# TRACSUM: A New Benchmark for Aspect-Based Summarization with Sentence-Level Traceability in Medical Domain

Anonymous ACL submission

#### Abstract

While document summarization with LLMs has enhanced access to textual information, concerns about the factual accuracy of these summaries persist (e.g., hallucination), especially in the medical domain. Identifying citations from which summaries are derived enables users to assess their accuracy, thereby alleviating this concern. In this paper, we introduce TRACSUM, a novel benchmark for traceable, aspect-based summarization, in which generated summaries are paired with sentence-level citations, enabling users to trace back to the original context. First, we annotate 500 medical abstracts<sup>1</sup> for seven key medical aspects, yielding 3.5K summary-citations pairs. We then propose a fine-grained evaluation framework for this new task, designed to assess the completeness and consistency of generated content using four metrics. Finally, we introduce a summarization pipeline, TRACK-THEN-SUM, which serves as a baseline method for comparison. In experiments, we evaluate both this baseline and a set of LLMs on TRACSUM, and conduct a human evaluation to assess the evaluation results. The findings demonstrate that TRACSUM can serve as an effective benchmark for traceable, aspect-based summarization tasks. We also observe that explicitly performing sentence-level tracking prior to summarization enhances generation accuracy, while incorporating the full context further improves summary completeness. The visualized dataset is anonymously available at https://www.tracsum.info.

# 1 Introduction

011

022

041

New findings observed in clinical trials are published in journal articles, which describe their design and outcomes (Hariton and Locascio, 2018), serving as a crucial foundation for evidence-based medicine (EBM) (Sackett, 1997; Joseph et al., 2024). Ideally, medical professionals would stay



Figure 1: Schematic diagram of the TRACSUM task, where aspect-based summaries are enriched with sentence-level citations linking back to their corresponding source sentences in the medical article.

current on all medical evidence from these articles to support their decision-making, but this is impractical due to the volume and growth of the evidence base (Marshall et al., 2021; Frihat and Fuhr, 2024).

044

045

047

049

051

052

060

061

063

Document summarization condenses the input document into a concise and coherent text that retains salient information (Narayan et al., 2018; Zheng et al., 2020; Wang et al., 2022; Zhang et al., 2023b). Recent advancements in document summarization methods have shown promising results in generating overall summaries (Rush et al., 2015; Cheng and Lapata, 2016; See et al., 2017; Paulus et al., 2018). However, when users refer to the same article, their areas of focus can vary significantly (Zhong et al., 2021; Goyal et al., 2022; Zhang et al., 2023b). Rather than an overall summary, they are often more interested in obtaining summaries focused on specific aspects (Yang et al., 2023; Takeshita et al., 2024; Guo and Vosoughi, 2024). Therefore, generating aspect-based summaries to meet diverse user preferences is a natural and important capability for modern summariza-

<sup>&</sup>lt;sup>1</sup>We focus on abstracts because they are always publicly accessible and typically include the key medical aspects.

065

087 088

088 089 090

09

09

096

099 100

101 102

103 104

105

# 106

106 107 108

109

110

tion systems (Xu et al., 2023; Kolagar and Zarcone, 2024; Takeshita et al., 2024).

Moreover, most current studies in this field (Zhang et al., 2023a,b; Takeshita et al., 2024) focus on unidirectional summarization with LLMs (i.e., article  $\Rightarrow$  summary). Despite their potential, stateof-the-art LLMs still struggle with factual inaccuracies (Mallen et al., 2023; Min et al., 2023), which pose significant risks when healthcare professionals rely on these summaries for treatment decisions (Burns et al., 2011; Xie et al., 2024). By providing referenced source texts from which summaries are derived (i.e., *article*  $\leftarrow$  *summary*), users can more easily locate relevant context and verify the generated content, thereby mitigating such concerns (Kambhamettu et al., 2024; Xie et al., 2024; Deng et al., 2024). Therefore, traceable summarization (i.e., *article*  $\Leftrightarrow$  *summary*) becomes especially crucial given that summarization systems can generate hallucinated content (Dhuliawala et al., 2024).

To address these two concerns, we introduce TRACSUM, a novel summarization task that generates structured summaries of clinical articles across seven key medical aspects, as shown in Figure 1. These structured summaries not only provide flexibility to meet diverse informational needs but also enable cross-study comparisons, supporting a more comprehensive synthesis of evidence for clinical decision-making. In addition, TRACSUM extends the task by identifying the sentences cited by the summary. In real-world scenarios, this sentencelevel traceable summarization enables users to locate the relevant context and verify the generation. Overall, our key contributions are as follows:

**Contribution 1:** We propose TRACSUM, a novel benchmark for generating structured summaries of clinical articles across seven key aspects, enriched with sentence-level citations for each summary. To support this task, we construct a new dataset by annotating 500 clinical abstracts, resulting in 3.5K summary–citations pairs (§3).

**Contribution 2:** We introduce a fine-grained automatic evaluation framework tailored for this task, which assesses the completeness and consistency of the system output by measuring the recall and precision of both generated facts and their corresponding sentence-level citations (§4).

111Contribution 3: Inspired by Chain-of-thought112(CoT) reasoning (Wei et al., 2022), we propose113a summarization pipeline, TRACK-THEN-SUM,114which consists of a tracker  $\mathcal{T}$  and a summarizer

S. The tracker  $\mathcal{T}$  identifies source sentences relevant to a specific aspect, and the summarizer S condenses them into a short summary (§5).

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

**Contribution 4:** We evaluate a diverse set of closed- and open-source LLMs on TRACSUM, and conduct a human evaluation to assess the outputs produced by our fine-grained evaluation method. The findings demonstrate that TRACSUM can serve as an effective benchmark for traceable, aspectbased summarization in the medical domain (§6).

## 2 Related Work

### 2.1 Aspect-Based Summarization

Articles describing clinical trials often present information aligned with fixed core aspects, such as PICO<sup>2</sup> elements (Richardson et al., 1995; Schardt et al., 2007; Schiavenato and Chu, 2021), which represent essential components of medical evidence (Jin and Szolovits, 2018; Joseph et al., 2024). Generating structured summaries for these elements offers flexibility to address diverse informational needs and facilitates cross-study comparisons (Yang et al., 2023; Takeshita et al., 2024), enabling a comprehensive synthesis of evidence for clinical decision-making. To support fine-grained summarization, this work focuses on generating structured summaries that cover seven medical aspects commonly reported in clinical articles.

## 2.2 Traceable Summarization

Identifying the citations that summaries rely on can help users verify their accuracy (Gao et al., 2023; Xie et al., 2024), particularly in high-stakes domains such as medicine. To support critical examination of summaries and their underlying sources, Kambhamettu et al. (2024) introduced a simple interaction primitive called "traceable text." In the domain of Question Answering (QA), Gao et al. (2023) showed that enabling LLMs to generate text with passage-level citations improves factual correctness and verifiability. Moreover, several studies on retrieval-augmented generation (RAG) approaches can support document- or paragraph-level traceability (Wang et al., 2024b; Xu et al., 2024; Wang et al., 2024a). Building on this prior work, our research introduces sentence-level traceability of summaries generated by summarization systems, allowing users to directly inspect the source content that supports each summarized aspect.

<sup>&</sup>lt;sup>2</sup>PICO: Participants/Problem (P), Intervention (I), Comparison (C), and Outcome (O).

## **3 TRACSUM Benchmark**

#### 3.1 Task Description

162

163

164

165

166

167

168

170

171

172

173

174

175

176

177

178

179

181

182

183

184

185

186

187

190

191

192

193

194

195

196

197

198

199

Given a clinical article and a specific medical aspect, TRACSUM requires summarization systems to generate an aspect-based summary along with the corresponding sentence-level citations from which the summary is derived. Formally, let the input article  $d = [c_1, c_2, ..., c_n]$  be a sequence of uniquely indexed sentences, and let a be a target aspect selected from predefined aspects  $\mathcal{A}$  (§3.2.1). The system  $\mathcal{M}(\mathcal{C}', sum' \mid d, a)$  is expected to generate an aspect-specific summary sum' and a set of cited sentences  $\mathcal{C}' = [c'_1, c'_2, ..., c'_k]$ , where  $c'_i$  refers to the index of a sentence in d that supports the summary. If the article contains no information relevant to the given aspect, the system should output  $sum' \leftarrow$  "Unknown" and  $\mathcal{C}' \leftarrow$  "Null".

