DIS2DIS: Explaining Ambiguity in Fact-Checking

Anonymous ACL submission

Abstract

Ambiguity is a linguistic tool for encoding information efficiently, yet it also causes misunderstandings and disagreements. It is particularly relevant to the domain of misinformation, as fact-checking ambiguous claims is difficult even for experts. In this paper we argue that instead of predicting a 007 veracity label for which there is genuine disagreement, it would be more beneficial to explain the ambiguity. Thus, this work introduces claim disambiguation, a novel constrained generation task, 011 for explaining ambiguous claims in fact-checking. This involves editing them to spell out an interpretation that can then be unequivocally supported by the given evidence. We collect a dataset of 1501 such claim revisions and conduct experiments with sequence-to-sequence models. The performance is 017 compared to a simple copy baseline and a Large 018 Language Model baseline. The best results are achieved by employing Minimum Bayes Decoding, with a BertScore F1 of 92.22. According to human evaluation, the model successfully disambiguates the claims 72% of the time.

1 Introduction

027

034

035

Ambiguity is a property of language that allows for utterances to have multiple possible meanings, which serves communicative purposes such as efficiency (Piantadosi et al., 2012). However, it also causes some complications. Ambiguity is not always perceived by listeners or readers (Rodd, 2018), with interpretations depending on context and motivation (Voss et al., 2008), and implicit meanings are difficult to argue with (Henderson and McCready, 2018). Recent work has also indicated that ambiguity is difficult not only for humans, but NLP models too. Liu et al. (2023) observed that Large Language Models (LLMs) are not good at detecting ambiguity in language, including very large models fine-tuned on human feedback such as GPT-4 (OpenAI, 2023).

039

040

041

043

044

045

047

050

051

053

054

057

059

060

061

062

063

064

065

066

067

068

069

070

071

073

074

075

076

078

079

Cognitive science research has shown that underspecified statements can lend themselves to misinformation due to the human cognitive predisposition to powerful inferences with little evidence (Cimpian et al., 2010). Misinformation refers to claims that are verifiably non-factual, however many claims lie in between the true/false dichotomy, due to the inherent ambiguity in language (Uscinski and Butler, 2013; Adams et al., 2023). Even expert fact-checkers often disagree on the factuality of claims, mainly in cases with ambiguous or partially true claims (Lim, 2018). Fact-checking is a particularly interesting domain for studying ambiguity, since claims are often presented for fact-checking out of context. In addition, annotation disagreement in fact-checking has been shown to be largely caused by ambiguous language (Glockner et al., 2024). To illustrate, disagreement in the top example in Table 1 stems from the vagueness of the term 'power', which could mean 'political power' or 'influence'. Under the former interpretation the claim is refuted, however under the latter it is neutral with regard to the evidence. Recent work has also shown that labels alone are not sufficiently informative to the end-users of automated fact-checking systems (Schlichtkrull et al., 2023a). Moreover, research on explainability in fact-checking provides explanations for the fact-checking labels (Kotonya and Toni, 2020; Stammbach and Ash, 2020; Krishna et al., 2022; Atanasova, 2024), however none of these works focuses on ambiguity.

In this work we generate explanations for ambiguous claims, which have been largely understudied. We propose the disambiguation of a claim as an explanation for why its factuality may be debatable, in the paradigm of elaborative simplification (Srikanth and Li, 2021), the idea that adding content can ease reasoning about the causal links in

Original claim: A Quiet Place has subtitles for the sign language. Evidence:[...] Producers Andrew Form and Bradley Fuller said that they initially planned not to provide on-screen subtitles for sign-language dialogue while providing only "context clues," but they realized that subtitles were necessary for the scene in which the deaf daughter and her hearing father argue about the modified hearing aid. [...] Revised claim: A Quiet Place has subtitles for the sign language. Original claim: Gold is the highest an album can go. R Evidence: [...] In 1975, the additional requirement of 500,000 units sold was added for Gold albums. Reflecting growth in record sales, the Platinum award was added in 1976, for albums able to sell one million units, and singles selling two million units. The Multi-Platinum award was introduced in 1984, signifying multiple Platinum levels of albums and singles. Reflecting additional growth in music sales, the Diamond award was instituted in 1999 for albums or singles selling ten million units. [...] Revised claim: Diamond is the highest an album can go. Original claim: The king of Cambodia does have power. Α Evidence: [...] Under the Constitution, the King has no political power, but as Norodom Sihanouk was revered in the country, his word often carried much influence in the government. [...] Revised claim: The king of Cambodia has no political power, but has had influence Original claim: No one died in the Tacoma Narrows Bridge collapse. U Evidence: [...] The weather system that caused the bridge collapse went on to cause the Armistice Day Blizzard that killed 145 people in the Midwest. [...] The Armistice Day storm and the strong winds that earlier had caused the Tacoma Narrows Bridge to oscillate, twist, and collapse into the waters below. [...]

Revised claim: It is not clear from the evidence whether anyone died in the Tacoma Narrows Bridge collapse.

Table 1: Examples of S(UPPORTED), R(EFUTED), A(AMBIGUOUS) and U(NSUBSTANTIATED) claims in DIS2DIS.

the text. In our context, the disambiguation makes the implicit interpretation that is supported by the evidence explicit. That is, a claim C is ambiguous, because it would only be supported by the evidence if we take the rewrite C' as its interpretation. The disambiguation is not intended to represent the intention of the speaker.

084

089

094

102

103

104

The claim and evidence pair is the input, and the unambiguously supported revised claim is the expected output. Annotator disagreement is used as signal for item ambiguity. We collect the DIS2DIS (Disagreement to Disambiguation) dataset, with annotators labeling claims as SUPPORTED, REFUTED, AMBIGUOUS or UNSUBSTANTIATED by the evidence, and then revising the claims to be supported. Multiple rounds of revisions are needed to reach consensus on a claim being supported. Sequenceto-sequence (seq2seq) models are trained on the ensuing dataset. The best results are achieved with Minimum Bayes Risk (MBR) decoding (Freitag et al., 2022) for finding the disambiguations that represent the model consensus. Our bestperforming model achieved 92.22 BertScore micro F1, and according to human evaluation, successfully disambiguates the claim 72% of the time ¹.

2 Related Work

Linguistic Phenomenon: Ambiguity Lexical am-106 biguity tasks have been successfully performed by 107 models for decades (Bunescu and Pasca, 2006; Ide 108 and Véronis, 1998; Mitkov, 2014; Ng and Cardie, 2002). While discourse information has also been 110 long and successfully applied to them (Asher and 111 Lascarides, 1995), disambiguation of entire dis-112 courses has only recently received attention. Some 113 recent work has studied the linguistic phenomena 114 that underpins ambiguity. Datasets of pragmatic 115 and discourse phenomena such as implicature have 116 been built (Nizamani et al., 2024), and models have 117 been proposed for performing such inference as a 118 task either in its own right (Pandia et al., 2021), 119 or as a by-product of other tasks such as natural 120 language inference (Jeretic et al., 2020). Other 121 work has focused on making implicit meanings ex-122 plicit. Quan et al. (2019) peform ellipsis and coref-123 erence resolution in dialogue turns, essentially dis-124 ambiguating utterances by making the omitted or 125 referred expressions explicit. Choi et al. (2021) de-126 fine the task of decontextualization, which consists 127 of rewriting sentences to be interpretable out of con-128 text. Similarly, Wu et al. (2023) generate Questions 129 under Discussion (QUDs) for sentences in dialogue 130 in order to make explicit the underlying drivers of 131 discourse, while Yu et al. (2023) edit loaded ques-132 tions in order to remove implicit or explicit presuppositions, and Min et al. (2020) disambiguate ques-134

¹https://anonymous.4open.science/r/ Dis2Dis-A6B3/

219

221

222

224

226

227

228

187

tions in open-domain question answering. Some recent work has also explored the ability of LLMs to detect ambiguity, and improved their near-random performance by instruction-tuning, showing that this task can benefit from specialised data (Ruis et al., 2024). However, to the best of our knowledge, such discourse-expounding methods have not yet been applied in the context of fact-checking.

135

136

137

138

140

141

142 Method: Text Editing Text simplification and er-143 ror correction both relate to disambiguation as they 144 use edits to clarify text. Both grammatical error 145 correction and text simplification are usually ap-146 proached with seq2seq or sequence-to-edit super-147 vised training methods (Chandrasekar et al., 1996; 148 Dahlmeier and Ng, 2012; Yuan and Briscoe, 2016; 149 Al-Thanyyan and Azmi, 2021). Most simplification models do not generate elaborative simplifications, and those that do, tend to hallucinate (Srikanth 152 and Li, 2021). Factual error correction is also ap-153 proached with seq2seq models (Cao et al., 2020), 154 distant supervision (Thorne and Vlachos, 2021), 155 and hyper-networks (Chen et al., 2023). The work in factual error correction has also replicated the 157 limited binary factuality judgment framework, and 158 159 is therefore limited to correcting **REFUTED** items to be SUPPORTED, without considering ambiguity. Domain: Fact-Checking Recent work has looked 161 into the insufficiency of the SUPPORTED, RE-FUTED and NEUTRAL label scheme. Schlichtkrull 163 et al. (2023b) add a category "conflicting ev-164 idence/cherry-picking" in order to characterise 165 cases where the evidence provides reasons to both 166 support and refute a claim. However, cherry pick-167 ing is only one particular type of ambiguity, which bears an intentional connotation. Glockner et al. 169 (2024) provide an analysis of the varied linguis-170 tic phenomena which cause disagreement over the 171 traditional ternary labels, showing a statistically sig-172 nificant correlation between various types of prag-173 matic, and discourse inference and annotator agree-174 ment over the labels. Consequently, they model 175 the fact-checking task with soft labels, predicting a 176 distribution rather than a single gold target, in order 177 to account for the difference in interpretations of 178 the ambiguous items. However, soft labels are not 179 easily interpretable.

