Document-level Claim Extraction and Decontextualisation for Fact-Checking

Anonymous ACL submission

Abstract

Selecting which claims to check is a timeconsuming task for human fact-checkers, especially from documents consisting of multiple sentences and containing multiple claims. 005 However, existing claim extraction approaches focus more on identifying and extracting claims from individual sentences, e.g., identifying 007 whether a sentence contains a claim or the exact 009 boundaries of the claim within a sentence. In this paper, we propose a method for *document*-011 level claim extraction for fact-checking, which aims to extract check-worthy claims from documents and decontextualise them so that they can be understood out of context. Specifically, we first recast claim extraction as extractive summarization in order to identify central sentences from documents, then rewrite them to 017 018 include necessary context from the originating document through sentence decontextualisation. Evaluation with both automatic metrics 021 and a fact-checking professional shows that our method is able to extract check-worthy claims from documents at a higher rate than previous 024 work, while also improving evidence retrieval.

1 Introduction

026

027

028

034

040

Human fact-checkers typically select a claim in the beginning of their day to work on for the rest of it. Claim extraction (CE) is an important part of their work, as the overwhelming volume of circulating claims means the choice of *what* to fact-check greatly affects the fact-checker's impact (Konstantinovskiy et al., 2021). Automated approaches to this task have been proposed to assist them in selecting check-worthy claims, *i.e.*, claims that the public has an interest in knowing the truth (Hassan et al., 2017a; Guo et al., 2022).

Existing CE methods mainly focus on detecting whether a sentence contains a claim (Reddy et al., 2021; Nakov et al., 2021b) or the boundaries of the claim within the sentence (Wührl and Klinger, 2021; Sundriyal et al., 2022; Bar-Haim et al., 2021). In real world scenarios though, claims often need to be extracted from documents consisting of multiple sentences and containing multiple claims, not all of which are relevant to the central idea of the document, and verifying all claims manually or even automatically would be inefficient. 042

043

044

045

046

047

048

050

051

053

054

060

061

062

063

064

065

066

067

068

069

070

071

072

074

075

076

077

078

079

081

Moving from sentence-level CE to documentlevel CE is challenging; we illustrate this with the example in Figure 1. Sentences in orange are claims selected by a popular sentence-level CE method (Claimbuster (Hassan et al., 2017b)) that are worth checking in principle but do not always relate to the central idea of the document, and multiple sentences with similar claims are selected, which would not all need to be fact-checked (*e.g.*, sentences 1 and 6).

Claims extracted for fact-checking are expected to be unambiguous (Lippi and Torroni, 2015; Wührl and Klinger, 2021), which means that they cannot be misinterpreted or misunderstood when they are considered outside the context of the document they were extracted from, consequently allowing them to be more easily fact-checked (Schlichtkrull et al., 2023). Figure 1 shows an example of claim decontextualisation, where the claim "Bird is scrapping thousands of e-scooters in the Middle East " requires coreference resolution to be understood out of context, e.g., "Bird" refers to "California scooter sharing start-up Bird". However, existing CE methods primarily focus on extracting sentence-level claims (i.e., extracting sentences that contain a claim) from the original document (Reddy et al., 2021) and ignore their decontextualisation, resulting in claims that are not unambiguously understood and verified.

To address these issues, we propose a novel method for *document-level* claim extraction and decontexualisation for fact-checking, aiming to extract salient check-worthy claims from documents that can be understood outside the context of the document. Specifically, assuming that salient

Document (CNBC News)

Document-level Claim Extraction			
Sentence-level: 8, 6, 1	Document-level: 4, 5, 7		
Gold Claim (Fact-checking Organization Misbar)			

Gold Claim (ract-checking Organization, wisbar).

Bird e-scooters are shutting down service in the Middle East, and scrapping as many as 10,000 scooters.

Claim extracted by decontextualising the 4th sentence:

<u>California scooter sharing start-up Bird</u> is scrapping thousands of e-scooters in the Middle East and shutting down its operations in the majority of the region as a result of the coronavirus pandemic, according to five people familiar with the matter.

Figure 1: An example of document-level claim extraction. Document¹ is a piece of news from CNBC. Gold Claim² is annotated by the fact-checking organization, Misbar. Sentences in orange denote check-worthy claims extracted by sentence-level CE (Claimbuster). Sentences in blue denote salient claims extracted by our document-level CE. The claim in green is a decontextualised claim derived from the 4th sentence obtained by our document-level CE.

claims are derived from central sentences, *i*) we recast the document-level CE task into the extractive summarization task to extract central sentences and reduce redundancy; *ii*) we decontextualise central sentences to be understandable out of context by enriching them with the necessary context; *iii*) we introduce a QA-based framework to obtain the necessary context by converting ambiguous information units in the sentence into declarative sentences.