#### **3.2 Dataset Collection**

## 3.2.1 Medical Aspects

Building on the PICO framework (§2.1), we define  $\mathcal{A}$  as a set of seven medical aspects commonly reported in clinical articles (as listed in Table 1).

Symbol	Aspect Description	
А	Aims	Objective
Ι	Intervention	Treatment Method
0	Outcomes	Results of Predefined Variables
Р	Participants	E.g., Diseases, Number
М	Medicine	E.g., Name, Dosage
D	Duration	Treatment Duration
S	Side Effects	Observed Adverse Events

Table 1: Definition of seven medical aspects.

## 3.2.2 Source Articles

We initially screened 741 medical abstracts from PubMed<sup>3</sup>, of which 500 were ultimately included. The screening criteria were as follows: (1) the study focuses on melanoma; (2) the publication date is within the past 10 years; (3) the article is written in English; (4) the study is classified as either a Clinical Trial or a Randomized Controlled Trial; and (5) the article is published in a journal ranked in Q1 or Q2 according to the Journal Citation Reports (JCR) (Clarivate Analytics, 2024).

#### 3.2.3 Initial Generation With Mistral Large

Manual dataset annotation is often costly and susceptible to stylistic inconsistencies. Consequently, leveraging LLMs to generate supervised datasets has gained popularity due to their strong zero-shot performance (Chen et al., 2024; Asai et al., 2024). In this work, we automatically constructed a draft dataset by prompting Mistral Large (Mistral AI, 2025) to summarize 500 included abstracts, resulting in 3.5K summary–citations pairs, which were subsequently evaluated by human experts using three qualitative metrics (§3.2.4). The prompt structure comprises an abstract, a target aspect, and a type-specific instruction, followed by two demonstration examples. If the abstract lacks relevant information for the specified aspect, the model is instructed to return "Unknown" without generating any alternative response. An example of prompt templates is illustrated in Table 15 in §G. 200

201

202

203

204

205

207

209

210

211

212

213

214

215

216

217

218

219

220

221

222

224

225

226

227

228

230

231

232

233

234

235

236

237

238

239

240

241

#### 3.2.4 Annotation Process

We recruited six annotators, including three medical students and three NLP researchers, who were compensated in accordance with minimum wage standards in Germany. The annotation process was carried out in two phases. In the first phase, annotators independently evaluated all data instances. In the second phase, data instances that received lower evaluation scores were manually revised. The full annotation guideline is described in §A.

**Phase I: Evaluation.** To ensure consistency in writing style, each data instance was independently evaluated by two independent annotators, one from the medical domain and one from the NLP domain. The annotators assessed each data instance using three qualitative evaluation metrics (as shown in Table 2) on a 5-point Likert scale, as detailed in §A.4. Evaluating a single article typically takes 10–15 minutes, depending on its complexity.

Metric	Description
Completeness	Does the generated summary include all facts for the given aspect?
Conciseness	Does the generated summary include any irrelevant or erroneous information?
Traceability	Do the citations accurately and sufficiently ground the generated summary?

Table 2: Qualitative evaluation metrics.

**Phase II: Revision.** Out of the 3.5K evaluated data instances, we filtered out 741 (21%) that required further revision. The filtering criteria were as follows: (1) the mean score for any of the three evaluation metrics was below 3.5, or (2) the score difference between annotators exceeded 2.0. Annotators were then instructed to revise both the summaries and their corresponding citations, as illustrated in Figure 8 in §A.

<sup>&</sup>lt;sup>3</sup>https://pubmed.ncbi.nlm.nih.gov/



Figure 2: Human evaluation results (5-point scale) across three qualitative metrics for the seven medical aspects. Completeness and Conciseness for summary evaluation, and Traceability for citation evaluation.

#### 3.3 Quality Analysis

242

243

245

247

249

250

254

261

262

264

265

267

269

270

271

274

276

277

281

To analyze the dataset's quality, we conducted a statistical analysis of the human evaluation results. Before filtering, the scores across all aspects and metrics are generally above 4.0 (as shown in Figure 2), indicating high overall quality. Of the 741 (21%) filtered instances, 197 concern the O (Outcomes), 174 the I (Intervention), and 171 the D (Duration), suggesting that Mistral Large's summaries diverge most from human judgment on these three aspects, possibly due to the relatively complex information in the source texts. To assess interannotator agreement (IAA), we report exact match accuracy, within-one accuracy, and mean absolute error, following prior work (Attali and Burstein, 2006; Zhang and Zhou, 2007). The statistical analysis revealed high agreement under the within-one accuracy metric (84.9%), despite a lower exact match accuracy (66.6%) and a mean absolute error of 0.56, indicating acceptable consistency with only minor scoring discrepancies.

3.4 Characteristics of the Dataset

Among the 500 abstracts, the average length is 319.89 tokens, with abstract lengths ranging from 25 to 1,104 tokens. Each abstract contains an average of 10.42 sentences, spanning from 1 to 32. In the dataset of 3.5K data instances, 2,862 are positive and 638 are negative<sup>4</sup>. The positive summaries average 28.06 tokens in length, with a range from 3 to 77 tokens. On average, each positive summary cites 1.78 sentences, with a range from 1 to 7. Example data instances are presented in Table 14 (see §F), and more characteristics are described in §B.

## 4 Automatic Evaluation Framework

Clinical texts have two essential characteristics: (1) *it must be entirely complete, with no omissions* and (2) *it must be fully accurate, without any errors* (Gao et al., 2023; Xie et al., 2024). In line with these considerations, we propose a finegrained evaluation framework for this new task by extending the methodology of Xie et al. (2024) and Gao et al. (2023), which evaluate completeness (§4.1) and conciseness (§4.2) of generated content through a suite of metrics, as illustrated in Figure 3. Unlike their original definitions, our approach incorporates citation recall and precision to evaluate completeness and conciseness. Before computing these metrics, we first check whether the cited sentences entail the generated summary.

282

283

284

285

287

290

292

294

295

296

297

298

299

300

301

302

303

304

305

306

307

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

#### 4.1 Completeness Evaluation

Building on characteristic (1) of clinical texts, we evaluate completeness — the extent to which clinically significant information is preserved in the system output. Unlike previous work (Van Veen et al., 2023), which assigns an overall score, our approach emphasizes identifying which specific salient information is retained or omitted. As described in §3.1, TRACSUM requires a summarization system to produce both a summary and its associated citations. To evaluate completeness, we introduce claim recall to assess summary content and citation recall to assess citation coverage.

**Claim Recall:** Following DOCLENS (Xie et al., 2024), we decompose each reference into a list of atomic subclaims using a decomposition model, where each subclaim represents a single factual statement from the reference. Let y denote the reference,  $\mathcal{L}_y$  the set of reference subclaims, and y' the system-generated summary. We employ a natural language inference (NLI) model to evaluate whether each subclaim  $l \in \mathcal{L}_y$  is entailed by y'. Claim recall is computed as  $\frac{1}{|\mathcal{L}_y|} \sum_{l \in \mathcal{L}_y} \mathbb{I}[y' \Rightarrow l]$ , where  $\mathbb{I}[y' \Rightarrow l]$  is an indicator function that returns 1 if y' entails l, and 0 otherwise.

**Citation Recall:** In contrast to previous approaches (Gao et al., 2023; Liu et al., 2023; Xie et al., 2024), which consider citations valid if the cited sentences collectively support the summary, our method assesses whether each cited sentence independently supports the output. Let C be the set of citations in the reference and C' the set in the system output. A citation is considered recalled if it satisfies the following two conditions: (1) the

<sup>&</sup>lt;sup>4</sup>Negative samples correspond to cases where both the summary and citation content are null.



Figure 3: Overview of the automatic evaluation framework. Completeness is assessed using Claim Recall and Citation Recall, while conciseness is measured by Claim Precision and Citation Precision. Decom. denotes the claim decomposition model, and Eval. refers to the entailment evaluator.

cited sentence supports the generated summary  $(c \rightarrow y')$ ; and (2) the citation is present in the reference  $(c \in C)$ . Citation recall is formally defined as  $\frac{1}{|C|} \sum_{c \in C'} \mathbb{I}[c \in C \land c \rightarrow y']$ .

