker.
740 741	• <i>Stance</i> is one of: Pro, Con, Neutral. It represents the stance of the speaker towards the main thesis/claim.
742	• <i>Prop. Relations List</i> is a list such that:
743	- Prop. Relations List = [NA]. In this case the Locution Pattern (LP) expresses the claim at the
744	root of the argument map, the map's main thesis. Or:
745	- Prop. Relations List = [Propositional Relation, ParentID], where Propositional Relation
746	can be one among Con, Pro, Disagreeing and Agreeing and ParentID is the parent claim
747	connected via the Propositional Relation. ¹²
748	• Illocutionary Connection (IC) is the illocutionary connection that is linked to the LP's sentence
749	$(L_{ID}).$
750	• L_{ID} is a unique identifier/label for the LP.
751	We use MainClaim to stand for the main claim/thesis of an argument map – this is the claim that sits at
752	the root of the map: it can have child claims (pro and con claims), but no parent claims. Finally, given a
753	list L, with $L[i]$ we mean the element at index i of L.
	¹² Note that 'Agreeing' and 'Disagreeing' are strictly speaking not relations between propositions. Rather they have to be understood either affirmation or denial of the proposition in question. In contrast, 'Pro' and 'Con' represent a relation between <i>two</i> propositions: one being in support or contradiction with the other.

C The DPCR Algorithm

Algorithm 1: DPCR Algorithm

```
Input: ADP, Arg<sub>map</sub>, Sentence<sub>templates</sub>
Output: Final<sub>dialogue</sub>
for LP \in ADP do
   Prop_{rel} = LP[PropRelList[0]];
   role = LP[Speaker_Role];
   SelectedTempl = Random(Sentence_{templates}[LP[IlloctionaryConnection]]);
   if Prop_{rel} = NA then
       C_x = MainClaim
   else
       Parent_{claim} = LP[PropRelList[1]];
       C_x = ChildClaim_{pro;con}(Arg_{map}, Prop_{rel}, Parent_{claim});
       if C_x = FinishedBranch then
           Final<sub>dialogue</sub> = empty list; /* if the end of a branch reached, discard the
            dialogue */
           End the algorithm;
       else
          Remove(Arg_{map}, C_x);
                                                        /* avoids sentence repetition */
       end
   end
   Loc_x = GenerateLocution(role, SelectedTempl, C_x);
   Add(Final_{dialogue}, Loc_x, x)
end
```

D Kialo maps used for generation

Kialo Terms of Service permit "crawling" and "use our export functionality to download debates for private use." (https://www.kialo.com/terms) Accordingly, we downloaded a set of debates for our experiments, but cannot redistribute the maps themselves with this paper. However, we can share the names of the specific maps that we used so other researchers can download these maps for their use in accordance with the aforementioned Terms of Service:

754

755

756

757

758

760

Brexit

1. Brexit: was it a good choice for the UK?	761
2. Should the UK remain in the EU if the only alternative is a hard Brexit?	762
3. Should the United Kingdom Remain A Member of the European Union?	763
Veganism	764
1. All humans should be vegan.	765
2. Is veganism a natural right?	766
3. Should people go vegan if they can?	767
4. The ethics of eating animals: Is eating meat wrong?	768
Vaccination	769
1. Do we need a vaccine to fight the Covid 19 pandemic?	770
2. Is Covid 19 more dangerous than regular flu viruses?	771
3. Is herd immunity for Covid 19 achievable?	772
4. It should be compulsory for those working with the elderly to take a Covid 19 vaccine.	773
5. Should Covid 19 vaccines be mandatory?	774
6. Should vaccinations be mandatory?	775

E Annotator Guidelines

776

778 779

791

803

804

809

810

811

812

813 814 Thank you for participating in this study. You are free to stop participating in the study at any time you want.

In the task, you will be presented with an argumentative dialogue. You will then be asked to carefully read it and judge it. In total, you will be presented with 15 dialogues.

Before starting the task, please read the following guidelines carefully. Do also feel free to refer back to these guidelines at any time during the annotation process. Indeed, we encourage you to read these guidelines anytime you have some doubts. The task should take you about fifteen-twenty minutes.

You will be asked to judge the coherence of the dialogue on a scale from 1 to 7 (1 being **incoherent** and 7 being **coherent**).

For this study, please try to use the following definition of coherence:

A dialogue is coherent if the following apply:

1) all sentences in the dialogue make sense by themselves and are clear (at the point in dialogue where they occur);

2) all sentences in the dialogue link together well with each other so that the dialogue is clear and sensible.

As a rule of thumb, if you believe that no sentences in the dialogue come out of the blue and the sentences in the dialogue are linked together well, then please rank the dialogue coherence as 7. Conversely, if you believe that all sentences in the dialogue come as out of the blue and the sentences in the dialogue are not linked together well, then please rank the dialogue coherence as 1. In the other cases, pick a number between 2 to 6 that you believe describes the level of coherence of that dialogue.