181Aim: ExplainabilityIn the field of explainability182of fact-checking, different types of explanations183have been proposed. Using saliency maps to indi-184cate the parts of the input that are relevant to the185assigned label is the most straigtforward approach186(Atanasova et al., 2022). Atanasova (2024) use

the explanations provided by fact-checkers themselves as justification for their judgment. Similarly, Kotonya and Toni (2020) collect expert data and generate free-form explanations, including explanations for and against a given claim if the evidence is mixed between SUPPORTED and REFUTED. However, they do not separate ambiguous items from those that have conflicting evidence. If an item has conflicting evidence from different sources, the claim itself may be unambiguous. Stammbach and Ash (2020) generate summaries of the evidence with regard to the given claim as explanations, and demonstrate their utility by predicting the veracity label from the summaries. More generally in the field of explainability, some work has studied the types of explanations that are helpful to the end user (Schuff et al., 2022). Jacovi et al. (2021) generate contrastive explanations for various NLP tasks, based on the idea that the cause for choosing a particular label is dependent on which alternative that label is being contrasted with.

Data signal: Disagreement Research on various NLP tasks has shown that disagreement over labels in classification tasks, as well as diversity of outputs in generation tasks, is informative of the difficulty of items (Uma et al., 2021), beneficial for training better models (Jiang and Marneffe, 2022), and valuable in evaluation (Pavlick and Kwiatkowski, 2019). However, disagreement has not been used as signal for disambiguation.

By and large, in the current paper we address the issues that have been raised by previous work, which have not been combined into one dataset as of yet, as summarised in Table 2.

(Thorne et al., 2018) (Stammbach and Ash, 2020) (Kotonya and Toni, 2020) (Thorne and Vlachos, 2021) (Schlichtkrull et al., 2023b) (Glockner et al., 2024)	× × × × × × ×	s / R / A / U Distinction	×	Explanations	× × ✓ × × × × ×	Ambiguity Explanations
--	---------------	------------------------------	---	--------------	--------------------------	---------------------------

Table 2: Dataset Comparison

3 DIS2DIS: Disagreement to Disambiguation

3.1 Task Definition

The task of disambiguation is, given a claim and evidence, to generate a disambiguated claim that is fully supported by that evidence. The expected disambiguation is different depending on the relation between the original claim and the evidence.

Claims can be SUPPORTED, REFUTED, AMBIGU-229 OUS or UNSUBSTANTIATED by the evidence. If 230 the claim is already SUPPORTED, then no changes are required, while REFUTED claims should be negated. The AMBIGUOUS class has items that could be either supported or refuted by the evidence depending on the interpretation, such as in the third example from the top in Table 1. If the claim is ambiguous, the revision should lay out an interpretation of the original claim which is supported 238 by the evidence, such as "The king of Cambodia has no political power, but has had influence" in 240 this case. The UNSUBSTANTIATED class contains 241 items where the evidence does not answer the Ques-242 tion Under Discussion (QUD) of the claim. For 243 instance, the claim in the bottom row of Table 1 is UNSUBSTANTIATED, because while the evidence 245 mentions the blizzard casualties, it does not specify whether anyone died in the bridge collapse. That 247 is, the evidence does not answer the question "Did anyone die at the Tacoma bridge collapse?" Thus the disambiguation should state that "It is not clear from the evidence whether the claim is true".

3.2 Annotation Scheme

252

253

256

261

265

266

267

269

270

273

274

275

277

We collected a dataset for this task by using claims and evidence from the AmbiFC (Glockner et al., 2024) dataset, which reportedly had a high annotator disagreement due to ambiguity. In order to get as many ambiguous items as possible, we mostly select claims from AmbiFC with the highest entropy of labels, motivated by the relationship between label entropy and annotator certainty shown in (Glockner et al., 2024).

The annotations were collected using the Prolific platform.² The open-source annotation tool 'Potato' (Pei et al., 2022) was used to design the interface. The annotators were provided with explanations and examples of all the possible label classes and the expected respective disambiguations. The annotators are asked to select a label for the original claim, revise the claim, and highlight the parts of the input that they deem the most informative for the label they selected. The annotation guidelines are presented in Appendix A. In addition, a pre-tester question was used to ensure the annotators understood and followed the instructions. The annotators were asked to label the pre-tester item in the second row of Table 1 in order to take part in the annotation task.

The main task for the annotators was to revise the claim to be unambiguously supported by the evidence. Interestingly, many revisions for the pretester item paraphrased the following undesirable claims: "Platinum is the highest an album can go" (15%), "Multi-Platinum is the highest an album can go'' (6%). This result shows that annotators were likely to stop reading once they reached the part of the evidence that was sufficient to reject the claim, namely the mention of the Platinum award, providing an insufficiently disambiguated revision. This provided an incentive to run multiple rounds of annotations of the same item by different annotators, as it indicated that single edits often do not suffice. A revised claim from the first annotator would be passed on to a second annotator as an original claim for a classification and disambiguation. This is repeated until an annotator labels the claim as SUPPORTED, which we take to mean that the claim has been fully disambiguated. If no consensus is reached after the third revision has been evaluated, we interpret this as an impossible disambiguation, therefore assigning it to the UNSUBSTANTIATED class. The flowchart in Figure 1 illustrates the iteration process. Some items are also annotated multiple times from scratch, in order to see the variation of disambiguations and acquire multiple references for a subset of the dataset.

278

279

280

281

282

283

284

285

287

289

290

291

292

293

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

327

328

The Sankey diagram in Figure 2 illustrates the paths through different labels that claims go through until a consensus is reached. A single edit is sufficient to disambiguate about half of the claims, however the remaining items require a few iterations until different annotators assign it the same label. The figure illustrates that the AMBIGU-OUS class is particularly difficult to tease apart from UNSUBSTANTIATED, since multiple rounds of annotations are sometimes required to reach consensus on an ambiguous item.

3.3 DIS2DIS Dataset

To generate a dataset for the task of disambiguation, the original claim, any intermediate claims, the evidence and the final revised claim are put together to form an instance of the DIS2DIS dataset. If a single edit was sufficiently disambiguating, the original claim and the first edit form a (*source, target*) pair. Otherwise, in the case that the original claim is agreed on by more than one annotator as SUPPORTED, then the (*source, target*) pair is (original claim, original claim), while if the original claim is agreed on by more than one annotator as

²https://www.prolific.co/



Figure 1: Flowchart illustrating multiple rounds of annotation.

	Train Dev Test		Original Label				Mean Cla	Mean # of		
	ITum	Dev	rest	S	R	A	U	Original	Revised	Revisions
Ambi	537	64	161	71	206	403	82	12.5	18.2	1.35
All	1128	136	237	219	317	403	562	12.6	18.6	1.96

Table 3: DIS2DIS dataset statistics. Original Label corresponds to Original Label in Figure 1.



Figure 2: The labels assigned to claim revisions as they are iteratively edited.

REFUTED, then it is (original claim, "*It is not true that* "+original claim). Alternatively, if multiple edits were required, the original as well as the intermediate claims are used as the *source* claim, while the final disambiguated claim is the *target*. Finally, if three edits still did not lead to agreement on the label, the original claim is treated as the *source*, while the *target* is formulated as "*It is not clear from the evidence whether* "+original claim.

329

330

335

338

339

341

343

346

The resulting dataset contains 1501 items (see Table 3 for dataset statistics, and Appendix H for the Dataset Datasheet). The split into training, development and test sets is performed by firstly retaining all the items with multiple references for the test set, and then applying stratified sampling to ensure that disambiguations that stem from the same AmbiFC (Glockner et al., 2024) claim or evidence do not get separated into different splits, to ensure no contamination of data from the training set in evaluation. The test set contains on average 1.48 references. The dataset has an 'AMBIguous' subset for experimenting only with items that are ambiguous, which contains 762 items that take at least one and no more than three edits to reach consensus on the veracity of the claim. This is the focus part of our study, however we include the other cases in the full dataset due to the fact that the model needs to learn different behaviors depending on the initial relationship of the claim and the evidence, which is not a given.

350

351

352

353

354

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

381

383

3.4 Agreement and Evaluation Metrics

For evaluating the quality of the collected dataset, as well as selecting the best automatic metrics for training and evaluating models on the data, we perform a blind human evaluation on the generated disambiguations. Two annotators with graduate training in language sciences review a set of 27 original claims and 108 of their revised versions, labeling each as SUPPORTED, REFUTED, AMBIGUOUS or UNSUBSTANTIATED by the evidence. The instructions to the evaluators provide the same information as the original annotators to keep the annotation scheme consistent, apart from the 'unsubstantiated' label. Due to the fact that the task of the evaluators is to judge the change between the original and revised claim, when asked about the revised claim the annotators are required to choose the UNSUB-STANTIATED label if the revised claim does not address the same QUD as either the evidence or the original claim. This difference is necessary due to the fact that disambiguations which drift away from the point being made in the original claim are not truly disambiguations, even if they are factual. The annotation guidelines for the human evaluation are presented in Appendix C. For example, if the

claim in the final row of Table 1 is revised to read "145 people died in the the Armistice Day Blizzard", it no longer answers the question of whether anyone died in the bridge collapse, and therefore is not a true disambiguation of the original claim.