## 4.2 Conciseness Evaluation

In line with characteristic (2), an ideal system output should avoid redundant or incorrect information. We evaluate conciseness as the proportion of generated content that is both factually accurate and salient. To this end, we use two metrics: claim precision, which assesses the informativeness and factual accuracy of the summary, and citation precision, which captures citation redundancy.

**Claim Precision:** Analogous to claim recall, we first decompose the generated summary into a list of subclaims, then use an evaluator to compute the proportion of these subclaims that are entailed by the reference. Claim precision is defined as  $\frac{1}{|\mathcal{L}'_y|} \sum_{l \in \mathcal{L}'_y} \mathbb{I}[y \Rightarrow l], \text{ where } \mathcal{L}'_y \text{ denotes the set of subclaims extracted from the generated summary.}$ 

**Citation Precision:** To assess whether the output includes unnecessary citations, we introduce citation precision. In line with citation recall, a citation is deemed valid if it satisfies both previously defined conditions ( $c \in C \land c \rightarrow y'$ ). Citation precision is then calculated as the proportion of system-generated citations that fulfill these criteria.

#### **5** Baseline Method

In this section, we introduce our baseline method, TRACK-THEN-SUM (TTS), which consists of a Algorithm 1: TRACK-THEN-SUM InferenceRequire: Tracker  $\mathcal{T}$ , Summarizer  $\mathcal{S}$ Input: article  $d = \{c_1, c_2, ..., c_n\}$  and aspect  $a \in \mathcal{A}$ Output: summary sum and its citations  $\mathcal{C}'$ 1:  $\mathcal{C}' \leftarrow \emptyset$ ;2: foreach  $c \in \{c_1, c_2, ..., c_n\}$ 3:  $\mathcal{T}$  predict relevance given (a, c);4: if relevance == Yes then append c to  $\mathcal{C}'$ ;5: summary sum  $\leftarrow \mathcal{S}(a, \mathcal{C}')$  or  $\mathcal{S}(a, (\mathcal{C}' \oplus f.))$ ;

Algorithm 1: TRACK-THEN-SUM inference process.

tracker  $\mathcal{T}$  and a summarizer S (available in two variants), as illustrated in Figure 10 in §C. The training procedure is detailed in §C.1.

#### 5.1 Inference Overview

The TRACK-THEN-SUM generation pipeline contains two phases: tracking and summarization. In the first phase,  $\mathcal{T}$  identifies the sentences most relevant to the given aspect. In the second phase,  $\mathcal{S}$ generates a concise summary based on the selected sentences. Finally, the summary and citations are merged into the output, as shown in Algorithm 1.

#### 5.2 Tracker T

5

**Data Collection:** We first applied sentence tokenization to each abstract in the training set. For each sentence, we generated (c, a) pairs by combining it with every predefined aspect  $a \in A$ . Each pair was labeled with a binary variable y based on the corresponding *citations* field: if the sentence index appeared in the *citations* associated with aspect a, we assigned y = 1; otherwise, y = 0. The resulting training dataset is denoted as  $D_T$ .

**Training:** Given the constructed dataset  $\mathcal{D}_{\mathcal{T}}$ , we initialized tracker  $\mathcal{T}$  using a pre-trained language

325

326

345

351

354

367

368

369

370

371

374

375

377

355

357

359

361

363

378model (LM) as the backbone. The model was sub-<br/>sequently fine-tuned on  $\mathcal{D}_{\mathcal{T}}$  using a standard binary<br/>classification objective which maximizes the log-<br/>likelihood of the observed labels:

$$\max_{\mathcal{T}} \mathbb{E}_{((c,a),y) \sim \mathcal{D}_{\mathcal{T}}} \log p_{\mathcal{T}}(y \mid (c,a))$$

#### **5.3** Summarizer S

387

390

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

**Data Collection:** For each summary *sum* in the training set, we extracted related sentences from the abstract based on the *citations* field to form the set C. Each C was paired with its associated aspect a, and combined with the *sum* to form ((C, a), sum). The resulting training dataset is denoted as  $\mathcal{D}_S$ .

**Training:** Similar to the training of  $\mathcal{T}$ , we initialized summarizer S using a pre-trained LM as the backbone. We then fine-tuned summarizer S on  $\mathcal{D}_S$ using a standard next-token prediction objective, which maximizes the likelihood of generating the target summary *sum* given the input (C, *a*) pair:

$$\max_{\mathcal{C}} \mathbb{E}_{((\mathcal{C},a),sum) \sim \mathcal{D}_{\mathcal{S}} \log p_{\mathcal{S}}(sum|C,a)}$$

To investigate the impact of incorporating full context (denoted as f.), we trained a variant S that generates a summary sum given the input ( $C \oplus f$ ., a).

#### 6 Experiment

In this section, we aim to address the following research questions: **RQ1:** How effective is TRACSUM as a benchmark for evaluating LLMs in aspect-based summarization with sentence-level traceability? **RQ2:** To what extent does the proposed evaluation method align with human judgment, and what role does the evaluator play in this process? **RQ3:** Which factors most significantly impact the accuracy of traceable summarization? To address these questions, we begin by conducting a preliminary evaluation of several LLMs, including both proprietary models (e.g., GPT-40 (Hurst et al., 2024)) and open-source models (e.g., LLaMA-3.1 (Grattafiori et al., 2024), Mistral (Jiang et al., 2024), and Gemma-3 (Team et al., 2025)).

#### 6.1 Experimental Setting

417Data Preparation: The TRACSUM dataset was418randomly split into training and test sets with an4198:2 ratio. We examined the distribution of samples420in the test set across the seven predefined aspects,421along with the proportion of positive and negative422instances for each, as shown in Figure 4. The re-423sults show that while nearly all abstracts contain



Figure 4: Distribution of test data across seven aspects.

information related to Aims (A), Intervention (I), and Outcomes (O), only 31% explicitly mention the Duration (D) aspect. The baseline model was fine-tuned on the training set, and both the baseline and LLMs were evaluated on the test set.

**Backbone Model Selection:** The TRACK-THEN-SUM (TTS) pipeline comprises two components (Tracker  $\mathcal{T}$  and Summarizer S) that can be initialized with any pre-trained LM. For consistency and ease of deployment, we adopt Llama-3.1-8B (Dubey et al., 2024) as the backbone for both components, with the training details provided in §C.1.

**LLMs and Prompt Setting:** We selected several widely used instruction-following LLMs for evaluation, as listed in Table 3. All models were evaluated using a two-shot prompting strategy, with each prompt containing one positive and one negative example. To ensure consistency, each model was prompted using its official input format with identical content (see Table 15 in §G), and a fixed temperature of 1.0 was used across all generations. Larger models were accessed via their official APIs, incurring additional usage costs (see §D).

Algorithm 2: Computation Process of Evaluation Metrics
<b>Require:</b> decomposition model: $\mathcal{E}$ , NLI model: $\phi$
<b>Input:</b> system output $(sum', C')$ , reference $(sum, C)$
Output: CLR, CIR, CLP, CIP
1: $\{s_1, s_2, \dots, s_n\} \leftarrow \mathcal{E}(sum); 0 \leftarrow n;$
2: foreach $s_i \in \{s_1, s_2,, s_n\}$
3: if $\phi(sum', s_i) == 1$ then $n++$ ;
4: <b>CLR</b> $\leftarrow n/ \{s_1, s_2,, s_n\} $
5: $0 \leftarrow n;$
6: foreach $c'_i \in \mathcal{C}'$
7: foreach $s'_i \in \{s'_1, s'_2,, s'_n\}$
8: if $\phi(c'_i, s'_i) == 1$ then $n++$ ; break;
9: $\mathbf{CIR} \leftarrow n/ \mathcal{C} ; \mathbf{CIP} \leftarrow n/ \mathcal{C}' ;$
10: $\{s'_1, s'_2,, s'_n\} \leftarrow \mathcal{E}(sum'); n \leftarrow 0;$
11: foreach $s'_i \in \{s'_1, s'_2,, s'_n\}$
12: if $\phi(sum, s'_i) == 1$ then $n++;$
13: <b>CLP</b> $\leftarrow n/ \{s'_1, s'_2,, s'_n\} ;$

Algorithm 2: Computation process of evaluation metrics. **CLR:** Claim Recall. **CIR:** Citation Recall. **CLP:** Claim Precision. **CIP:** Citation Precision.