To evaluate our method we derive a dataset containing decontextualised claims from AVeriTeC (Schlichtkrull et al., 2023), a recently proposed dataset containing decontextualised claims from real-world fact-checking articles³. Our method achieves a Precision@1 score of 47.8 on identifying central sentences, with 10% improvement over Claimbuster. This was further assessed manually by a fact-checking professional, as the sentences returned by our method were more often central to the document, and more often check-worthy than those extracted by Claimbuster. Additionally, our method achieved a chrf score of 26.4 on gold claims decontextualised by the fact-checkers, outperforming all baselines. When evaluated for evidence retrieval potential, the decontextualised claims obtained by enriching original sentences with the necessary context, are better than the original claim sentences, with an average 1.08 improvement in precision.

096

100

101

102

103

105

106

107

108

109

110

bird-circ-scooters-middle-east.html

²https://misbar.com/en/factcheck/2020/06/18/ are-bird-e-scooters-leaving-the-middle-east

³Our code and data will be released on Github.

2 Related Work

Claim Extraction Claim extraction is typically framed either as a classification task or claim boundary identification task. The former framing focuses on detecting whether a given sentence contains a check-worthy claim. Claimbuster (Hassan et al., 2017b), the most popular method in this paradigm, computes the score of how important a sentence is to be fact-checked. More similar to our work are studies that formulate the checkworthy claim detection task as the sentence ranking task. For example, Zhou et al. (2021) present a sentence-level classifier by combining a fine-tuned hate-speech model with one dropout layer and one classification layer to rank sentences. However, these methods were not able to handle the challenges of document-level claim extraction, e.g., reduce redundant claim sentences. The framing of claim extraction as boundary identification focuses on detecting the exact claim boundary within the sentence. Nakov et al. (2021a) propose a BERTbased model to refine the claim detection service (Levy et al., 2014) for identifying the boundaries of the claim within the sentence. Sundrival et al. (2022) tackle claim span identification as a token classification task for identifying argument units of claims in the given text. Unlike the above methods, where the claims are extracted from given sentences, our work aims to extract salient checkworthy claims from documents, thus addressing the limitations of sentence-level methods in extracting salient claim sentences and avoiding redundancy.

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

138

139

140

141

¹https://www.cnbc.com/2020/06/03/



Figure 2: An overview of our document-level claim extraction framework. Given an input document, we first use extractive summarization to rank all sentences and select summary sentences as central sentences. Then, we describe a QA-based framework to generate a specific high-quality context for important information units in the sentence. Next, we use a seq2seq generation model to decontextualise sentences by enriching them with their corresponding context. Finally, a claim check-worthiness classifier is used to select salient check-worthy claim sentences based on the score that reflects the degree to which sentences belong to the check-worthy claim.

Decontextualisation Choi et al. (2021) propose 143 two different methods for decontextualisation, 144 145 based on either a coreference resolution model or a seq2seq generation model. Both methods use the 146 sentences in the paragraph containing the target 147 sentence as context to rewrite the target sentence. 148 Newman et al. (2023) utilize an LLM to generate 149 QA pairs for each sentence by designing specific 150 prompts, and then use an LMM with QA pairs to 151 rewrite each sentence. Sundriyal et al. (2023) propose to combine chain-of-thought and in-context 153 learning for claim normalization. Unlike the above methods, we generate declarative sentences for po-155 tentially ambiguous information units in the tar-156 get sentence based on the whole document, and 157 combine them into context to rewrite the target sentence, thereby improving the sensitivity to am-159 biguity when decontextualising.

3 Method

161

162As illustrated in Figure 2, our proposed document-163level claim extraction framework which consists164of four components: i) Sentence extraction ($\S3.1$);165extracts the sentences related to the central idea the166document as candidate claim sentences; ii) Context

generation (§3.2), extracts context from the document for each candidate sentence; *iii*) Sentence decontextualisation (§3.3), rewrites each sentence with its corresponding context to be understandable out of context; *iv*) Check-worthiness estimation (§3.4), selects the final check-worthy claims from candidate decontextualised sentences. 167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

184

186

187

188

190

3.1 Sentence Extraction

The claims selected by human fact-checkers are typically related to the central idea of the document considered. Thus we propose to model sentence extraction as extractive summarization. For this purpose, we concatenate all the sentences in the document into an input sequence, which is then fed to BertSum (Liu and Lapata, 2019), a commonly-used method for document-level extractive summarization, to extract central sentences. Specifically, given a document consisting of n sentences $\mathcal{D} = \{s_1, s_2, ..., s_n\}$, we first formulate the input sequence C as "[CLS] s_1 [SEP] [CLS] s_2 $[SEP] \dots [CLS] s_n [SEP]$ ", where [CLS] and [SEP]denote the start and end token for each sentence, respectively, and then feed them into a pre-trained encoder BERT to obtain the sentence representation s.