A heuristic combines judgments on individual claims into an overall score for the quality of the edit, as shown in Figure 3. The agreement between the two evaluators on their individual labels assigned to original and revised claims, as well as the binary score between 0 (not disambiguated or poorly disambiguated) and 1 (disambiguated), is measured with Cohen's κ . We observe substantial agreement at κ values of .66 and .69 respectively.

S A R U R A U U U S S S R U U I

Figure 3: The overall score of 1 (green arrow) for disambiguated items, 0 (red arrow) for not disambiguated or poorly disambiguated, depending on the label of the original claim (top) and the revised claim (bottom) being S(UPPORTED), A(MBIGUOUS), R(EFUTED) or U(NSUBSTANTIATED).

The overall scores of the evaluators are then compared to automated metric scores in order to select the most appropriate metric for the task. The metrics commonly used in text generation tasks such as machine translation or text simplification are tested: ROUGE (Lin, 2004), BLEU (Papineni et al., 2002), BertScore (Zhang et al., 2019), Comet (Rei et al., 2020) and SARI (Xu et al., 2016). Both neural and token-matching metrics are used, some of which support the inclusion of the source into the evaluation, which is valuable in a task such as disambiguation, where the original claim as well as the evidence text are relevant to evaluating the quality of the generated sequence. Table 4 presents the correlation scores for ALL items as well as the AMBIguous subset, using Pearson (Sedgwick, 2012) correlation coefficient. The correlation is strongest with BertScore values.

The neural metrics are better correlated with human judgments than the token-based metrics, especially for the AMBIguous subset. This is expected given that the ambiguous items have more nuanced edits which are less tied to exact token matches. It is interesting to note that BertScore Precision (p) stands out with a much higher correlation coefficient with human judgments for the AMBIguous subset compared to the rest of the items. This could be explained by the fact that for items which are already unambiguous, adding false positive tokens (words which are not in the reference), might not affect the relationship between the claim and the evidence. That is, if the claim is already supported, adding irrelevant details to make the claim more specific does not affect the ambiguity or the factuality. On the other hand, for ambiguous items, adding the exact disambiguating details is key.

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

4 Experiments

4.1 Baselines

4.1.1 Copy Baseline

A common text editing baseline is keeping the input as is. This baseline would provide correct disambiguations for the items which are already supported, and in the other cases the copied claims would superficially be very similar to the reference claims, as the disambiguations are comprised of a few token changes only. This could be expected to yield relatively high evaluation scores, especially on automatic metrics.

4.1.2 LLM Baseline

We run a zero-shot and few-shot experiments with the Llama3, 8 billion parameter model (AI@Meta, 2024), in order to evaluate the out-of-the-box LLM performance on the task of disambiguation. For the few-shot scenario we provide the model with 4 or 8 examples, covering all 4 class types (SUPPORTED, REFUTED, AMBIGUOUS, UNSUBSTANTIATED). The model is given the instruction to "please make the following claim less ambiguous with regard to the following evidence." The examples are presented in random order, separated by newlines, and punctuated with 'Claim:', 'Evidence:' and 'Revised Claim:' tags. The reader may find the inputs in Appendix B. The model output is constrained to generate a single sentence by stopping generation at a newline token, replicating the way that the few-shot examples are fed to the model. The model is expected to perform well due its large size and large training set, however we hypothesize that the task is still hard enough for the model to incur some errors, given the scarcity of direct examples of disambiguation during training and the linguistic complexity of the ambiguity relations.

4.2 Sequence-to-Sequence Model

In order to evaluate how well the collected data can serve for training seq2seq models, we finetuned a

421

422

423

3	8	4	
3	8	5	
3	8	6	
3	8	7	
3	8	8	
3	8	9	
3	9	0	
3	9	1	
3	9	2	
3	9	3	
3	9	4	
3	9	5	
3	9	6	
0	0	7	

	DIEU		RO	UGE			SARI			BertScore			Comot
	DLEU	1	2	L	L-sum	mean	keep	add	delete	F1 (micro)	р	r	Comet
Neural	ıral X X			X			✓ ✓			 Image: A start of the start of			
Source	X			X			•	/		, א	x		1
						A	MBI						
PCC	0.41	0.52	0.48	0.52	0.52	0.41	0.44	0.33	0.31	0.54	0.58	0.44	0.55
ALL													
PCC	0.44	0.50	0.47	0.50	0.50	0.45	0.45	0.42	0.33	0.52	0.52	0.46	0.50

Table 4: Pearson's Correlation Coefficient (PCC) between automatic metrics and human judgment across ALL items and for the AMBIguous subset. 'Neural' marks scores which are based on neural models, and 'Source' marks whether the score takes the source input into account.

Flan-T5 base model with 250 million parameters 473 (Chung et al., 2024). The model input is the same 474 as for the LLM baseline for consistency. Addi-475 tionally, decoding techniques are used to improve 476 model performance by guiding it to select the spe-477 cific tokens in the evidence which would help dis-478 ambiguate a claim if added to the revised version. 479 480 Length penalty, vocabulary forcing, and MBR (Freitag et al., 2022) are experimented with. The sim-481 plest of such methods is length penalty, which pe-482 nalises the model for short generations. This is 483 expected to improve results as disambiguations typ-484 ically require an addition of a modifier, conditional 485 clause or other specifying details, which make the 486 reference length longer than the source. Alterna-487 488 tively, using vocabulary forcing on this model leverages the fact that the modifying phrases needed for 489 disambiguation can be generally found in the ev-490 idence. We therefore constrain the generation to 491 include at least one of the tokens that appears in the 492 evidence but does not appear in the source claim. 493

495 496

497

498

499

501

502

506

507

509

510

511

512

513

Finally, the application of MBR to this task is inspired by the idea that disambiguation is tied to decreasing disagreement, which is reflected in the way the dataset was collected. Intrinsic uncertainty and ambiguity are related to the inadequacy of the mode sought by greedy and beam search decoding (Stahlberg et al., 2022). The MBR method generates a number of hypothesis sequences as well as pseudo references, and uses a utility function to find the best hypothesis. In our case, the best hypothesis would be the one that would reach highest agreement amongst humans, therefore we try to find the hypothesis which has the highest BertScore value when compared to the pseudo references.

The model is trained on a single NVIDIA TU102 GPU with batch size 8, for a maximum of 30 epochs, using early stopping by monitoring the BertScore metric, which has the highest correlation with human judgment. A hyperparameter search is performed (please see Appendix E for the values searched).

5 Results and Analysis

5.1 Results

Table 5 presents the best single run results of the models described in Section 4.2 after the hyperparameter search, and the baselines from Sections 4.1.1 and 4.1.2. The copy baseline predictions receive the highest scores on automatic metrics, however a careful inspection of the outputs of the models indicates that this result is not representative of the real ranking. The length of the generations also indicates a discrepancy between the appropriateness of the generation and its BertScore values. Length penalty, Vocabulary forcing and the 0-shot LLama3-8B model all overshoot the target by generating lengthy claims which are not actually helpful disambiguations. The models perform relatively on par across the different classification labels, with the largest differences between approaches seen in the 'ambiguous' class. Please see Appendix F for a breakdown of the results by class, and Appendix G for the statistical significance test results.

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

544

545

546

547

548

549

550

551

552

Table 7 presents the results of selected models on the test set, with multiple references, including a human evaluation on a random subset of 50 items.We perform a human evaluation with crowdworkers on the Prolific platform. The annotators are asked whether the revised claim is a good disambiguation of the original claim. The annotation guidelines are presented in Appendix D. As unreliability of the automatic metrics for evaluating disambiguations is corroborated by the results in Table 7 as well, which shows that the ranking order based on automatic metrics does not match the ranking order of human evaluation at all. Interestingly, the models trained on DIS2DIS perform better on the AMBIguous dataset than overall exhibiting specialised knowledge, while the LLM

-			Conv	L	lama3-8	В		Flan-T:	5-250M		
		Human	Desalina	0 shot	1 abot	9 abot	Daga	Length	Vocab	MDD	
			Dasenne	0-81101	4-81101	0-8110t	Dase	Penalty	Forcing	MDK	
Bert F1 (m	nicro)	100	94.17	91.39	93.29	94.42	93.98	92.72	92.76	93.36	А
Score p		100	95.50	91.16	94.52	95.28	94.92	92.42	93.00	93.88	ź
len (tokens)		17.33	12.5	86.53	14.53	15.95	15.86	56.08	19.89	16.30	BI
Bert F1 (m	nicro)	100	94.38	90.91	94.39	94.25	94.13	94.18	93.59	94.81	~
Score p		100	95.99	91.23	95.79	95.27	93.66	93.52	93.13	94.87	AL
len (tokens)		18.43	12.6	86.01	15.82	21.46	21.61	24.39	45.13	18.74	Г

Table 5: Model performance and baseline scores on the development set, for the AMBIguous subset and ALL items.