**Evaluation Setting:** In the preliminary experiment, we adopt Mistral Large (Mistral AI, 2024) as the decomposition model  $\mathcal{E}$ , which is used to break

447 448 449

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

down both the system-generated and reference summaries into a set of atomic subclaims. For the entailment evaluation, we utilize TRUE (Honovich et al., 2022) as the evaluator  $\phi$ . Let  $\phi(p, h)$  denote the output of the NLI model, where the value is 1 if the premise p entails the hypothesis h, and 0 otherwise. The computation process of the evaluation metrics is presented in Algorithm 2.

#### 6.2 Preliminary Results

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

481

482

483

484

485

486

487

488

489

490

491

492

493

494

**Comparison of LLMs:** Table 3 shows the evaluation results of various LLMs along with our proposed method (in two variants). We observe the following: (1) Larger open-source models (e.g., LLaMA-3.1-70B, Mistral-8x7B) consistently outperform smaller ones across all metrics. (2) Proprietary models like GPT-40 and GPT-40-mini also perform well, with only small differences between them. (3) Our proposed method, fine-tuned from LLaMA-3.1-8B, shows clear improvements over both the base model and other LLMs, particularly on the two citation-based metrics CIR and CIP ( $\geq$  74.0%), demonstrating their strength in identifying supporting source sentences.

**Performance on Completeness and Conciseness:** 473 As shown in Table 3, LLMs generally perform bet-474 475 ter on completeness than on conciseness, suggesting a tendency to generate content that exceeds the 476 scope of the reference data. This may be due to 477 full context visibility during generation, which can 478 cause the models to include content only loosely 479 related to the target aspects. 480

**Does Full Context Help?** In the TTS pipeline, we extend the input to the summarizer S by including not only the tracked sentences but also the full context (i.e., the abstract). This modification allows the TTS  $\oplus$  f. variant to improve the claim recall CLR ( $67.1\% \rightarrow 79.8\%$ ) of the generated summaries without substantially compromising performance on other metrics. With the tracker T output unchanged, the observed gains may stem from the full context offering useful explanations for abbreviations or domain-specific terminology, thereby helping S better interpret the tracked sentences. A detailed case analysis is provided in §E.1.

#### 6.3 Agreement with Human Evaluation

495To address first sub-question of RQ2, we conducted496a human evaluation and measured the agreement be-497tween human judgments and the automatic evalua-498tion scores produced by the NLI model (TRUE) us-

	Completeness		Conciseness		F1 Score	
Method	CLR	CIR	CLP	CIP	$F_1^{\text{cl.}}$	$F_1^{\text{ci.}}$
Llama-3.1-8B	59.2	62.5	63.6	54.8	61.3	58.4
Llama-3.1-70B	74.7	77.9	71.3	67.7	72.9	72.4
Mistral-7B	59.1	59.5	55.5	48.4	57.4	53.4
Mistral-8x7B	61.1	62.1	58.9	58.4	60.0	60.2
Gemma3-12B	62.8	66.0	58.3	55.3	60.5	60.2
Gemma3-27B	64.6	66.4	57.7	59.6	61.0	63.0
GPT-40	74.0	78.2	66.2	63.8	69.9	70.3
GPT-4o-mini	67.8	76.0	<u>67.6</u>	<u>68.4</u>	67.7	72.0
TTS	67.1	76.2	<u>68.4</u>	77.0	67.8	76.6
$TTS \oplus f.$	79.8	74.6	67.2	<u>75.0</u>	73.0	<u>74.8</u>

Table 3: Preliminary evaluation results (%). **Bold** values indicate the best performance in each metric, <u>underlined</u> values indicate the second-best, and <u>wave underlined</u> values indicate the third-best.  $\oplus$  *f*. denotes the configuration where the full context is concatenated to the input of the summarizer S.  $F_1^{\text{cl.}}$  and  $F_1^{\text{ci.}}$  represent the F1 scores for claim and citation prediction, respectively.

ing Spearman's correlation coefficient ( $\rho$ ) (Kendall and Gibbons, 1990) and Pearson's correlation coefficient (r) (Sheskin, 2003). We randomly sampled ten abstracts from the test set, and the annotator followed the procedure in Algorithm 2 to evaluate outputs from our TTS  $\oplus$  f., as shown in Table 4. The results show an average Spearman's  $\rho = 0.612$ and Pearson's r = 0.577, indicating a moderate positive correlation between automatic evaluation and human judgments. This suggests that our proposed evaluation framework aligns reasonably well with human assessments, while still leaving room for improvement. A detailed comparison of the final evaluation results is provided in §E.2.

<b>Keletence.</b> Subclaims $\neg$ Citations $\rightarrow$ 1, 5
1. The study included 533 patients.
2. The patients were treatment-naive.
3. The patients had unresectable stage III-IV melanoma.
<b>TTS</b> $\oplus$ <i>f</i> . <b>Output:</b> Subclaims $\neg$ Citations $\rightarrow$ 1, 3, 5
<ul><li>1'. The study involved treatment-naive patients.</li><li>2'. The patients had unresectable stage III-IV melanoma.</li><li>3'. 533 patients received nivolumab plus ipilimumab.</li></ul>
NLI: $reference \rightarrow s1', s2' \checkmark \rightarrow s3' \bigstar$ CLR: 66.7% Human: $reference \rightarrow s1', s2', s3' \checkmark$ CLR: 100% Become "523 patiente" is found in the reference

Table 4: An example comparing automatic and human evaluation of claim recall (PMID: 37307514, Aspect: Patients).

#### 6.4 Aspect-Wise Performance Analysis

To analyze the performance of the TTS  $\oplus$  f. variant across the seven aspects, we grouped the data514ant across the seven aspects, we grouped the data515by aspect and computed the four evaluation metrics516for each group, as shown in Table 5. We observed517substantial variation in the model's performance518across different aspects. Notably, aspects O (Outcomes) and I (Intervention) received lower scores520

510

511

512

513

499



Figure 5: Spearman ( $\rho$ ) and Pearson (r) correlations between evaluators and human scores across four metrics.

	Compl	eteness	Conci	seness	F1 S	core
Aspect	CLR	CIR	CLP	CIP	$F_1^{\text{cl.}}$	$F_1^{\text{ci.}}$
A	86.3	83.2	71.8	89.8	78.4	86.4
Ι	69.8	61.4	51.0	47.6	58.9	53.4
0	61.4	50.2	48.7	50.1	54.2	50.1
Р	87.7	78.4	80.2	84.2	83.7	81.3
Μ	85.4	71.9	75.1	73.3	79.9	72.6
D	92.2	93.6	81.4	93.2	86.4	93.3
S	75.8	<u>83.5</u>	62.2	86.8	68.3	85.0
Avg.	79.8	74.6	67.2	75.0	73.0	74.8

Table 5: Aspect-wise performance of method TTS  $\oplus f$ .

across all four evaluation metrics, likely because the corresponding abstracts often contain a large number of relevant sentences, making precise extraction more challenging. In contrast, aspect D (Duration) achieved relatively higher scores, possibly due to the fact that 69% of its test instances are negative cases (i.e., both the summary and citation are null), which simplifies the task and makes correct predictions easier for the model.

#### 6.5 Ablation Studies

Comparison of Entailment Evaluators: To address the second sub-question of RQ2, we experiment with two additional instruction-following LLMs as entailment evaluators: the proprietary GPT-40 (Hurst et al., 2024) and the open-source Mistral-Large (Mistral AI, 2025). Building on the experimental setup described in  $\S6.3$ , we replace the TRUE model with each of these evaluators to assess the outputs generated by the TTS  $\oplus f$ . variant. The experiment procedure and results are described in §E.3. We then compute Spearman's  $\rho$  and Pearson's r to quantify their agreement with human judgments in four metrics, as presented in Figure 5. Our findings reveal that: (1) both GPT-40 ( $\rho = 0.80; r = 0.77$ ) and Mistral-Large ( $\rho = 0.71; r = 0.70$ ) show substantially stronger alignment with human judgments compared to TRUE ( $\rho = 0.61; r = 0.57$ ); and (2) GPT-40 achieves a higher correlation with human judgments than Mistral-Large. We found that GPT-40 is better at understanding abbreviations. For instance, it correctly infers that the reference "50 participants were randomized: 23 to observation and 27 to radiation therapy" entails the subclaim "27 participants were assigned to the RT group",

whereas Mistral and TRUE do not.