236

237

238

250 251

249

252 253

254

255

257

258

259

260

261

262

263

264

265

267

269

270

271

272

273

274

275

276

277

191 192

193

210

213

Finally, a linear layer on sentence representations is used to score sentences.

$$\mathbf{s}_{i} = \text{BERT}(s_{i})$$

$$score_{i} = \sigma(W\mathbf{s}_{i} + b_{0})$$
(1)

where σ is a sigmoid function, s_i denotes the rep-194 resentation of the *i*-th [CLS] token, *i.e.*, the rep-195 resentation of the *i*-th sentence in the document, 196 and $score_i$ denotes the score of the *i*-th sentence. 197 All sentences are constructed into an ordered set S198 according to their scores. Since all sentences are 199 ranked by sentence-level scoring, some top-scoring sentences may have the same or similar meaning. 201 To avoid redundancy, we add a textual entailment model DocNLI (Yin et al., 2021) on the top of BertSum output to remove redundant sentences by calculating the entailment scores between sen-205 tences. Following previous work (Liu and Lapata, 206 2019), we only select the top-k summary sentences with the highest scores in Equation 1 as candidate central sentences.

$$S' = \text{DocNLI}(S) \tag{2}$$

where $S' = \{s'_1, s'_2, ..., s'_i, ..., s'_k\}$ is a set of central 211 sentences that do not contain the same meaning. 212

3.2 **Context Generation**

After sentence extraction, the next step is to clarify the (possibly) ambiguous sentences in S' by rewrit-215 ing them with their necessary context. Unlike Choi 216 217 et al. (2021) where the context consists of a sequence of sentences in the paragraph containing the 218 ambiguous sentence, we need to consider the whole 219 document, *i.e.*, sentences from *different* paragraphs. To improve the sensitivity to ambiguity when decontextualising, we propose a QA-based context generation framework to produce a specific highquality context for each ambiguous sentence, which contains three components: *i*) Question Generation: extracts potentially ambiguous information units from the sentence and generates questions with 227 them as answers; ii) Question Answering: finds 228 more information about ambiguous information units by answering generated questions with the whole document; *iii*) QA-to-Context Generation: converts question-answer pairs into declarative sentences and combines them into context specific to the sentence. In the following subsections we describe each component in detail.

Question Generation. To identify ambiguous information units in candidate central sentences, we first use Spacy⁴ to extract named entities, pronouns, nouns, noun phrases, verbs and verb phrases in the sentence i as the potentially ambiguous information units $U_i = \{u_i^1, u_i^2, ..., u_i^j, ..., u_i^m\}, i \in [1, 2, ..., k],$ where u_i^j denotes the *j*-th information unit of the *i*-th candidate sentence s'_i .

Once the set of information units for a sentence U_i is identified, we then generate a question for each of them. Specifically, we concatenate u_i^j and s'_i in which u^j_i is located as the input sequence and feed it into a question generator QG (Murakhovs'ka et al., 2022) to produce the question q_i^j with u_i^j as the answer.

$$Q_i = \{q_i^j\}_{j=1}^m = \{\text{QG}(s'_i, u_i^j)\}_{j=1}^m$$
(3)

where Q_i denotes the set of questions corresponding to U_i in the sentence s'_i .

Question Answering. After question generation, our next step is to clarify ambiguous information units by answering corresponding questions with the whole document \mathcal{D} . Specifically, we first use BM25 (Robertson et al., 2009) to retrieve relevant evidence E related to the question q_i^j from \mathcal{D} , and then answer q_i^j with the retrieved evidence E using an existing QA model (Khashabi et al., 2022):

$$E = BM25(\mathcal{D}, q_i^j)$$

$$a_i^j = QA(E, q_i^j)$$
(4)

where a_i^j denotes a more complete information unit corresponding to u_j^i , e.g., a complete coreference. We denote all question-answer pairs of the *i*-th sentence as $P_i = \{(q_i^1, a_i^1), (q_i^2, a_i^2), \dots, (q_i^m, a_i^m)\}.$

QA-to-Context Generation. After question answering, we utilize a seq2seq generation model to convert QA pairs P_i into the corresponding highquality context C'_i . Specifically, we first concatenate the question q_i^j and the answer a_i^j as the input sequence, and then output a sentence using the BART model (Lewis et al., 2019) finetuned on the QA2D (Demszky et al., 2018) dataset.

$$\tilde{s}_i^j = \text{BART}(q_i^j, a_i^j) \tag{5}$$

where \tilde{s}_i^j is a declarative sentence corresponding to the information unit u_i^j . Finally, all generated

⁴https://spacy.io

Data	#sample	med.sent	avg.sent	len.claim	len.document
Train	830	9	1.09	17	274
Dev	149	6	1.01	16	120
Test	252	5	1.05	16	63
All	1231	7	1.07	17	17

Table 1: Descriptive statistics for AVeriTeC-DCE. #sample refers to the number of samples available for claim extraction, *i.e.*, the total number of accessible *source url*, med.sent refers to the median number of sentences in documents. *avg.sent* refers to the average number of sentences in claims, len.claim refers to the median length of claims in words, len.document refers to the median length of documents in words.

sentences are combined into high-quality context $C'_i = {\tilde{s}_i^1, \tilde{s}_i^2, ..., \tilde{s}_i^m}$ corresponding to the information units U_i in sentence s'_i , which is then used in the next decontextualisation step to enrich the ambiguous sentences.