Origin	al claim: You can na	me your kid anything in America.							
Evider	Evidence: [] Traditionally, the right to name one's child or oneself as one chooses has been upheld by court rulings and is								
rooted	rooted in the fourteenth Amendment and the First Amendment, but a few restrictions do exist. Several states limit the number								
of char	acters that can be use	ed. A few states ban the use of obscenity. Restrictions vary by state, Kentucky for instance, has no							
naming	g laws whatsoever.[
	Llama3-8B 8-shot:	You can name your kid anything in America.							
ise	Flan-T5 MBR:	You can name your kid anything in America, but restrictions exist.							
lai	Human:	You can name your kid anything in Kentucky, while other states have some restrictions on length							
R 0		or obscenity.							

Table 6: Example target, baseline and model outputs.

baseline shows the reverse. The annotator agreement is 0.56, measured with Cohen's κ .

	Сору	Llama3-8B	Flar	Human					
	Baseline	8-shot Base		MBR	пишап				
	Амві								
Bert F1	93.85	93.12	93.51	92.22	97.57				
Score p	95.49	93.47	94.35	92.36	97.51				
Human	0.10	0.60	0.62	0.74	0.85				
		All							
Bert F1	94.10	93.30	93.07	93.05	98.05				
Score p	95.65	94.01	93.19	93.24	98.03				
Human	0.27	0.67	0.60	0.72	0.82				

Table 7: Model performance and baseline scores with human evaluation on ALL items in the test set, and its AMBIguous subset.

5.2 Analysis

Table 6 presents different generations to the same 556 original claim containing underspecification. The LLama3-8B model baseline fails to disambiguate the claim, leaving it as is. Human and MBR model generations both provide suitable disambiguations, where the nuance of restrictions to naming are mentioned in both, with the human generation providing a more detailed explanation. This represents the general tendency observed in qualitative analysis, with MBR model generations providing better 565 disambiguations than other models and baselines, 566 however not reaching the full potential of human revisions. The elaborations for disambiguating un-568 derspecified claims take the form of relative clauses and subordinate clauses (e.g. conditional or contrastive). On the other hand, hyponyms or shorter

modifiers such as adjectives can be sufficient to disambiguate a vague claim, such as the one in row 3 of Table 1.

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

Based on a qualitative analysis, the most common errors for all models include a) not changing the claim at all when a revision is required, b) mixing up the types of edits needed for the 'unsubstantiated' and the 'ambiguous' classes, c) hallucinating details not present in the evidence, d) missing or superfluous negation. The Base model exhibits the highest number of a), b) and d) type errors, while the Llama3-8B baseline suffers the most from c).

6 **Discussion and Conclusion**

The results of this study provide evidence that ambiguity is difficult to detect and remove for humans as well as language models. We argue that the fact that humans find detecting ambiguity and disambiguation difficult, calls for work on disambiguation. Apart from the application to explainability in fact-checking, disambiguation could also be applied to assisting in writing less ambiguously, or providing less ambiguous summaries.

Future research may involve experimenting with multi-step disambiguation as well as exploring the utility of highlighted inputs for model training. Future directions could also include exploring the link between disagreement and ambiguity by directly using disagreement as feedback for disambiguations through reinforcement learning strategies.

553 554

567

571

Limitations

601

616

617

618

621

622

624

628

631

633 634

636

639

644

645

647

Our approach is limited in handling certain types of ambiguity, namely the ones which are prominent in the fact-checking data we used: underspec-604 ification, vagueness, implicature, presupposition, probabilistic enrichment, coreference. This may 606 not cover other types of ambiguity that could be more common in different domains. In addition, we only focused on English due to dataset availability. Our work was limited to claims with evi-610 dence from Wikipedia, however fact-checking and 611 ambiguity are pervasive in various platforms of 612 communication. This study is also limited to the 613 fact-verification step of fact-checking, studying the 614 impact of ambiguity when the evidence is given. 615

> The study shows that automatic evaluation metrics are not reliable in evaluating the performance of different methods of disambiguation. As a result, a human evaluation is required, which is labourintensive and time-consuming. In addition, the LLM baseline performance depends on the prompts used.

We recognise the potential risk that a disambiguation dataset, when misused, could be used to obscure rather than clarify claims, which could contribute to the spread of misinformation. We believe, however, that the benefits of learning about misinformation and ambiguity detection outweigh the drawbacks.

630 Ethics Statement

The annotators in this study were selected on the basis of residing in the UK and being native English language speakers, and were paid an hourly rate above the minimum wage in the UK (\pounds 11.44), averaging at \pounds 13.28. The annotation protocol was approved by an ethics review board. The annotation instructions contained a disclaimer that the topics appearing in the claims in the study would contain content comparable to what one might encounter while browsing the internet, as the claims are sourced from common search engine queries (Clark et al., 2019). No personal information of the annotators was collected.

Acknowledgments

References

Zoë Adams, Magda Osman, Christos Bechlivanidis, and Björn Meder. 2023. (why) is misinformation a problem? *Perspectives on Psychological Science*, 18(6):1436–1463.

AI@Meta. 2024. Llama 3 model card.	650
Suha S Al-Thanyyan and Aqil M Azmi. 2021. Auto- mated text simplification: a survey. <i>ACM Computing</i> <i>Surveys</i> (<i>CSUR</i>), 54(2):1–36.	651 652 653
Nicholas Asher and Alex Lascarides. 1995. Lexical	654
disambiguation in a discourse context. <i>Journal of Semantics</i> , 12(1):69–69.	655 656
Pepa Atanasova. 2024. Generating fact checking ex-	657
ods for Complex Reasoning over Text, pages 83–103. Springer.	658 659 660
Pepa Atanasova, Jakob Grue Simonsen, Christina Li- oma, and Isabelle Augenstein. 2022. Diagnostics- guided auglengtion generation. In Proceedings of the	661 662
AAAI Conference on Artificial Intelligence, volume 36, 10, pages 10445–10453.	664 665
Razvan Bunescu and Marius Pasca. 2006. Using ency-	666
clopedic knowledge for named entity disambiguation.	667
In 11th conference of the European Chapter of the	668
Association for Computational Linguistics, pages 9– 16.	669 670
Meng Cao, Yue Dong, Jiapeng Wu, and Jackie Chi Kit	671
Cheung. 2020. Factual error correction for abstrac-	672
tive summarization models. In <i>Proceedings of the</i>	673
Language Processing (EMNLP), pages 6251–6258.	674 675
Raman Chandrasekar, Christine Doran, and Srinivas	676
Bangalore. 1996. Motivations and methods for text	677
simplification. In COLING 1996 Volume 2: The 16th	678
International Conference on Computational Linguis- tics.	679 680
Jiangjie Chen, Rui Xu, Wenxuan Zeng, Changzhi Sun,	681
Lei Li, and Yanghua Xiao. 2023. Converge to	682
the truth: Factual error correction via iterative con-	683
strained editing. In Proceedings of the AAAI Confer-	684
ence on Artificial Intelligence, volume 37, 11, pages	685
12616–12625.	686
Eunsol Choi, Jennimaria Palomaki, Matthew Lamm,	687
Tom Kwiatkowski, Dipanjan Das, and Michael	688
Collins. 2021. Decontextualization: Making sen-	689
for Computational Linguistics, 9:447–461.	690
Hyung Won Chung Le Hou Shavne Longpre Barret	602
Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi	693
Wang, Mostafa Dehghani, Siddhartha Brahma, et al.	694
2024. Scaling instruction-finetuned language models.	695
Journal of Machine Learning Research, 25(70):1–53.	696
Andrei Cimpian, Amanda C. Brandone, and Susan A.	697
Gelman. 2010. Generic statements require little evi-	698
dence for acceptance but have powerful implications. <i>Cognitive science</i> , 34 8:1452–1482.	699 700
Christopher Clark, Kenton Lee, Ming-Wei Chang	701
Tom Kwiatkowski, Michael Collins, and Kristina	702
Toutanova. 2019. BoolQ: Exploring the surprising	703

816

difficulty of natural yes/no questions. In *Proceedings* of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 2924–2936, Minneapolis, Minnesota. Association for Computational Linguistics.

704

705

710

711

712

713

714

715

716

718

719

720

721

722

723

725

726

727

730

732

733

734

735

736

737

739

740

741

743

744

745

747

748

750

751

752

753

754

755

756

757

- Daniel Dahlmeier and Hwee Tou Ng. 2012. Better evaluation for grammatical error correction. In *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 568–572.
- Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022. High quality rather than high model probability: Minimum bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.
- Max Glockner, Ieva Staliūnaitė, James Thorne, Gisela Vallejo, Andreas Vlachos, and Iryna Gurevych. 2024.
 Ambifc: Fact-checking ambiguous claims with evidence. *Transactions of the Association for Computational Linguistics*, 12:1–18.
- Robert Henderson and Elin McCready. 2018. How dogwhistles work. In New Frontiers in Artificial Intelligence: JSAI-isAI Workshops, JURISIN, SKL, AI-Biz, LENLS, AAA, SCIDOCA, kNeXI, Tsukuba, Tokyo, November 13-15, 2017, Revised Selected Papers 9, pages 231–240. Springer.
- Nancy Ide and Jean Véronis. 1998. Introduction to the special issue on word sense disambiguation: the state of the art. *Computational linguistics*, 24(1):1–40.
- Alon Jacovi, Swabha Swayamdipta, Shauli Ravfogel, Yanai Elazar, Yejin Choi, and Yoav Goldberg. 2021.
 Contrastive explanations for model interpretability. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 1597–1611.
- Paloma Jeretic, Alex Warstadt, Suvrat Bhooshan, and Adina Williams. 2020. Are natural language inference models imppressive? learning implicature and presupposition. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 8690–8705.
- Nan-Jiang Jiang and Marie-Catherine de Marneffe. 2022. Investigating Reasons for Disagreement in Natural Language Inference. *Transactions of the Association for Computational Linguistics*, 10:1357– 1374.
- Neema Kotonya and Francesca Toni. 2020. Explainable automated fact-checking for public health claims. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 7740–7754.
- Amrith Krishna, Sebastian Riedel, and Andreas Vlachos. 2022. Proofver: Natural logic theorem proving for fact verification. *Transactions of the Association for Computational Linguistics*, 10:1013–1030.