The Effect of Tracking Order: To address RO3, we design two variants by modifying the position of the tracker  $\mathcal{T}$ : (i) SUM-THEN-TRACK (STT) places  $\mathcal{T}$  after the summarizer  $\mathcal{S}$ , where  $\mathcal{S}$  first generates an aspect-based summary, and  $\mathcal{T}$  then retrieves source sentences relevant to that summary; (ii) END-TO-END (ETE) removes the tracker entirely and fine-tunes a single model  $\mathcal{M}$  to generate both summary and citations. The experimental procedures are detailed in §C. We evaluated STT and ETE on the test set, with results shown in Table 6. We observe that: (1) removing the tracker results in a decline in citation-based performance, highlighting the importance of explicit sentence tracking; and (2) while STT improves claim recall, it performs worse on other metrics, likely due to its dependence on pre-generated summaries, which may introduce noise or inaccuracies. These findings emphasize the importance of incorporating tracking early in the summarization process.

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

584

585

586

587

589

590

591

592

593

	Compl	eteness	Conci	seness	F1 S	Score
Method	CLR	CIR	CLP	CIP	$F_1^{\text{cl.}}$	$F_1^{\text{ci.}}$
$TTS \oplus f.$	79.8	74.6	67.2	75.0	73.0	74.8
ETE	80.1	72.6	64.1	71.2	71.2	71.9
STT	81.2	62.2	58.1	66.4	67.7	64.1

Table 6: Comparison of the three tracking order variants.

## 7 Conclusion

Motivated by growing concerns over the factual accuracy of system-generated summaries in the medical domain, we present TRACSUM, a novel benchmark for aspect-based summarization that incorporates sentence-level citations. This enables users to trace source content and verify the factual consistency of generated information. Experimental results, which show strong alignment with human judgments, demonstrate that TRACSUM can serve as a reliable benchmark for assessing both the completeness and conciseness of summaries and their citations. Furthermore, we also observe that explicitly performing sentence-level tracking prior to summarization enhances generation accuracy, while incorporating the full context further improves summary completeness.

540

542

543

544

545

551

552

553

555

## Limitations

594

616

621

627

632

637

638

640

641

Our research marks a significant step toward evaluating sentence-level traceability in aspect-based 596 summarization. Nonetheless, it has certain limi-597 tations. 1). The dataset used in TRACSUM was initially generated by Mistral Large. While this approach helped reduce time and cost, it may also introduce model-specific biases. To address this concern, we implemented two mitigation strategies: (i) we conducted two rounds of human evaluation, followed by manual revision of samples with low scores or inconsistent annotations; and (ii) we excluded Mistral Large from the list of evaluated models to avoid unfair advantages or confirmation bias. 2). The structure and content of prompts can significantly influence the outputs of LLMs and, in turn, their evaluation scores. Although our prompt 610 template was designed to be general and broadly 611 applicable, it may not elicit the best performance from every model. To reduce potential bias and ensure fair comparison, we used a standardized 614 615 prompt format across all models.

#### References

- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. Self-RAG: Learning to retrieve, generate, and critique through self-reflection. In *The Twelfth International Conference on Learning Representations*.
- Yigal Attali and Jill Burstein. 2006. Automated essay scoring with e-rater® v. 2. *The Journal of Technology, Learning and Assessment*, 4(3).
  - Jacob Benesty, Jingdong Chen, Yiteng Huang, and Israel Cohen. 2009. Pearson correlation coefficient. *Noise reduction in speech processing*, pages 1–4.
  - Patricia B Burns, Rod J Rohrich, and Kevin C Chung. 2011. The levels of evidence and their role in evidence-based medicine. *Plastic and reconstructive surgery*, 128(1):305–310.
- Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. 2024. Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17754–17762.
- Jianpeng Cheng and Mirella Lapata. 2016. Neural summarization by extracting sentences and words. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 484–494, Berlin, Germany. Association for Computational Linguistics.
- Clarivate Analytics. 2024. Journal Citation Reports. Accessed: 2024-03-12.

Zhenyun Deng, Michael Schlichtkrull, and Andreas Vlachos. 2024. Document-level claim extraction and decontextualisation for fact-checking. In *Proceedings* of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 11943–11954, Bangkok, Thailand. Association for Computational Linguistics. 645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

689

690

691

692

693

694

695

696

697

698

699

- Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. 2024. Chain-of-verification reduces hallucination in large language models. In *Findings* of the Association for Computational Linguistics: ACL 2024, pages 3563–3578, Bangkok, Thailand. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Sameh Frihat and Norbert Fuhr. 2024. Supporting evidence-based medicine by finding both relevant and significant works. *arXiv preprint arXiv:2407.18383*.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023. Enabling large language models to generate text with citations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6465–6488, Singapore. Association for Computational Linguistics.
- Tanya Goyal, Nazneen Rajani, Wenhao Liu, and Wojciech Kryscinski. 2022. HydraSum: Disentangling style features in text summarization with multidecoder models. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 464–479, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The Ilama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Xiaobo Guo and Soroush Vosoughi. 2024. Disordered-DABS: A benchmark for dynamic aspect-based summarization in disordered texts. In *Findings of the Association for Computational Linguistics: EMNLP* 2024, pages 416–431, Miami, Florida, USA. Association for Computational Linguistics.
- Eduardo Hariton and Joseph J Locascio. 2018. Randomised controlled trials—the gold standard for effectiveness research. *BJOG: an international journal of obstetrics and gynaecology*, 125(13):1716.
- Janusz Hauke and Tomasz Kossowski. 2011. Comparison of values of pearson's and spearman's correlation coefficients on the same sets of data. *Quaestiones Geographicae*, 30(2):87–93.
- Or Honovich, Roee Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansy, Vered Cohen, Thomas