3.3 Sentence Decontextualisation

279

281

284

287

291

294

295

296

297

302

Sentence decontextualisation aims to rewrite sentences to be understandable out of context, while retaining their original meaning. To do this, we use a seq2seq generation model T5 (Raffel et al., 2020) to enrich the target sentence with its corresponding high-quality context generated in the previous step. Specifically, we first formulate the input sequence as "[CLS] \tilde{s}_i^1 [SEP] \tilde{s}_i^2 [SEP] \tilde{s}_i^m [SEP] s_i' ", where s_i' denotes the potential ambiguous sentence and [SEP] is a separator token between declarative sentences. We then feed the input sequence to a decontextualisation model D (Choi et al., 2021) to rewrite the sentence. Similarly, we set the output sequence to be [CAT] [SEP] y.

$$y_{i} = \begin{cases} D(s'_{i}, C'_{i}), if CAT = feasible\\ s'_{i} , if CAT = infeasible\\ s'_{i} , if CAT = unnecessary \end{cases}$$
(6)

where CAT = feasible or infeasible denotes that s'_i can or cannot be decontextualised, CAT = unnecessary denotes that s'_i can be understood without being rewritten, y_i denotes the *i*-th decontextualised sentence.

3.4 Check-Worthiness Estimation

305Unlike existing CE methods that determine whether306a sentence is worth checking without considering307the context, we estimate the check-worthiness of a308sentence after decontextualisation because some309sentences may be transformed from not check-310worthy into check-worthy ones after decontextualisation.311alisation.

trained on the ClaimBuster dataset (Arslan et al., 2020) to classify sentences into three categories: Check-worthy Factual Sentence (CFS), Unimportant Factual Sentence (UFS) and Non-Factual Sentence (NFS). 312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

$$score(y_i) = \text{DeBERTa}(class = \text{CFS} \mid y_i)$$

$$\text{claim} = \operatorname{argmax}\{score(y_i)\}_{i=1}^k$$
(7)

where $score(y_i)$ reflects the degree to which the decontextualised sentence y_i belongs to CFS, and claim denotes that the final salient check-worthy claim that can be understood out of context.

4 Dataset

We convert AVeriTeC, a recently proposed dataset for real-world claim extraction and verification (Schlichtkrull et al., 2023), into a dataset AVeriTeC-DCE for the document-level CE task. AVeriTeC is collected from 50 different fact-checking organizations, which contains 4568 real-world claims. Each claim is associated with attributes such as its type, source and date. In this work, we mainly focus on the following attributes relevant to claim extraction: i) claim, the claim as extracted by the fact-checkers and decontextualised by annotators, and *ii*) source *url*: the URL linking to the original web article of the *claim*. The task of this work is to extract the claims as extracted from the source url. We also consider whether claims need to be decontextualised when extracting them from documents, as they will directly affect the subsequent evidence retrieval and claim verification.

To extract claim-document pairs from AVeriTeC that can be used for document-level CE, we performed the following steps: 1) Since we focus on the extraction of textual claims, we do not include *source urls* containing images, video or audio; 2) To extract the sentences containing the claims from *source urls*, we build a web scraper to extract text 348data in the source url as the document. We found349that the attribute source url is not always avail-350able in samples, thus we only select those samples351where the web scraper can return text data from352source urls. We obtain a dataset AVeriTeC-DCE,353containing 1231 available samples, for document-354level CE. We will make this dataset available on-355line. Statistics for AVeriTeC-DCE are described in356Table 1.

5 Experiments

357

364

367

371

374

375

379

380

Our approach consists of two components, *i.e.*, sentence extraction and decontextualisation. We conduct separate experiments to evaluate them in detail, as well as an overall evaluation for document-level claim extraction.

5.1 Sentence Extraction

We compare our sentence extraction method, the combination of BertSum and DocNLI stated in Section 3.1, against other baselines through automatic evaluation and human evaluation.

Baselines 1) Lead sentence: the lead (first) sentence of most documents is considered to be the most salient, especially in news articles (Narayan et al., 2018); 2) Claimbuster (Hassan et al., 2017b): we use this well-established method to compute the check-worthiness score of each sentence and we rank sentences based on their scores; 3) LSA (Gong and Liu, 2001): a common method of identifying central sentences of the document using the latent semantic analysis technique; 4) TextRank (Mihalcea and Tarau, 2004): a graph-based ranking method for identifying important sentences in the document; 5) BertSum (Liu and Lapata, 2019): a BERT-based document-level extractive summarization method for ranking sentences.