- Chloe Lim. 2018. Checking how fact-checkers check. *Research & Politics*, 5(3):2053168018786848.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Alisa Liu, Zhaofeng Wu, Julian Michael, Alane Suhr, Peter West, Alexander Koller, Swabha Swayamdipta, Noah A Smith, and Yejin Choi. 2023. We're afraid language models aren't modeling ambiguity. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 790–807.
- Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. AmbigQA: Answering ambiguous open-domain questions. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 5783– 5797, Online. Association for Computational Linguistics.

Ruslan Mitkov. 2014. Anaphora resolution. Routledge.

- Vincent Ng and Claire Cardie. 2002. Improving machine learning approaches to coreference resolution. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 104– 111.
- Rashid Nizamani, Sebastian Schuster, and Vera Demberg. 2024. Siga: A naturalistic nli dataset of english scalar implicatures with gradable adjectives. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 14784– 14795.
- OpenAI. 2023. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.
- Lalchand Pandia, Yan Cong, and Allyson Ettinger. 2021. Pragmatic competence of pre-trained language models through the lens of discourse connectives. In *Proceedings of the 25th Conference on Computational Natural Language Learning*, pages 367–379.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Ellie Pavlick and Tom Kwiatkowski. 2019. Inherent Disagreements in Human Textual Inferences. *Transactions of the Association for Computational Linguistics*, 7:677–694.
- Jiaxin Pei, Aparna Ananthasubramaniam, Xingyao Wang, Naitian Zhou, Apostolos Dedeloudis, Jackson Sargent, and David Jurgens. 2022. Potato: The portable text annotation tool. In *Proceedings of the* 2022 Conference on Empirical Methods in Natural Language Processing: System Demonstrations.

924

925

926

870

Steven T. Piantadosi, Harry Tily, and Edward Gibson. 2012. The communicative function of ambiguity in language. *Cognition*, 122(3):280–291.

817

818

819

820

822

824

825

826

828

829

831

832

834

838

839

840

842

844

845

848

852

862

864

- Jun Quan, Deyi Xiong, Bonnie Webber, and Changjian Hu. 2019. Gecor: An end-to-end generative ellipsis and co-reference resolution model for task-oriented dialogue. 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 4547–4557.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Comet: A neural framework for mt evaluation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 2685–2702.
 - Jennifer Rodd. 2018. Lexical ambiguity. Oxford handbook of psycholinguistics, pages 120–144.
- Laura Ruis, Akbir Khan, Stella Biderman, Sara Hooker, Tim Rocktäschel, and Edward Grefenstette. 2024. The goldilocks of pragmatic understanding: Finetuning strategy matters for implicature resolution by Ilms. Advances in Neural Information Processing Systems, 36.
- Michael Schlichtkrull, Nedjma Ousidhoum, and Andreas Vlachos. 2023a. The intended uses of automated fact-checking artefacts: Why, how and who.
 In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 8618–8642.
- Michael Sejr Schlichtkrull, Zhijiang Guo, and Andreas Vlachos. 2023b. AVeritec: A dataset for real-world claim verification with evidence from the web. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track.*
- Hendrik Schuff, Alon Jacovi, Heike Adel, Yoav Goldberg, and Ngoc Thang Vu. 2022. Human interpretation of saliency-based explanation over text. In 2022 ACM Conference on Fairness, Accountability, and Transparency, pages 611–636.
- Philip Sedgwick. 2012. Pearson's correlation coefficient. *Bmj*, 345.
- Neha Srikanth and Junyi Jessy Li. 2021. Elaborative simplification: Content addition and explanation generation in text simplification. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP* 2021, pages 5123–5137.
- Felix Stahlberg, Ilia Kulikov, and Shankar Kumar. 2022. Uncertainty determines the adequacy of the mode and the tractability of decoding in sequence-to-sequence models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (*Volume 1: Long Papers*), pages 8634–8645, Dublin, Ireland. Association for Computational Linguistics.

- Dominik Stammbach and Elliott Ash. 2020. e-fever: Explanations and summaries for automated fact checking. *Proceedings of the 2020 Truth and Trust Online (TTO 2020)*, pages 32–43.
- James Thorne and Andreas Vlachos. 2021. Evidencebased factual error correction. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 3298–3309.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.
- Alexandra N Uma, Tommaso Fornaciari, Dirk Hovy, Silviu Paun, Barbara Plank, and Massimo Poesio. 2021. Learning from disagreement: A survey. *Journal of Artificial Intelligence Research*, 72:1385–1470.
- Joseph E Uscinski and Ryden W Butler. 2013. The epistemology of fact checking. *Critical Review*, 25(2):162–180.
- Andreas Voss, Klaus Rothermund, and Jochen Brandtstädter. 2008. Interpreting ambiguous stimuli: Separating perceptual and judgmental biases. *Journal of Experimental Social Psychology*, 44(4):1048–1056.
- Yating Wu, Ritika Mangla, Greg Durrett, and Junyi Jessy Li. 2023. Qudeval: The evaluation of questions under discussion discourse parsing. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5344– 5363.
- Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, and Chris Callison-Burch. 2016. Optimizing statistical machine translation for text simplification. *Transactions of the Association for Computational Linguistics*, 4:401–415.
- Xinyan Yu, Sewon Min, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Crepe: Open-domain question answering with false presuppositions. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 10457–10480.
- Zheng Yuan and Ted Briscoe. 2016. Grammatical error correction using neural machine translation. In Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 380–386.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*.

977

978

A Annotation Guidelines for Data Collection

A.1 Instructions

927

928

929

930

931

932

933

934

937

941

942

943

Procedures In this study, you will be presented with pairs of claims and evidence, and your task will be to (1) label the the given claim as 'supported', 'refuted', 'ambiguous' or 'unsubstantiated' with regard to the provided evidence, (2) highlight parts of the text that support your choice, and (3) edit the claim to make the claim match the supported label better. You should make your decisions based on the information provided more than on your world knowledge. You can expect a larger portion of the provided items to be ambiguous, so please read carefully.

Risks The risks and discomfort associated with participation in this study are no greater than those ordinarily encountered in daily life, such as when surfing the internet.

Benefits There may be no personal benefit from
your participation in the study but the knowledge
gained may have academic or industrial value.

Confidentiality By participating in this research, 949 you understand and agree that the researcher may be required to disclose your consent form, data, 952 and other personally identifiable information as required by law, regulation, subpoena, or court order. Otherwise, your confidentiality will be maintained 955 in the following manner: To protect your identity, the researchers will take the following steps: (1) Each participant will be assigned a number; (2) The 957 researchers will record any data collected during the study by number, not by name; (3) Any original recordings or data files will be stored in a secured location accessed only by authorized researchers. 961

Voluntary Participation Your participation in this
research is voluntary. You may discontinue participation at any time during the research activity.

Navigating You can use the right arrow to move 965 forward, but you are not allowed to go backward. 966 For highlighting text, multiple selections are allowed and encouraged, however they all have to 968 correspond to the same veracity label that you have 969 selected. If you wish to remove highlights that you 970 have made, you can do so by clicking on them. The 972 edits you are asked to do should be as minimal as possible, however they should not simply negate 973 the original claim. The idea is to specify the claim 974 more or change some details in the claim, to more closely match the 'supported' label. 976

A.2 Annotation Scheme

Here are the explanations and examples of the 'supported', 'refuted', 'unsubstantiated' and 'ambiguous' labels, please read them carefully before moving forward.

Supported A claim is supported by its evidence if the evidence is sufficient to draw the conclusion that the claim is true. For instance, the evidence "The NATO summit will be hosted in Vilnius, Lithuania to discuss Ukraine" supports the claim "The NATO summit will be held in Eastern Europe ".

When highlighting the important parts of the input, you could emphasize the georgraphical references to the country and the larger region it belongs to.

In order to match the 'supported' label even better, the claim could be edited to read "The NATO summit will be held in <u>Vilnius</u>, which is in Eastern Europe" in order to remove any uncertainty about the georgraphical classification of the country. The edited claim is now even more supported by the evidence, because it clarifies the location of Vilnius for the readers who may not be aware of it.

Refuted In contrast, a claim is refuted by if the evidence is its evidence sufficient to draw the conclusion that the claim is false. For instance. the claim "Ukraine has a timeline for joining NATO " is refuted by the evidence stating that "Ukraine will not be offered timeline for NATO membership at the summit in July ".