7	0	2
7	n	3
1	2	2
ſ	U	4
7	0	5
7	0	6
7	n	7
1	С	6
ſ	U	Q
_		~
ſ	U	9
7	1	0
7	1	1
7	1	2
<u>_</u>	i	2
1	ł	J
_		
1	ł	4
7	1	5
7	1	6
7	1	7
_	i	
1	1	ö
_		~
ſ	1	9
7	2	0
7	2	1
7	2	2
÷	0	2
ſ	2	J
_	0	л
[	2	4
7	2	5
7	2	6
7	2	7
-	-	0
_	_	0
7	2	9
7	$\sim$	$\sim$
1	3	U
′ 7	33	U 1
7	3	U 1
′ 7 7	3 3 3	0 1 2
7 7 7	333	0 1 2
7 7 7	3 3 3 3	1 2 3
7 7 7 7	3 3 3 3 3	0 1 2 3 4
'7 7 7 7 7	3 3 3 3 3 3 3	0 1 2 3 4 5
7 7 7 7 7	3 3 3 3 3 3	0 1 2 3 4 5
7 7 7 7 7 7	3 3 3 3 3 3 3 3	0 1 2 3 4 5 6
, 7 7 7 7 7 7 7 7	3 3 3 3 3 3 3 3	0 1 2 3 4 5 6 7
, 7 7 7 7 7 7 7 7 7 7	3 3 3 3 3 3 3 3 3 3 3	0 1 2 3 4 5 6 7 9
7 777777777777	3 3 3 3 3 3 3 3 3 3 3	0 1 2 3 4 5 6 7 8
7 7777 7777	33333333	0 1 2 3 4 5 6 7 8 0
, 7 7 7 7 7 7 7 7 7	3333333333	0 1 2 3 4 5 6 7 8 9
7 7777 7777 77777	3 3 3 3 3 3 3 3 3 3 4	0 1 2 3 4 5 6 7 8 9 0
, 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	3 3 3 3 3 3 3 3 3 4 4	0 1 2 3 4 5 6 7 8 9 0 1
, 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	333333333444	0 1 2 3 4 5 6 7 8 9 0 1 2
, 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	3 3 3 3 3 3 3 3 3 3 3 4 4 4 4 4	U1234567890123
, 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	3 3 3 3 3 3 3 3 3 3 3 4 4 4 4 4	01 2345 678 90123
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	3 3 3 3 3 3 3 3 3 3 3 3 3 4 4 4 4 4 4	01 2345 678 901234
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	3 3 3 3 3 3 3 3 3 3 3 3 3 4 4 4 4 4 4 4	01 2345 678 9012345
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	3 3 3 3 3 3 3 3 3 3 3 3 3 3 4 4 4 4 4 4	01 2345 678 901 2345
7 777777777777777777777777777777777777	3 3 3 3 3 3 3 3 3 3 3 3 3 3 4 4 4 4 4 4	01 2345 678 901 2345 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	3 3 3 3 3 3 3 3 3 3 4 4 4 4 4 4 4 4	01 2345 678 9012345 67
,7777777777777777777777777777777777777	3 3 3 3 3 3 3 3 3 4 4 4 4 4 4 4 4	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	3 3 3 3 3 3 3 3 3 3 3 4 4 4 4 4 4 4 4 4	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 0
	3 3 3 3 3 3 3 3 3 3 3 4 4 4 4 4 4 4 4 4	01 2345 678 9012345 67892
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	3 3 3 3 3 3 3 3 3 4 4 4 4 4 4 4 4 4 4 4	01 2345 678 9012345 67890
	3 3 3 3 3 3 3 3 3 4 4 4 4 4 4 4 4 5	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	3 3 3 3 3 3 3 3 3 3 4 4 4 4 4 4 4 5 5	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 7 8 9 0 1 2 3 4 5 7 8 9 0 1 2 3 4 5 7 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 3 4 5 5 7 8 1 2 3 1 2 3 4 5 7 8 9 0 1 2 3 4 5 1 2 3 4 5 1 2 3 4 5 1 2 3 4 5 1 2 3 4 5 1 2 3 4 5 5 7 8 9 0 1 2 3 4 5 1 2 3 4 5 1 2 3 1 2 3 1 2 3 4 5 1 2 3 4 5 1 2 3 4 5 1 2 3 4 5 1 2 3 1 2 3 1 1 2 3 4 5 5 1 2 3 1 1 2 3 1 2 3 1 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 1 2
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7	33333333334444444445555	01 2345 678 9012345 67890 12345 67890
7 7777 777 7777777777777777777777777777	333333333444444444444455555	01 2345 678 9012345 67890 123
7 7777 777 7777777777777777777777777777	333333333344444444445555555555555555555	01 2345 678 9012345 67890 1234
777777777777777777777777777777777777777	333333333444444445555555555555555555555	01 2345 678 9012345 67890 12345
7 7777 777 7777777777777777777777777777	333333333444444455555555555555555555555	0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 7 8 9 0 1 2 3 4 5 7 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 8 9 0 1 2 3 4 5 1 2 3 4 5 9 0 1 2 3 4 5 1 2 3 4 5 1 2 3 4 5 1 2 3 2 3 4 5 5 7 8 9 0 1 2 3 4 5 5 7 8 9 0 1 2 3 4 5 1 2 3 1 2 3 1 2 3 1 1 2 3 4 5 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 1 2 3 2 3
7 7777 777 77777777777777777777777777	3 3 3 3 3 3 3 3 4 4 4 4 4 4 4 4 4 5 5 5 5	01         2345         678         9012345         67890         123456
7 7777 77777777777777777777777777777777	3 3 3 3 3 3 3 3 4 4 4 4 4 4 4 4 4 5 5 5 5	01 2345 678 9012345 67890 1234567

701

Scialom, Idan Szpektor, Avinatan Hassidim, and Yossi Matias. 2022. TRUE: Re-evaluating factual consistency evaluation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3905–3920, Seattle, United States. Association for Computational Linguistics.

- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. 2024. Gpt-40 system card. *arXiv preprint arXiv:2410.21276*.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. arXiv preprint arXiv:2401.04088.
- Di Jin and Peter Szolovits. 2018. PICO element detection in medical text via long short-term memory neural networks. In *Proceedings of the BioNLP 2018 workshop*, pages 67–75, Melbourne, Australia. Association for Computational Linguistics.
- Sebastian Joseph, Lily Chen, Jan Trienes, Hannah Göke, Monika Coers, Wei Xu, Byron Wallace, and Junyi Jessy Li. 2024. FactPICO: Factuality evaluation for plain language summarization of medical evidence. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 8437–8464, Bangkok, Thailand. Association for Computational Linguistics.
- Hita Kambhamettu, Jamie Flores, and Andrew Head. 2024. Traceable text: Deepening reading of aigenerated summaries with phrase-level provenance links. *arXiv preprint arXiv:2409.13099*.
- Maurice G. Kendall and Jean Dickinson Gibbons. 1990. *Rank Correlation Methods*, 5th edition. Oxford University Press, New York.
- Zahra Kolagar and Alessandra Zarcone. 2024. Hum-Sum: A personalized lecture summarization tool for humanities students using LLMs. In Proceedings of the 1st Workshop on Personalization of Generative AI Systems (PERSONALIZE 2024), pages 36–70, St. Julians, Malta. Association for Computational Linguistics.
- Nelson Liu, Tianyi Zhang, and Percy Liang. 2023. Evaluating verifiability in generative search engines. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7001–7025, Singapore. Association for Computational Linguistics.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023.
  When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.

Iain James Marshall, Veline L'Esperance, Rachel Marshall, James Thomas, Anna Noel-Storr, Frank Soboczenski, Benjamin Nye, Ani Nenkova, and Byron C Wallace. 2021. State of the evidence: a survey of global disparities in clinical trials. *BMJ global health*, 6(1):e004145. 759

760

763

765

766

767

768

771

774

775

779

782

783

784

785

787

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

- Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 12076–12100, Singapore. Association for Computational Linguistics.
- Mistral AI. 2024. Mistral large. Accessed: 2025-05-02.
- Mistral AI. 2025. Mistral: Introducing the large language model 2407. Accessed: 2025-02-10.
- Shashi Narayan, Shay Cohen, and Maria Lapata. 2018. Don't give me the details, just the summary! topicaware convolutional neural networks for extreme summarization. In 2018 Conference on Empirical Methods in Natural Language Processing, pages 1797–1807. Association for Computational Linguistics.
- Romain Paulus, Caiming Xiong, and Richard Socher. 2018. A deep reinforced model for abstractive summarization. In *International Conference on Learning Representations*.
- W Scott Richardson, Mark C Wilson, Jim Nishikawa, and Robert S Hayward. 1995. The well-built clinical question: a key to evidence-based decisions. *ACP journal club*, 123(3):A12–3.
- Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive sentence summarization. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 379–389, Lisbon, Portugal. Association for Computational Linguistics.
- David L Sackett. 1997. Evidence-based medicine. In *Seminars in perinatology*, volume 21, pages 3–5. Elsevier.
- Connie Schardt, Martha B Adams, Thomas Owens, Sheri Keitz, and Paul Fontelo. 2007. Utilization of the pico framework to improve searching pubmed for clinical questions. *BMC medical informatics and decision making*, 7:1–6.
- Martin Schiavenato and Frances Chu. 2021. Pico: What it is and what it is not. *Nurse education in practice*, 56:103194.
- Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get to the point: Summarization with pointergenerator networks. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1073– 1083, Vancouver, Canada. Association for Computational Linguistics.

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

872

873

874

875

David J. Sheskin. 2003. *Handbook of Parametric and Nonparametric Statistical Procedures*. Chapman and Hall/CRC.