383 Automatic Evaluation Since the central sentences of documents are not given in AVeriTeC, we cannot evaluate the extracted central sentences 385 by exact matching. Thus, we instead rely on the sentence that has the highest chrf (Popović, 2015) 387 score with the *claim*, as the *claim* is the central claim annotated by human fact-checkers. We use Precision@k as the evaluation metric, which denotes the probability that the first k sentences in the extracted sentences contain the central sentence. Table 2 shows the results of different sentence extraction methods. We can see that our method outperforms all baselines in identifying the central

Method	P@1	P@3	P@5	P@10
Claimbuster	37.8	59.1	65.7	71.4
Lead Sentence	42.3	-	-	-
LSA	38.4	55.3	62.1	70.4
TextRank	42.7	60.6	65.1	71.2
BertSum	43.4	61.6	67.5	72.3
Ours	47.8	63.1	68.6	73.8

Table 2: Results with different sentence extraction methods. P@k denotes the probability that the first k sentences in the ranked sentences contain the central sentence.

	Claimbuster	Ours
IsCheckWorthy	0.36	0.44
IsCentralClaim	0.24	0.68

Table 3: Human Evaluation of sentence extraction ontwo different dimensions.

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

sentence, achieving a P@1/P@3/P@5/P@10 score of 47.8/63.1/68.6/73.8, which indicates that the combination of the extractive summarization (Bert-Sum) and entailment model (DocNLI) can better capture the central sentences and reduce redundant sentences. We found that the common extractive summarization methods (e.g., Lead Sentence, TextRank and BertSum) are better than Claimbuster on P@1, confirming what we had stated in the introduction, that sentence-level CE methods have limitations when they are applied at the documentlevel CE. Moreover, we observe that the lead sentence achieves a P@1 score of 42.3, indicating that there is a correlation between the sentences selected for fact-checking and the lead sentence that often served as the summary. We list the source URLs of the samples for claim extraction in Appendix A1.

Human Evaluation To further compare sentences extracted by our method and Claimbuster, we asked a fact-checking professional to evaluate the quality of extracted sentences on the following two dimensions: 1) **IsCheckWorthy**: is the sentence worth checking? 2) **IsCentralClaim**: is the sentence related to the central idea of the article? We randomly select 50 samples, each containing at least 5 sentences. For simplicity, we only select the top-1 sentence returned by each method for comparison. As shown in Table 3, we observed that 68% of the central sentences extracted by our method are related to the central idea of the doc-

Coreference Resolution
Sentence: <u>He</u> has pubicly stated that he sympathized with their cause and even hinted that he would provide them with American resources should they be in need during his 2008 State of the Union address. Decontextualized sentence: <u>President Obama</u> has pubicly stated that he sympathized with their cause and even hinted that he would provide them with American resources should they be in need during his 2008 State of the Union address.
Global Scoping
Sentence: <u>During the attack</u> , Capitol Police made the request again. Decontextualized sentence: During the attack on Washington D.C., Capitol Police made the request again.
Bridge Anaphora
Sentence: <u>The government</u> does not have proper storage facilities for stocking such a large amount of excess grain. Decontextualized sentence: The government of India does not have proper storage facilities for stocking such a large amount of excess grain.

Figure 3: Case studies of sentence decontextualisation solving linguistic problems, such as coreference resolution, global scoping and bridge anaphora.

426 ument compared to 24% of Claimbuster, which further supports our conclusion obtained by auto-427 matic evaluation, *i.e.*, the sentences extracted by 428 our method were more often central to the doc-429 ument, and more often check-worthy than those 430 that extracted by Claimbuster. This indicates that 431 432 when identifying salient check-worthy claims from documents, it is not enough to consider whether a 433 sentence is worth checking at the sentence level, 434 but also whether the sentence is related to the cen-435 tral idea of the document. Thus, we believe that the 436 claims related to the central idea of the document 437 are the ones that the public is more interested in 438 knowing the truth. 439

5.2 Decontextualisation

440

441 442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

To evaluate the effectiveness of decontextualisation on evidence retrieval, for a fair comparison, we select the sentence that has the highest chrf score with the *claim* as the best sentence as we considered in Section 5.1, and conduct a comparison between it and its corresponding decontextualised sentence. The evidence set used for evaluation is retrieved from the Internet using the Google Search API given a claim, each containing gold evidence and additional distractors (Schlichtkrull et al., 2023). We use Precision@k as the evaluation metric.

Baselines 1) Coreference model: decontextualisation by replacing unresolved coreferences in the target sentence, *e.g.*, (Joshi et al., 2020); 2) Seq2seq model: decontextualisation by rewriting the target sentence with necessary context (Choi et al., 2021).

Retrieval-based Evaluation Following previous work (Choi et al., 2021), we compare our QA-based decontextualisation method against other baselines through retrieval-based evaluation. We use BM25 as the retriever to find evidence with different sentences as the query. Table 4 shows the results for evidence retrieval with different decontextualised sentences on the dev set of AVeriTeC-DCE. We found that our method outperforms all baselines, and improves the P@3/P@5/P@10 score over the original sentence by 1.42/0.82/0.99, achieving an average 1.08 improvement in precision. After further analysis, we found that only 21/149 sentences are decontextualised by our method, and 17/21 of these sentences obtain better evidence retrieval, with an average improvement of 1.21 in precision over the original sentences, proving that decontextualisation enables evidence retrieval more effectively.