When highlighting the relevant parts of the claim and evidence, you may want to consider what makes the claim and evidence contrast, such as the different time references and the negation.

The claim can be edited to match the supported label as such: "Ukraine's has a timeline for joining NATO has not been determined yet."

Unsubstantiated Alternatively, if the claim is neither supported nor refuted by the evidence, the evidence may not provide enough information to draw a conclusion. For instance, the evidence "The NATO summit will be hosted in Vilnius, Lithuania to discuss Ukraine" is not enough to determine the veracity of the claim "Ukraine has a timeline for joining NATO". While the claim and the evidence discuss the same

1074

topic, the evidence here does not provide any answer as to whether the claim is true or false, therefore it should be marked as unsubstantiated.

In the case of an unsubstantiated claim, it would be good to highlight the parts of the input that refer to different aspects of the topic, such as the location of the summit in the evidence, and the NATO membership timeline in the claim.

The claim could be rewritten to be supported by stating that "<u>The status of</u> Ukraine's has a timeline for joining NATO is not clear from the evidence"

Ambiguous

In contrast, there is an ambiguous relationship between the claim that "Ukraine's application to join NATO is being supported"

and the evidence that says "France resolves to support Ukraine's NATO membership bid". The claim is partially true, as the application is supported by some countries, but it is not known whether it is supported by everyone. The generic statement in the claim is too broad.

In order to show the source of ambiguity, you should highlight the parts of the input that make the claim vaguer than the evidence, such as specifically naming France in the evidence in this case.

In order to match the claim to the 'supported' label, the claim could be rewritten as "Ukraine's application to join NATO is being supported by France.", as this removes the ambiguity from the original claim by specifying the country.

A.3 Examples

The examples shown to the annotators are shown in Figures 4 and 5.

B Instruction and Examples for LLM Baseline

B.1 0-shot Instruction and Input Format

"Please make the following claim less ambiguous with regard to the following evidence. Claim: [CLAIM], Evidence: [EVIDENCE], Revised Claim:"

B.2 4-shot Instruction, Examples and Input Format

"Please make the following claim less ambiguous with regard to the following evidence, as in the examples below. Claim: bridges of madison county is a true story.

1075

1076

1077

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1125

1126

Evidence: The Bridges of Madison County (also published as Love in Black and White) is a 1992 best-selling romance novella by American writer Robert James Waller that tells the story of a married Italian-American woman (WW2 War bride) living on a Madison County, Iowa, farm in the 1960s. While her husband and children are away at the State Fair, she engages in an affair with a National Geographic photographer from Bellingham, Washington, who is visiting Madison County to create a photographic essay on the covered bridges in the area. The novel is presented as a novelization of a true story, but it is in fact entirely fictional. The novel is one of the bestselling books of the 20th century, with 60 million copies sold world-wide. It has also been adapted into a feature film in 1995 and a musical in 2013.

Revised claim: bridges of madison county is a fictional story

Claim: you can keep a gray wolf as a pet.

Evidence: Some wildlife centers housing cap-1098 tive wolves prohibit handlers from entering wolf 1099 enclosures if they happen to have a cold or other 1100 vulnerability which the wolves can detect. Captive 1101 wolves are generally shy and avoid eye contact with 1102 humans other than their owner, as well as not listen-1103 ing to any commands made by any other humans. 1104 They usually vacate rooms or hide when a new 1105 person enters the establishment. Even seemingly 1106 friendly wolves need to be treated with caution, 1107 as captive wolves tend to view and treat people 1108 as other wolves, and will thus bite or dominate 1109 people in the same situation in which they would 1110 other wolves. Ordinary pet food is inadequate, as 1111 an adult wolf needs 12.5 kg (25 lbs) of meat daily 1112 along with bones, skin and fur to meet its nutri-1113 tional requirements. Wolves may defend their food 1114 against people, and react violently to people trying 1115 to remove it. The exercise needs of a wolf exceed 1116 the average dog's demand. Because of this, captive 1117 wolves typically do not cope well in urban areas. 1118 Due to their talent at observational learning, adult 1119 captive wolves can quickly work out how to escape 1120 confinement, and require constant reinforcement by 1121 caretakers or owners, which makes raising wolves 1122 difficult for people who raise their pets in an even, 1123 rather than subordinate, environment. 1124

Revised claim: it is difficult to raise a wolf as a pet.



Figure 4: An example annotation with an 'Unsubstantiated' label.



Figure 5: An example annotation with an 'Ambiguous' label.

Claim: it is illegal to flash your headlights to warn off the police in the uk.

1127

1128

1129

1131

1133

1137

1138

1140

1141

1143

1145

1147

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

Evidence: Though not all of its rules represent 1130 law, the Highway Code states "Only flash your headlights to let other road users know that you are 1132 there. Do not flash your headlights in an attempt to intimidate other road users". Drivers warning 1134 others about speed traps have been fined in the past 1135 for "misuse of headlights". Headlight flashing in 1136 the United Kingdom is often used as a signal that the driver flashing you is offering to let you go first. Such use is however strongly discouraged because 1139 it can lead to accidents where the driver flashing has not seen the approach of another road user. Using it to indicate that you are coming through 1142 and the other driver must wait, could lead to an accident. Drivers should also be aware of the so-1144 called "Flash-for-Cash" scam, in which criminals flash their lights to let other drivers out of a junc-1146 tion, then crash into them on purpose in order to make fraudulent insurance claims for damage and 1148 whiplash injury. 1149

> Revised claim: In the UK, you should only flash your headlights to let other drivers know you are there.

> Claim: you do need intent to commit a crime Evidence: In criminal law, intent is a subjective state of mind that must accompany the acts of certain crimes to constitute a violation. A more formal, generally synonymous legal term is scienter: intent or knowledge of wrongdoing. Revised claim: you do need intent to commit some crimes

[CLAIM], Evidence: Claim: [EVIDENCE], Revised Claim:"

B.3 8-shot Instruction, Examples and Input Format

"Please make the following claim less ambiguous with regard to the following evidence, as in the examples below.

Claim: bridges of madison county is a true story.

1173 Evidence: The Bridges of Madison County (also published as Love in Black and White) is a 1992 1174 best-selling romance novella by American writer 1175 Robert James Waller that tells the story of a married 1176 Italian-American woman (WW2 War bride) living 1177

on a Madison County, Iowa, farm in the 1960s. While her husband and children are away at the State Fair, she engages in an affair with a National Geographic photographer from Bellingham, Washington, who is visiting Madison County to create a photographic essay on the covered bridges in the area. The novel is presented as a novelization of a true story, but it is in fact entirely fictional. The novel is one of the bestselling books of the 20th century, with 60 million copies sold world-wide. It has also been adapted into a feature film in 1995 and a musical in 2013.

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

1203

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

Revised claim: bridges of madison county is a fictional story

Claim: you can keep a gray wolf as a pet.

Evidence: Some wildlife centers housing captive wolves prohibit handlers from entering wolf enclosures if they happen to have a cold or other vulnerability which the wolves can detect. Captive wolves are generally shy and avoid eye contact with humans other than their owner, as well as not listening to any commands made by any other humans. They usually vacate rooms or hide when a new person enters the establishment. Even seemingly friendly wolves need to be treated with caution, as captive wolves tend to view and treat people as other wolves, and will thus bite or dominate people in the same situation in which they would other wolves. Ordinary pet food is inadequate, as an adult wolf needs 12.5 kg (25 lbs) of meat daily along with bones, skin and fur to meet its nutritional requirements. Wolves may defend their food against people, and react violently to people trying to remove it. The exercise needs of a wolf exceed the average dog's demand. Because of this, captive wolves typically do not cope well in urban areas. Due to their talent at observational learning, adult captive wolves can quickly work out how to escape confinement, and require constant reinforcement by caretakers or owners, which makes raising wolves difficult for people who raise their pets in an even, rather than subordinate, environment.

Revised claim: it is difficult to raise a wolf as a pet.

Claim: it is illegal to flash your headlights to warn off the police in the uk.

Evidence: Though not all of its rules represent law, the Highway Code states "Only flash your headlights to let other road users know that you are there. Do not flash your headlights in an attempt

to intimidate other road users". Drivers warning 1230 others about speed traps have been fined in the past 1231 for "misuse of headlights". Headlight flashing in 1232 the United Kingdom is often used as a signal that 1233 the driver flashing you is offering to let you go first. Such use is however strongly discouraged because 1235 it can lead to accidents where the driver flashing 1236 has not seen the approach of another road user. 1237 Using it to indicate that you are coming through 1238 and the other driver must wait, could lead to an 1239 accident. Drivers should also be aware of the so-1240 called "Flash-for-Cash" scam, in which criminals 1241 flash their lights to let other drivers out of a junc-1242 tion, then crash into them on purpose in order to 1243 make fraudulent insurance claims for damage and 1244 whiplash injury. 1245

> Revised claim: In the UK, you should only flash your headlights to let other drivers know you are there.

Claim: you do need intent to commit a crime

1247

1248

1249

1250

1252

1253

1254

1255

1256

1257

1258

1259

1260

1261

1262

1264

1265

1266

1267

1268

1270

1271

1272

1273

1274

1277

1278

1279

1281

Evidence: In criminal law, intent is a subjective state of mind that must accompany the acts of certain crimes to constitute a violation. A more formal, generally synonymous legal term is scienter: intent or knowledge of wrongdoing.