815

816

817

818

819

820

833

838

844

847

849

851

852

853

854

859

861

870

- Sotaro Takeshita, Tommaso Green, Ines Reinig, Kai Eckert, and Simone Ponzetto. 2024. ACLSum: A new dataset for aspect-based summarization of scientific publications. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 6660–6675, Mexico City, Mexico. Association for Computational Linguistics.
- Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. 2025. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*.
  - Dave Van Veen, Cara Van Uden, Louis Blankemeier, Jean-Benoit Delbrouck, Asad Aali, Christian Bluethgen, Anuj Pareek, Malgorzata Polacin, Eduardo Pontes Reis, Anna Seehofnerova, et al. 2023. Clinical text summarization: adapting large language models can outperform human experts. *Research square*, pages rs–3.
- Fei Wang, Kaiqiang Song, Hongming Zhang, Lifeng Jin, Sangwoo Cho, Wenlin Yao, Xiaoyang Wang, Muhao Chen, and Dong Yu. 2022. Salience allocation as guidance for abstractive summarization. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 6094–6106.
- Weixuan Wang, Barry Haddow, and Alexandra Birch. 2024a. Retrieval-augmented multilingual knowledge editing. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 335–354, Bangkok, Thailand. Association for Computational Linguistics.
- Zheng Wang, Shu Teo, Jieer Ouyang, Yongjun Xu, and Wei Shi. 2024b. M-RAG: Reinforcing large language model performance through retrieval-augmented generation with multiple partitions. In *Proceedings* of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1966–1978, Bangkok, Thailand. Association for Computational Linguistics.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Yiqing Xie, Sheng Zhang, Hao Cheng, Pengfei Liu, Zelalem Gero, Cliff Wong, Tristan Naumann, Hoifung Poon, and Carolyn Rose. 2024. DocLens: Multiaspect fine-grained medical text evaluation. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 649–679, Bangkok, Thailand. Association for Computational Linguistics.

- Hongyan Xu, Hongtao Liu, Zhepeng Lv, Qing Yang, and Wenjun Wang. 2023. Pre-trained personalized review summarization with effective salience estimation. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 10743–10754, Toronto, Canada. Association for Computational Linguistics.
- Shicheng Xu, Liang Pang, Mo Yu, Fandong Meng, Huawei Shen, Xueqi Cheng, and Jie Zhou. 2024. Unsupervised information refinement training of large language models for retrieval-augmented generation. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 133–145, Bangkok, Thailand. Association for Computational Linguistics.
- Xianjun Yang, Kaiqiang Song, Sangwoo Cho, Xiaoyang Wang, Xiaoman Pan, Linda Petzold, and Dong Yu. 2023. OASum: Large-scale open domain aspectbased summarization. In *Findings of the Association* for Computational Linguistics: ACL 2023, pages 4381–4401, Toronto, Canada. Association for Computational Linguistics.
- Min-Ling Zhang and Zhi-Hua Zhou. 2007. Ml-knn: A lazy learning approach to multi-label learning. *Pattern recognition*, 40(7):2038–2048.
- Nan Zhang, Yusen Zhang, Wu Guo, Prasenjit Mitra, and Rui Zhang. 2023a. FaMeSumm: Investigating and improving faithfulness of medical summarization. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 10915–10931, Singapore. Association for Computational Linguistics.
- Yusen Zhang, Yang Liu, Ziyi Yang, Yuwei Fang, Yulong Chen, Dragomir Radev, Chenguang Zhu, Michael Zeng, and Rui Zhang. 2023b. MACSum: Controllable summarization with mixed attributes. *Transactions of the Association for Computational Linguistics*, 11:787–803.
- Chujie Zheng, Kunpeng Zhang, Harry Jiannan Wang, and Ling Fan. 2020. A two-phase approach for abstractive podcast summarization. *arXiv preprint arXiv:2011.08291*.
- Ming Zhong, Da Yin, Tao Yu, Ahmad Zaidi, Mutethia Mutuma, Rahul Jha, Ahmed Hassan Awadallah, Asli Celikyilmaz, Yang Liu, Xipeng Qiu, and Dragomir Radev. 2021. QMSum: A new benchmark for querybased multi-domain meeting summarization. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 5905–5921, Online. Association for Computational Linguistics.

## A Annotation Guideline

#### A.1 Annotation Tool

924

925

926

928

929

930

931

932

935

937

938

939

941

942

945

947

949

951

We developed a custom interactive annotation tool to support efficient and user-friendly dataset annotation, which is accessible online. The backend was implemented in the Go programming language<sup>5</sup>, chosen for its performance and simplicity. The frontend was built using the Vue.js framework<sup>6</sup>, which enabled a responsive and intuitive user interface, and PostgreSQL<sup>7</sup> served as the database.

#### A.2 Consent Statement

Users first register on the tool by providing their email address and selecting their role (medical domain or NLP domain). Registration is subject to approval by an administrator. During the session, only non-personal cookies are collected, and users can choose whether to accept them, as shown in Table 7. Access to the annotation interface is granted only after the user has provided explicit consent.

I agree to the use of the collected data for research purposes. I agree to the use of functional cookies on this site.

 Table 7: Consent Statement.

#### A.3 Task Assignment

Both evaluation and annotation tasks are randomly assigned by administrators, as illustrated in Figure 6. Each data sample is assigned to two annotators from different domains—one from the medical domain and one from the NLP domain. Annotators were instructed not to communicate with each other to maintain data quality and ensure the authenticity of their responses.

		HOME		
	PMID: 30908472	PMID: 30965384 : ser	PMID: 31371315 :	PMID: 31385109: 100
4	AUSTRALIAN AND NEW ZEALAND STUDY OF PHOTODYNAMIC THERMPY IN CHOROIDAL AMELANOTIC MELANOMA	Benefits of a brief psychological intervention targeting fear of cancer recurrence in people at high risk of developing another melanoma: 1	A prospective trial of adjustrat therapy for high-risk uncell melenomic assessing 5-year servicel extremes	Phase I study of the anti-endothelin B receptor antibody-drug conjugate DEDN6585A in patients with metastatic or unresectable
2	PMID: 31403867 :==	PMID: 31425478 ; 100	PMID: 31425480	PMID: 31535302
	Autishib in Combination With Toripalimab. a Hunsmitzed Immunoglabulin G(4) Monoclonal Antibody Against Programmed Cell Deeth-1,	Necodywant intravitreal ranibizumab treatment in high-risk coular molonoma patients: a two-stage single-centre phase IL.	Validation of datordenic trametinity prognossic groups in patients treated with venuraterial and cobinestinith for advanced BRAF-mutated	The MELFO Study: A Multicenter, Prospective, Randomized Clinical Trial on the Effects of a Reduced Stage Adjusted Follow-Up Schedule
	PMID: 31553661	PMID: 31558480	PMID: 31581055	PMID: 31564722.cm
	Adjuvant Whole-Brain Radiation Therapy Compared With Observation //hirr Local Treatment of Melanoma Brain Melastases: A	Randomized Phase II Trial and Tumor Muzational Spectrum Analysis from Cabazantinib versus Chemotherapy in	Safety and efficacy of nivolumab in challenging subgroups with advanced melanoms who progressed on or after joilinsumab treatment	Site-specific response patterns, pseudoprogression, and acquired resistance in patients with melanona treated with

Figure 6: List of tasks in the annotation tool.

#### A.4 Evaluation Phase

n the evaluation phase, the evaluator is required to assess two components of the system output based

on three aspects: Completeness (Comprehensiveness), Conciseness (Faithfulness), and Traceability. Each aspect is rated using a 5-point Likert scale, with detailed scoring guidelines provided in Table 8. On the evaluation page, the left panel displays the content of the article (specifically, the abstract section), while the right panel presents summary cards corresponding to seven medical aspects. When the user hovers over a summary card, the relevant sentences in the abstract on the left are highlighted, as illustrated in Figure 7. The highlight remains visible until the user hovers over another summary card, enabling easy traceability to the corresponding source sentences in the article.