Please note, if the speakers (who will be labelled as Chair, Witness and Panellist) are in disagreement with each other, this does not mean that the dialogue is incoherent. Speakers can have a coherent dialogue although there is a disagreement between them. Remember, a dialogue is coherent if all its sentences are clear, make sense and go well with each other.

Judging the coherence of a dialogue is not straightforward. Many factors can impact dialogue coherence, and make a dialogue more or less coherent. We ask you to judge the coherence of a dialogue based on a seven-point scale which ranges from incoherent to coherent. Please try to be consistent with your judgements throughout the evaluation.

Finally, when judging the coherence of a dialogue please do not be influenced by whether you agree with the arguments in the dialogue. Remember, coherence is independent of what you think about the topic under discussion.

Dialogue source Example of a generated dialogue 1 Algorithm Moral maze_{original} Witness: Actually, we need to think about, how we have a more paternalist system for people like that. If you took the people around this table this evening, and you took away all of our contacts, our qualifications our great jobs and so on, we'd still have more internal resources than that guy. Panellist: Isn't the reality that in the last ten or twenty years there's been a massive transfer of wealth from the poor to the rich, and from the young to the old, and that what you're really trying to do is to justify that by blaming the poor for their position? Witness: No, I'm not trying to justify anything. Moral maze^{+context} Witness: Actually, we need to think about, how we have a more paternalist system for people like, for example, a young guy that lost his job at Tesco, mainly because he wasn't turning up to work on time, which is mainly because he was smoking a lot of spliff and he was basically very disorganised. If you took the people around this table this evening, and you took away all of our contacts, our qualifications our great jobs and so on, we'd still have more internal resources than that guy. Panellist: Isn't the reality that in the last ten or twenty years there's been a massive transfer of wealth from the poor to the rich, and from the young to the old, and that what you're really trying to do is to justify that by blaming the poor for their position? Witness: No, I'm not trying to justify anything. DPCR Witness: I think that all humans should be vegan. In the sense that a world of veganism would be a more ethical world: its morals would bring benefits to human society. Panellist: Don't you think that killing animals for food is a survival instinct, and so not inherently unethical or morally blameworthy? Witness: I don't think so, I think that instinctive, natural behavior is counterproductive can create problems, both for the individual and society, and both might want it removed. If both deem it immoral and unethical, then it is it as such and the unwanted behaviors should be shied away from and hopefully removed if possible.

F Examples of original and generated dialogues

Table 20: Examples of the original dialogues MORAL MAZE_{original} and context-enhanced original dialogues MORAL MAZE_{original}, as well as a dialogue generated with the full DPCR algorithm.

Algorithm	Example of a generated dialogue
DPCR _{-templ}	Witness: All humans should be vegan. A world of veganism would be a more
-	ethical world: its morals would bring benefits to human society.
	Panellist: Killing animals for food is a survival instinct, and so not inherently
	unethical or morally blameworthy.
	Witness: Instinctive, natural behavior is counterproductive can create problems,
	both for the individual and society, and both might want it removed. If both
	deem it immoral and unethical, then it is it as such and the unwanted behaviors
	should be shied away from and hopefully removed if possible.
DPCR ^{-rel}	Witness: I believe that Veganism is a natural right. In the sense that humans
	sit in the greatest position of control on earth, to rule it and shape it as though
	the highest power in it. Since we are considering that inalienable rights are
	endowed by natural law, we must be inferring that there is a natural preference
	for how justice is shaped. Nature, (particularly the expression of life), is most at
	peace when ruled in fairness, so it follows that natural law should direct humans
	to be benevolent. Humans are meant to be vegan.
	Panellist: Don't you think that there is no evidence proving that humans were
	created by a mindless force of evolution, and there is overwhelming evidence
	that many have found the mind of the creator can be reasonably discerned?
	Witness: I don't think that's true. I think that there is an overwhelming
	consensus in the scientific community to support the claim all life on earth is
	the result of Evolution.
DPCR ^{<math>-rel$-templ$</math>}	Witness: Veganism is a natural right. How humans are meant to behave is not
	necessarily defined by what is best for their human health.
	Panellist: An abnormal health condition can result in a risk to a person's life
	if they were to live a normal lifestyle. In context of veganism, if a person's
	digestive system has become unable to sustain life without eating meat, there is
	an unnatural conflict between the human's right to live versus the animal's right
	to live, where the vegan cannot ultimately choose to preserve the life of both.
	Witness: It should be argued that the right to live with good conscience qualifies
	the right to take one's own life.

Table 21: Examples of generated dialogues for the algorithms \mathbf{DPCR}_{-templ} , \mathbf{DPCR}^{-rel} and $\mathbf{DPCR}^{-rel}_{-templ}$. These three algorithms are ablated version of the full DPCR algorithm.