Method	P@3	P@5	P@10
Sentence	35.45	44.72	61.31
Coreferece	36.02	44.98	61.79
Seq2seq (Context)	36.17	45.04	61.99
Seq2seq (Context*)	36.87	45.54	62.30

Table 4: Results for evidence retrieval with different decontextualised sentences. Context consists of a sequence of sentences in the paragraph containing the target sentence. Context* consists of declarative sentences generated by our context generation module.

Case Study Figure 3 illustrates three case studies of sentence decontextualisation. The first case is an example that requires coreference resolution. To make the sentence understandable out of context, these words (*e.g.*, "*He*", "*their*") need to be rewritten with the context. After decontextualisation, we can see that "*He*" is rewritten to "*President Obama*", which helps us understand the sentence better without context. As for "*their*", we cannot decontextualise it because there is no information about this word in the document. This supports

475

476

477

478

479

480

481

482

483

484

485

462

463

464

465

466

467

468

469

470

471

472

the claim that providing a high-quality context is 486 necessary for better decontextualisation. The sec-487 ond case is an example that requires global scoping, 488 which requires adding a phrase (e.g., prepositional 489 phrase) to the entire sentence to make it better un-490 derstood. In this case, we add "on Washington D.C." 491 as a modifier to "During the attack" to help us un-492 derstand where the attack took place. The third 493 case is an example that requires a bridge anaphora, 494 where the phrase noun "The government" becomes 495 clear by adding a modifier "India". In summary, 496 decontextualising the claim is helpful for humans 497 to better understand the claim without context. 498

5.3 Document-level Claim Extraction

499

500

504

505

506

507

509

510

511

512

513

514

515

In this section, we put two components together to conduct an overall evaluation for document-level CE, e.g., extracting the claim in green from the document in Figure 1. We first select top-3 sentences returned by different sentence extraction methods as candidate central sentences, and then feed them into our decontextualisation model to obtain decontextualised claim sentences, finally use a claim check-worthiness classifier to select the final claim. We evaluate performance by calculating the similarity between our final decontextualised claim sentence and the *claim* decontextualised by factcheckers. We use the chrf as the evaluation metric. The chrf computes the similarity between texts using the character n-gram F-score. Other metrics are reported in Appendix A2.

Statistics for decontextualisation are described 516 in Table 5, we observe that 122 out of 1231 (10%)517 sentences can be decontextualised (feasible); 122 out of 1231 (10%) sentences cannot be decontex-519 tualised (infeasible); and 987 out of 1231 (80%) 520 sentences can be understood without being decon-521 textualised (unnecessary), including the lead sentence. Since the lead sentence of the document is 523 often considered to be the most salient, we do not decontextualise the lead sentence. As shown in Ta-525 ble 6, our method achieves a chrf score of 26.4 on 526 gold claims decontextualised by the fact-checkers, 527 outperforming all baselines. We observe that the performance of claim extraction and sentence ex-529 traction is positively correlated, *i.e.*, the closer the extracted sentence is to the central sentence, the 531 more similar the extracted claim is to the *claim*, 532 which supports our assumption that salient claims 533 are derived from central sentences. For this reason, the performance of our document-level CE is

Data	#claim	#fea.	#infea.	#unnec.
All	1231	122	122	987

Table 5: Statistic of decontextualisation. #fea./infea. denotes the number of sentences that can/cannot be decontextualised, #unnec. denotes the number of sentences that can be understood without context.

Method	chrf			
	Sentence*	Dec. Sentence*		
Claimbuster	24.3	24.5		
Lead Sentence	23.8	-		
LSA	24.1	24.3		
TextRank	24.5	25.4		
BertSum	25.6	25.9		
Ours	25.9	26.4		

Table 6: Results of Document-level CE. Sentence^{*} denotes the best sentence returned by different sentence extraction methods. Dec. Sentence^{*} denotes the decontextualised Sentence^{*}.

limited by the performance of sentence extraction, *i.e.*, if our sentence extraction method cannot find the gold central sentence, decontextualisation may not improve the performance of CE, and may even lead to a decrease in the performance of CE due to noise caused by decontextualisation.

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

6 Conclusions and Future work

This paper presented a document-level claim extraction framework for fact-checking, aiming to extract salient check-worthy claims from documents that can be understood out of context. To extract salient claims from documents, we recast the claim extraction task as the extractive summarization task to select candidate claim sentences. To make sentences understandable out of context, we introduce a QA-based decontextualisation model to enrich them with the necessary context. The experimental results show the superiority of our method over previous methods, including document-level claim extraction and evidence retrieval, as indicated by human evaluation and automatic evaluation. In future work, we plan to extend our document-level claim extraction method to extract salient checkworthy claims from multimodal web articles.