Revised claim: you do need intent to commit some crimes

Claim: running with scissors is based on a true story.

Evidence: In 2005, the family of Dr. Rodolph H. Turcotte (191920132000), of Massachusetts filed suit against Burroughs and his publisher, alleging defamation of character and invasion of privacy. They stated that they were the basis for the Finch family portrayed in the book but that Burroughs had fabricated or exaggerated various descriptions of their activities. It's still a memoir, it's marketed as a memoir, they've agreed one hundred percent that it is a memoir. The case was later settled with Sony Pictures Entertainment in October 2006, prior to the release of the film adaptation. Burroughs and his publisher, St. Martin's Press, settled with the Turcotte family in August 2007. The Turcottes were reportedly seeking damages of \$2 million for invasion of privacy, defamation, and emotional distress; the Turcottes alleged Running with Scissors was largely fictional and written in a sensational manner. Burroughs defended his work as "entirely accurate", but agreed to call the work a "book" (instead of a "memoir") in the author's note, to alter

the acknowledgments page in future editions to rec-1282 ognize the Turcotte family's conflicting memories 1283 of described events, and express regret for "any un-1284 intentional harm" to the Turcotte family. Burroughs 1285 felt vindicated by the settlement. "I'm not at all 1286 sorry that I wrote [the book]. And you know, the 1287 suit settled - it settled in my favor. I didn't change 1288 a word of the memoir, not one word of it. It's 1289 still a memoir, it's marketed as a memoir, they've 1290 agreed one hundred percent that it is a memoir". 1291 Future printings of Running with Scissors will con-1292 tain modified language in the Author's Note and 1293 Acknowledgments pages. Where the Acknowledg-1294 ments page had read: "Additionally, I would like 1295 to thank each and every member of a certain Mas-1296 sachusetts family for taking me into their home 1297 and accepting me as one of their own," the follow-1298 ing was substituted: "Additionally, I would like 1299 to thank the real-life members of the family por-1300 trayed in this book for taking me into their home 1301 and accepting me as one of their own. I recognize 1302 that their memories of the events described in this 1303 book are different than my own. They are each 1304 fine, decent, and hard-working people. The book 1305 was not intended to hurt the family. Both my pub-1306 lisher and I regret any unintentional harm resulting 1307 from the publishing and marketing of Running with 1308 Scissors" 1309

Revised claim: running with scissors is somewhat based on the recollections of part of the author's life 1310

1311

1312

1313

1314

1315

Claim: you can drink at any age in wisconsin.

Evidence: The drinking age in Wisconsin is 21. 1316 Those under the legal drinking age may be served, 1317 possess, or consume alcohol if they are with a par-1318 ent, legal guardian, or spouse who is of legal drink-1319 ing age. Those age 18 to 20 may also possess (but 1320 not consume) alcohol as part of their employment. 1321 In the early 70s the sale of alcohol was reduced 1322 to the age of 18. The 1983 Wisconsin Act 74, 1323 effective July 1, 1984, created a drinking age of 19. Meeting in special session at the call of the 1325 governor, the legislature enacted 1985 Wisconsin 1326 Act 337, which raised the drinking age to 21 and 1327 brought the state into compliance with the NMDA 1328 (National Minimum Drinking Age) on September 1329 1, 1986. The NMDA law was amended to permit 1330 an exception for those persons who were between 1331 ages 18 and 21 on the effective date of the law. Wis-1332 consin 19- and 20-year-olds were grandfathered in 1333

1383

1385

by this exception after enactment of Act 337. In effect, the state did not have a uniform age of 21 until September 1, 1988.

Revised claim: you can drink at any age in wisconsin with someone who is of legal drinking age.

Claim: it is normal for your second toe to be longer.

Evidence: Morton's toe is the condition of having a first metatarsal which is short in relation to the second metatarsal (see diagram). It is a type of brachymetatarsia. The distal metatarsal bones vary in relative length compared to the proximal. For most feet, a smooth curve can be traced through the joints at the bases of the toes (the metatarsalphalangeal, or MTP, joints). But in Morton's foot, the line has to bend more sharply to go through the base of the big toe, as shown in the diagram. This is because the first metatarsal, behind the big toe, is short compared to the second metatarsal, next to it. The longer second metatarsal puts the MTP joint at the base of the second toe further forward. If the big toe and the second toe are the same length (as measured from the MTP joint to the tip, including only the toe bones or phalanges), then the second toe will protrude farther than the big toe, as shown in the photo.

Revised claim: your second toe can be longer than your big toe.

Claim: baby sign language is the same as regular sign language.

Evidence: Baby sign involves enhanced gestures and altered signs that infants are taught in conjunction with spoken words with the intention of creating richer parent-child communication. The main reason that parents use baby sign is with hope that it will reduce the frustration involved in trying to interpret their pre-verbal child's needs. It can be considered a useful method of communication in the early developmental stages, since speech production follows children's ability to express themselves through bodily movement. Baby sign is distinct from sign language. Baby sign is used by hearing parents with hearing children to improve communication. Sign languages, including ASL, BSL, ISL and others, are natural languages, typically used in the Deaf community. Sign languages maintain their own grammar, and sentence structure. Because sign languages are as complex to learn as any spoken language, simplified signs are often used with infants in baby sign. Teaching baby signs allows for greater flexibility in the form of sign and does not require the parent to learn the grammar of a sign language. Baby signs are usually gestures or signs taken from the sign language community and modified to make them easier for an infant to form.

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

1404

1405

1406

1407

1408

1409

1410

1411

1412

1413

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

Revised claim: baby sign language is distinct from regular sign language.

Claim: [CLAIM], Evidence: [EVIDENCE], Revised Claim:"

C Annotation Guidelines for the Initial Human Evaluation

The following are examples of claims which might be ambiguous with regard to the given evidence (original claims). Attempts have been made to disambiguate the claims by editing them to be more supported by the evidence (edited claims). Please read the original claims, the evidence, and the edited claims, and assess whether the original claim is a) supported by the evidence, b) refuted by the evidence, c) ambiguous with regard to the evidence, or d) the evidence does not address the question that is implied in the original claim. Then, assess the revised claim with regard to the evidence and determine whether it is a) supported by the evidence, b) refuted by the evidence, c) ambiguous with regard to the evidence, or d) the revised claim does not address the same question as the original claim or the evidence.

For example, the claim "you can keep a gray wolf as a pet" is ambiguous with regard to the following evidence: "[...] Captive wolves are generally shy and avoid eye contact with humans other than their owner, as well as not listening to any commands made by any other humans. [...] Ordinary pet food is inadequate, as an adult wolf needs 12.5 kg (25 lbs) of meat daily along with bones, skin and fur to meet its nutritional requirements. [...] The exercise needs of a wolf exceed the average dog's demand. Because of this, captive wolves typically do not cope well in urban areas. Due to their talent at observational learning, adult captive wolves can quickly work out how to escape confinement, and require constant reinforcement by caretakers or owners, which makes raising wolves difficult for people who raise their pets in an even, rather than subordinate, environment." Depending if the claim is taken to mean that it is possible, le-

1522

1523

1524

1525

1526

1527

1529

1530

1531

1532

1533

1534

1535

1537

1488

gal or practical to keep a wolf as a pet, the reader might reach different conclusions about whether the evidence supports it. The evidence provides conflicting reasons for both a supported and a refuted label. Therefore, the annotator should choose option c) ambiguous with regard to the evidence.

1436

1437

1438

1439

1440

1441

1442

1443

1444

1445

1446

1447

1448

1449

1450

1451

1452

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

1465

1466

1467

1468

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

Alternatively, if the original claim read "it is legal to keep a gray wolf as a pet", the annotator would have to choose d) the evidence does not address the question that is implied in the original claim. The original claim addresses the question of the legality of keeping a wolf as a pet, which the evidence does not cover. Differently from option c) above, this case does not provide conflicting evidence, but rather does not provide enough evidence to choose either way.

When it comes to determining the ambiguity of the original claim with regard to the evidence, feel free to rely on common knowledge to determine whether the claim and the evidence are talking about the exact same entities. For instance, while the claim above is about gray wolves, and the evidence talks more generally about wolves, the annotator may make the assumption that a gray wolf is a wolf, based on their general knowledge. Similar assumptions can be made about names, such as the name Lopez referring to Jennifer Lopez if the evidence mentions the song "Jenny from the block", or any other information that the annotator deems sufficient to determine the referrent. If the annotator does not feel confident about such coreferences, please treat such items as ambiguous with regard to the evidence.

A revised claim "keeping a gray wolf as a pet is very difficult" is supported by the evidence, as the evidence states that raising wolves requires constant reinforcement by the caretakers, which makes it difficult to keep them as pets. The claim addresses the same question as the original claim, as it still implicitly answers the question "can one keep a wolf as a pet?", just like the original claim. Therefore, the annotator should choose option a) supported by the evidence.

If, alternatively, the edited claim read "captive wolves are shy", it would not be addressing the same question as the original claim anymore, as it is not about whether one can keep a wolf as a pet, but about wolf personalities. Even though it would still be supported by the evidence and unambiguous, in this case the annotator should choose d) the revised claim does not address the same question as the original claim (or the evidence). Similarly, if the revised claim might be true but is not related to the evidence anymore, such as "wolves are predators", option d) is again the right choice.