610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

608

609

Limitations

560

While our method has demonstrated superiority in extracting salient check-worthy claims and im-562 proving evidence retrieval, we recognize that our method is not able to decontextualise all ambiguous sentences, particularly those that lack the nec-565 566 essary context in the source url. Also, human fact-checkers have different missions, thus checkworthiness claims to one fact-checking organiza-568 tion may not be check-worthiness to another organization (*i.e.*, some organizations check parody 570 claims or claims from satire websites, while oth-571 ers do not). Furthermore, since the documents we 572 use are extracted from the *source url*, a powerful web scraper is required when pulling documents 574 from source urls. Moreover, our method assumes salient claims are derived from central sentences. Although this assumption is true in most cases, it 577 may be inconsistent with central claims collected by human fact-checkers. Besides, we use the chrf 579 metric to calculate the similarity between claims extracted from the source url and the gold claim, while gold claims are decontextualised by the fact-582 checkers with fact-checking articles and may con-583 tain information that is not in the original article, 584 thus the metrics used to evaluate document-level 585 claims are worth further exploring.

Ethics Statement

588

589

590

592

593

594

595

597

598

604

606

We rely on fact-checks from real-world factcheckers to develop and evaluate our models. Nevertheless, as any dataset, it is possible that it contains biases which influenced the development of our approach. Given the societal importance of factchecking, we advise that any automated system is employed with human oversight to ensure that the fact-checkers fact-check appropriate claims.

References

- Fatma Arslan, Naeemul Hassan, Chengkai Li, and Mark Tremayne. 2020. A benchmark dataset of checkworthy factual claims. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 14, pages 821–829.
- Roy Bar-Haim, Yoav Kantor, Elad Venezian, Yoav Katz, and Noam Slonim. 2021. Project debater apis: decomposing the ai grand challenge. *arXiv preprint arXiv:2110.01029*.
- Eunsol Choi, Jennimaria Palomaki, Matthew Lamm, Tom Kwiatkowski, Dipanjan Das, and Michael

Collins. 2021. Decontextualization: Making sentences stand-alone. *Transactions of the Association for Computational Linguistics*, 9:447–461.

- Dorottya Demszky, Kelvin Guu, and Percy Liang. 2018. Transforming question answering datasets into natural language inference datasets. *arXiv preprint arXiv:1809.02922*.
- Yihong Gong and Xin Liu. 2001. Generic text summarization using relevance measure and latent semantic analysis. In *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 19–25.
- Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking. *Transactions of the Association for Computational Linguistics*, 10:178–206.
- Naeemul Hassan, Fatma Arslan, Chengkai Li, and Mark Tremayne. 2017a. Toward automated fact-checking: Detecting check-worthy factual claims by claimbuster. In *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1803–1812.
- Naeemul Hassan, Gensheng Zhang, Fatma Arslan, Josue Caraballo, Damian Jimenez, Siddhant Gawsane, Shohedul Hasan, Minumol Joseph, Aaditya Kulkarni, Anil Kumar Nayak, et al. 2017b. Claimbuster: The first-ever end-to-end fact-checking system. *Proceedings of the VLDB Endowment*, 10(12):1945–1948.
- Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S Weld, Luke Zettlemoyer, and Omer Levy. 2020. Spanbert: Improving pre-training by representing and predicting spans. *Transactions of the association for computational linguistics*, 8:64–77.
- Daniel Khashabi, Yeganeh Kordi, and Hannaneh Hajishirzi. 2022. Unifiedqa-v2: Stronger generalization via broader cross-format training. *arXiv preprint arXiv:2202.12359*.
- Lev Konstantinovskiy, Oliver Price, Mevan Babakar, and Arkaitz Zubiaga. 2021. Toward automated factchecking: Developing an annotation schema and benchmark for consistent automated claim detection. *Digital threats: research and practice*, 2(2):1–16.
- Ran Levy, Yonatan Bilu, Daniel Hershcovich, Ehud Aharoni, and Noam Slonim. 2014. Context dependent claim detection. In *Proceedings of COLING* 2014, the 25th International Conference on Computational Linguistics: Technical Papers, pages 1489– 1500.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.

751

752

753

754

755

756

757

758

718

- Marco Lippi and Paolo Torroni. 2015. Contextindependent claim detection for argument mining. In *IJCAI*.
 - Yang Liu and Mirella Lapata. 2019. Text summarization with pretrained encoders. *arXiv preprint arXiv:1908.08345*.
 - Rada Mihalcea and Paul Tarau. 2004. Textrank: Bringing order into text. In *Proceedings of the 2004 conference on empirical methods in natural language processing*, pages 404–411.
- Lidiya Murakhovs'ka, Chien-Sheng Wu, Philippe Laban, Tong Niu, Wenhao Liu, and Caiming Xiong.
 2022. MixQG: Neural question generation with mixed answer types. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1486–1497, Seattle, United States. Association for Computational Linguistics.

673

674

675

676

677

678

679

691

692

697

700

701

707

710 711

712

714

716

- Preslav Nakov, David Corney, Maram Hasanain, Firoj Alam, Tamer Elsayed, Alberto Barrón-Cedeño, Paolo Papotti, Shaden Shaar, and Giovanni Da San Martino. 2021a. Automated fact-checking for assisting human fact-checkers. *arXiv preprint arXiv:2103.07769*.
- Preslav Nakov, Giovanni Da San Martino, Tamer Elsayed, Alberto Barrón-Cedeno, Rubén Míguez, Shaden Shaar, Firoj Alam, Fatima Haouari, Maram Hasanain, Nikolay Babulkov, et al. 2021b. The clef-2021 checkthat! lab on detecting check-worthy claims, previously fact-checked claims, and fake news. In *ECIR 2021, Virtual Event, March 28–April 1, 2021, Proceedings, Part II 43*, pages 639–649. Springer.
- Shashi Narayan, Shay B Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. *arXiv preprint arXiv:1808.08745*.
- Benjamin Newman, Luca Soldaini, Raymond Fok, Arman Cohan, and Kyle Lo. 2023. A controllable qa-based framework for decontextualization. *arXiv preprint arXiv:2305.14772*.
- Maja Popović. 2015. chrf: character n-gram f-score for automatic mt evaluation. In *Proceedings of the tenth workshop on statistical machine translation*, pages 392–395.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551.
- Revanth Gangi Reddy, Sai Chinthakindi, Zhenhailong Wang, Yi R Fung, Kathryn S Conger, Ahmed S Elsayed, Martha Palmer, and Heng Ji. 2021. Newsclaims: A new benchmark for claim detection from news with background knowledge. *arXiv preprint arXiv:2112.08544*.

- Stephen Robertson, Hugo Zaragoza, et al. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends*® *in Information Retrieval*, 3(4):333–389.
- Michael Schlichtkrull, Zhijiang Guo, and Andreas Vlachos. 2023. Averitec: A dataset for real-world claim verification with evidence from the web. *arXiv preprint arXiv:2305.13117*.
- Megha Sundriyal, Tanmoy Chakraborty, and Preslav Nakov. 2023. From chaos to clarity: Claim normalization to empower fact-checking. *arXiv preprint arXiv:2310.14338*.
- Megha Sundriyal, Atharva Kulkarni, Vaibhav Pulastya, Md. Shad Akhtar, and Tanmoy Chakraborty. 2022. Empowering the fact-checkers! automatic identification of claim spans on Twitter. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7701–7715, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Amelie Wührl and Roman Klinger. 2021. Claim detection in biomedical twitter posts. *arXiv preprint arXiv:2104.11639*.
- 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.
- Wenpeng Yin, Dragomir Radev, and Caiming Xiong. 2021. Docnli: A large-scale dataset for document-level natural language inference. *arXiv preprint arXiv:2106.09449*.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.
- Xinrui Zhou, Bohuai Wu, and Pascale Fung. 2021. Fight for 4230 at checkthat! 2021: Domain-specific preprocessing and pretrained model for ranking claims by check-worthiness. In *CLEF (Working Notes)*, pages 681–692.

A1 Statistic of Source URLs

URL	#sample
twitter.com	241
facebook.com	235
perma.cc	63
channelstv.com	37
aljazeera.com	34
president.go.ke	27
gov.za	25
instagram.com	18
c-span.org	16
factba.se	13
axios.com	12
youtu.be	12
rumble.com	12
abcnews.go.com	11
rev.com	11
cnn.com	10
news24.com	10
punchng.com	10
washingtonpost.com	10
cbsnews.com	8
foxnews.com	8
misbar.com	7
thegatewaypundit.com	7
politifact.com	7
nypost.com	7
nbcnews.com	7
telegraph.co.uk	7
wisn.com	6
tatersgonnatate.com	6
bustatroll.org	6
dailymail.co.uk	5
whitehouse.gov	5

Table A1: Statistic of the source URLs of the samples for document-level CE. We only list URLs with a total number number greater than 5.

A2 Evaluation Metrics

We use the following metrics to assess the similarity between the claim decontextualised by our method and the *claim* decontextualised by fact-checkers. **SARI** (Xu et al., 2016) is developed to compare the claim with the reference claim by measuring the goodness of words that are added, deleted and kept. **BERTScore** (Zhang et al., 2019) is utilized to compute the semantic overlap between the claim and the reference claim by sentence representation. Since most *claims* in AVerTeC are decontextualised by fact-checkers with fact-checking articles, they may contain some information that is not in the *source url*, making it challenging for SARI and BERTScore to be used as evaluation metrics in this task. Thus, we use the chrf as our main evaluation metric for claim extraction.

Method		Sentence*			Dec. Sentence*		
	SARI	BERTScore	chrf	SARI	BERTScore	chrf	
Claimbuster	6.23	82.7	24.3	6.24	82.8	24.5	
Lead Sentence	6.41	83.4	23.8	-	-	-	
LSA	5.56	83.2	24.1	5.57	83.2	24.3	
TextRank	6.60	83.1	24.5	6.61	83.1	25.4	
BertSum	6.54	83.6	25.6	6.55	83.6	25.9	
DCE(Ours)	6.56	83.7	25.9	6.70	83.8	26.4	

Table A2: Results of Document-level CE on three different metrics.