In a different scenario, if the edited claim stated that "wolves are easy to care for as pets", this would still address the question of the difficulty of raising wolves as pets, but the annotator would have to mark it as b) refuted by the evidence. Similarly, if the revised claim states that "it is not clear from the evidence whether wolves are difficult to care for", the annotator should also choose b) refuted by the evidence, as the evidence does in fact specify the difficulty of care for wolves.

Finally, the edited claim might also not be properly disambiguated, for example, if it says that "you should not keep a wolf as a pet". This claim is even more ambiguous than the original one, as the evidence does not provide directives on whether one should attempt keeping such a pet. In this case, it should be noted as c) ambiguous with regard to the evidence.

If the original claim stated that "gray wolves eat meat", the annotator should choose a) supported by the evidence. If it was revised to the claim "gray wolves eat meat, including bones, skin and fur", this edited claim would still be addressing the question of what gray wolves eat, and it would still be a) supported by the evidence. If instead the claim was edited to read "gray wolves eat food", this would still address the question of wolf diets, but it would be c) ambiguous with regard to the evidence, as some of the elements of wolf diets may not be considered food by some readers.

D Annotation Guidelines for the Final Human Evaluation

The following are examples of claims which might be ambiguous with regard to the given evidence (original claims). Attempts have been made to disambiguate the claims by editing them to be more supported by the evidence (edited claims). Please read the original claims, the evidence, and the edited claims, and assign a score of 0 or 1, where 1 means that the revised claim is now fully supported by the evidence, whereas 0 means that it is still either ambiguous, unsubstantiated, irrelevant or refuted by the evidence.

For example, you are given the following Evidence: "[...] Captive wolves are generally shy and avoid eye contact with humans other than their owner, as well as not listening to any commands

made by any other humans. [...] Ordinary pet food 1538 is inadequate, as an adult wolf needs 12.5 kg (25 lbs) of meat daily along with bones, skin and fur to meet its nutritional requirements. [...] The exercise needs of a wolf exceed the average dog's demand. Because of this, captive wolves typically do not cope well in urban areas. Due to their talent at observational learning, adult captive wolves can quickly work out how to escape confinement, and require constant reinforcement by caretakers or owners, which makes raising wolves difficult for people who raise their pets in an even, rather than subordinate, environment."

> If the original claim is already unambiguously supported, then the only correct revision would be to keep it as is (minor changes in phrasing would be no problem):

1.

1539

1540

1541

1542

1543

1544

1545

1546

1547

1549

1550

1551

1552

1553

1555

1556

1559

1560

1561

1562

1563

1564

1565

1566

1567

1569

1570

1571

1572

1573

1574

1575

1576

1577

Original claim: "Captive wolves are shy"

Revised claim: "Captive wolves are generally shy" Score: 1

If the original claim is refuted by the evidence, then the disambiguation should simply negate it: 2.

Original claim: "Captive wolves are not shy" Revised claim: "Captive wolves are generally shy" Score: 1

On the other hand, if the original claim is ambiguous, then the revision only gets a score of 1 if it is better supported by the evidence:

3.

Original claim: "You can keep a gray wolf as a pet" Revised claim: "It may be possible to keep a gray wolf as a pet, but they are very difficult to manage" Score: 1

Anything that makes the claim unsupported by the evidence, or change the main point of the original claim, would get a score 0:

Original claim: "You can keep a gray wolf as a pet"

Revised claim: "Gray wolves avoid eye contact"

4.

5.

1578 1579

1581

1583

1584

1585

1586

1587

1588

Original claim: "You can keep a gray wolf as a pet" Revised claim: "Ordinary pet food is adequate for wolves"

Score: 0 (explicitly refuted by the evidence)

Score: 0 (irrelevant to the original claim)

In other cases the evidence might not provide enough information to disambiguate the claim, in which case that should be stated:

6.

Original claim: "It is legal to keep a gray wolf as a pet"

Revised claim: "It is not clear from the evidence whether it is legal to keep a gray wolf as a pet" Score: 1

However, if it is possible to disambiguate the claim, like in the ambiguous example 2. above, then it is not sufficient to say that it is not clear from the evidence, as there could be a better disambiguation:

7.

Original claim: "You can keep a gray wolf as a pet" Revised claim: "It is not clear from the evidence whether it is legal to keep a gray wolf as a pet" Score: 0 (the disambiguation should be provided as in example 3.)

Е **Hyperparameters**

Parameter	Values	Model
lr	$[5x10^{-4}, 5x10^{-5}, 5x10^{-6}]$	all
beam size	[1 ,5,10]	base
penalty	[1,2,3]	lp
# pseudo ref	[32,64, 128]	MBR
# hypotheses	[32,64, 128]	MBR
top p	[0.85, 0.9 ,0.95]	MBR
top k	[40, 50 ,60]	MBR
epsilon	[0.01, 0.02 ,0.03]	MBR

Table 8: Model hyperparameter values searched, for base, length penalty (lp) and Minimum Bayes Risk (MBR) models (best in bold).

F **Breakdown of Results**

Table 9 presents the model performance on the test set of DIS2DIS, separated by prediction class as well as evaluation metric. 1609

Statistical Significance Test G

Table 10 presents the results of the T-test statistical significance test, comparing each model pair.

Datasheet for Dataset н

Why was this dataset created? **H.1**

The DIS2DIS dataset was created for a novel task 1615 of disambiguation. Disambiguation is intended 1616 as an alternative method to existing explainabil-1617 ity approaches in fact-checking. The dataset was 1618 collected for training and testing models for this 1619 task. The intended use of the data is to study the 1620 phenomenon of ambiguity in the domain of factchecking. 1622

1605

1589

1590

1591

1592

1593

1594

1595

1596

1597

1598

1599

1600

1601

1602

1603

1608

1606

1610

1611

1612

Metric	Class	Base	MBR	Llama3-8B	Copy Baseline	Human
	supported	94.40	93.66	93.87	95.58	98.67
lre	refuted	93.00	93.03	93.31	93.79	98.14
ž	ambiguous	91.99	92.35	92.53	93.85	97.35
rtS	unsubstantiated	94.36	94.05	94.42	93.97	98.91
Be	all	93.07	93.05	93.30	94.10	98.05
	supported	0.90	0.90	0.85	1.00	0.90
an I	refuted	0.44	0.44	0.69	0.06	0.75
Xa Na	ambiguous	0.57	0.67	0.36	0.10	0.79
Η	unsubstantiated	0.82	0.86	0.77	0.09	0.86
	all	0.60	0.72	0.67	0.27	0.82

Table 9: The breakdown of the results by class.

Model 1	Model 2		BertScore		Hum	an Evaluatio	n
		t-statistic	p-value	< 0.05	t-statistic	p-value	< 0.05
Base	Human	-2.46	0.015	1	-17.46	5.03e - 53	1
Llama3-8B	Human	-3.52	0.000	1	-17.88	5.66e - 55	1
MBR	Human	-1.68	0.094	X	-18.52	5.60e - 58	1
Llama3-8B	Base	-1.03	0.306	X	0.71	0.466	X
MBR	Base	0.77	0.445	X	-0.04	0.965	X
MBR	Llama3-8B	1.8	0.073	X	-0.79	0.43	X
Base	Copy Baseline	-3.74	0.000	1	8.277	9.65e - 15	1
Llama3-8B	Copy Baseline	-2.92	0.004	1	6.90	4.80e - 11	1
MBR	Copy Baseline	-3.51	0.000	1	9.37	6.56e - 18	1
Human	Copy Baseline	16.82	4.47e - 50	1	11.96	4.78e - 26	1

Table 10: The results of the statistical significance test comparing different disambiguation methods.

1623 H.2 Who funded the creation of the dataset?

1624 [Anonymised]

1626

1627

1628

1629

1630

1631

1625 H.3 What preprocessing/cleaning was done?

Removal of instances was performed by a manual inspection of random samples from the dataset to ensure high quality annotations.

H.4 If it relates to people, were they told what the dataset would be used for and did they consent?

No personal data was collected. The annotators 1632 consented to the following confidentiality terms: 1633 "By participating in this research, you understand 1634 and agree that the researcher may be required to 1635 disclose your consent form, data, and other person-1636 ally identifiable information as required by law, regulation, subpoena, or court order. Otherwise, your 1638 confidentiality will be maintained in the following 1639 manner: To protect your identity, the researchers 1640 will take the following steps: (1) Each participant 1641 1642 will be assigned a number; (2) The researchers will record any data collected during the study by 1643 number, not by name; (3) Any original recordings 1644 or data files will be stored in a secured location accessed only by authorized researchers." 1646

H.5 If so, how? Were they provided with any mechanism to revoke their consent in the future or for certain uses?

1647

1648

1649

1650

1651

1652

1653

1654

1655

1656

1657

The annotators were provided the opportunity to revoke their consent at any point during the study. No option to revoke consent in the future was offered, due to the fact that the data collected was completely anonymized and it would not be possible to trace back the responses of a particular annotator.

H.6 Will the dataset be updated? How often, by whom?

The dataset may be updated in the future, to include1659other domains or languages. This would be done1660by the authors of the paper.1661