	номе	GUIDE			
	TITLE: Randomized Phase II Trial and Tumor Mutational Spectrum Analysis from Cabozantinib versus Chemotherapy in Metastatic Uveal Melanoma (Alliance AD91201	ı Î	Research Aims   [ 1 ]		(18)
al Overview Tasks ^ FT Todo Cone	MMID: 31534480 Year: 2007 Author: Lule JJ Jonané Clin Cancer Inn     MID: 31534480 Year: 2007 Author: Lule JJ Jonané Clin Cancer Inn     Middle State receptor MET is highly expressed on primary uneal     metanoma, MET inhibiton demonstrated asky diricial signals of efficacy is showing     uneal metanoma growth, [1] The primary objective of our study was to compare the     programisci-fie autrivial rate alt -morth (1954) of palatient with swall metanoma     result in alt alt -morth (1954) of palatient with swall metanoma	1	The study aims to compa with useal melanoma treat Politifulous ***	re the progression-free survi ated with cabozantinib or ch Corporatersiveness ****	nal rate at 4 months of patients imotherapy. Troublity ***
Revise	treated with caloozantinib or chemotherapy. [2] PATIENTS AND METHODS: Patients with metastatic useal melanoma and RECIST measurable disease were randomized 2: to receive either caloozantinib (zmr 1) versus temozoloniside or dacarbatine (zmr 2)	1	Research Methods   [2]	I	1
	with restaging imaging every two cycles. (a) Loss-over norm alm z to calcolation to after progression was allowed time 72.1. (a) Available tumor specifications were avalyce by whole-exome sequencing (WES) and results were correlated with outcome. [5] RESULTS: Forty-site eligible patients were accrued with 31, 15, and 9 in arms 1, 2, and 22, respectively. [6] Median inters of prior themapy, including hepdie embediation.	1	Patients with metastatic u cabozantinib or chemoth every two cycles. Fottigulous	oveal melanoma were randor erapy (ternozolornide or clac Comprehensivenesr	rized to receive either arbazine) with restaging imaging Topochility
	were two. [7] Bates of PIS4 in arm 1 and arm 2 seee 32.2% and 26.2% [9 = 0.35] respectively, with medias PIS time of 60 and 50 days (P = 0.061 HR = 0.590 [8] Median overall survival (GS) was 64 medits and 7.3 months (P = 0.506 HR = 1.21).		****	*****	*****
	respectively. [9] Grade 3-4 Common Terminology Criteria for Adverse Events were present in 61.3%, 46.7%, and 37.5% in arms 1, 2, and 2K, respectively. [10] WES demonstrated a mean tumor mutational burden of 1.53 mutations/Mb and did not		Research Results   [7, 9,	. 12 ]	۲
	separate OS s or >1 year (P = 0.14), (11) Known mutations were identified by WES and novel mutations were nominated. (12) CONCLUSIONS: METATORIR blockade with calcurateiths demonstrated on simprovement in STS but an increase in twoichy relative to tempolomide/dacabasine is metastatic useal melanoma.		The study found that cab increased toxicity compar melanoma.	ozantinib did not improve pr red to temozolomide/dacarb	ogression-free survival but azine in metastatic uveal
	Open Is PubMed		*****	*****	****

Figure 7: Evaluation page in the annotation tool.

#### A.5 Revision Phase

Out of the 3.5K evaluated data instances, 741 (21%) were filtered for further revision. The filtering criteria were as follows: (1) the mean score for any of the three evaluation metrics was below 3.5, or (2) the score difference between annotators exceeded 2.0. Annotators were then instructed to revise both the summaries and their corresponding citations based on the evaluation results. On the revision page, as illustrated in Figure 8, the left panel displayed the document content, while the right panel showed the summary along with evaluation results from two annotators. Annotators revised the summaries and updated the sentence indices according to the evaluation feedback.

TITLE: Safety, Clinical Activity, and Biological Correlates of Response in Patients with Metastatic Malanemy: Dentity from a Phase I Trial of Manolinumah	Research Results	[ 9, 10 ]		
PMID: 31358540 Year: 2019 Author: Jourani:	Atezolizumab show overall response rat	ed durable responses ar e of 30% and a median	d survival in patients w overall survival of 23 m	ith melanoma, with an ortha. 🕒
(0) PURPOSE: Attectionable (and programmed death-ligand 1 (PD-L1)) selectively targets PD-L1 to block its interaction with receptors programmed death-ligand 17.1, thereby reiningconting antitumer - Icell activity. (19) We evaluated the long-terms attery and activity	1 (Martines	Comprehensiveness	Traceability	Remark
of atezolizumab, along with biological correlates of clinical activity endpoints, in a cohort of patients with restances in an ongoing phase Ia study (NCI01375842) [2] (NJIENIS AND METHODS: Patients with unreasciable or metatatic metanome were encreded to receive	****	*****	*****	12 13 Information of biological correlates is mire
atezolizumab 0.1 to 20 mg/kg or >10 mg/kg every 3 weeks. [3] Primary study objectives	± 06830950	研究での設め		
were sarety and toestaulity. (4) secondary objectives included investigator-assessed efficacy measures; pharmatodynamic and predictive biomaters of antibumor activity were explored. (5) RESULTS: Forty-five patients were enrolled and were evaluable for safety. (6) More and the safety of the safety.	Fattifulress	Comprehensiveness ★★★★★	Traceability ★★★★★	Remark Main Conclusion: 12, 13 ]
much contributions and a service of the service of	1. If necessary, pleat information on this	ie rewrite this summary aspect, enter "Unknown	according evaluation re	sults. If there is no
95% CL 35-not estimable (VEI) [10] Clinically meaningful long-term survival was observed, with a median overall survival of 23 months (95% CL 9-466, 111) Baseline biomarkers of turnor immunity (PD-L1 expression on immune cells, T effector (Teff), and antigen	Atezolizumab sho an overall respons median overall su	wed durable responses ie rate of 30%, a mediar rvival of 23 months.	and survival in patients induration of response	with melanoma, with of 62 months, and a
presentation gene signatures) and turnor mutational burden (TMB) were associated with improved response, progression-free survival, and overall survival. [12] CONCLUSIONS: AtteoDizamab was well tolerated, with durable responses and survival in putients with	2. If necessary, plea to be separated by	ie rewrite the positions i commas. If there is no ir	according evaluation re dormation on this aspe	sults, and numbers nee ct, enter "[]".
melanoma. [13] PD-L1 expression, TMB, and Teff signatures may indicate improved benefit	[9,10,12]			

Figure 8: Revision page in the annotation tool.

955

956

957

958

959

960

961

962

963

964

965

966

967

968

978

979

980

981

982

<sup>&</sup>lt;sup>5</sup>https://go.dev/ <sup>6</sup>https://vuejs.org/

<sup>&</sup>lt;sup>7</sup>https://www.postgresql.org/

Aspect	Likert Score	Score Description
	****	All key relevant information from the article is accurately captured.
	★★★★☆	Most key relevant information from the article is present, with minor omissions.
Completeness	★★★☆☆	Some key relevant information from the article is present, but some is missing.
•	★★☆☆☆	Most key relevant information from the article is missing.
	★☆☆☆☆	All key relevant information from the article is missing.
	****	In the generated summary, all content is relevant to this aspect.
	★★★★☆	In the generated summary, most content is relevant to this aspect, with minor irrelevant parts.
Completeness	★★★☆☆	In the generated summary, some content is relevant to this aspect, while some is irrelevant.
	★★☆☆☆	In the generated summary, most content is irrelevant to this aspect.
	★☆☆☆☆	In the generated summary, all content is irrelevant to this aspect or contains errors.
	****	All relevant sentences have been accurately traced (highlighted).
	★★★★☆	Most relevant sentences have been accurately traced (highlighted).
Traceability	★★★☆☆	Some relevant sentences have been accurately traced, but some are missing or irrelevant.
	★★☆☆☆	Most relevant sentences have not been accurately traced.
	★☆☆☆☆	None of the relevant sentences have been accurately traced.

Table 8: Evaluation Criteria and Scoring Guidelines.



Figure 9: Characteristics of the TRACSUMdataset.

984

985

987

991

992

993

997

998

999

1000

1001

1002

1003 1004

1005

1006

1008

# **B** Characteristics of the Dataset

## **B.1** Source Article Length

Among the 500 abstracts, the average length per abstract was 319.89 tokens, with the longest containing 1,104 tokens and the shortest containing only 25. The distribution of token counts across abstracts is illustrated in Figure 9a. Additionally, each abstract contained an average of 10.42 sentences, with sentence counts ranging from 1 to 32. The distribution of sentence counts is shown in Figure 9b.

## **B.2** Aspect Coverage in Abstracts

All 500 documents contained information on at least three aspects. Among them, 118 documents covered all seven aspects, and 211 documents covered six aspects, as illustrated in Figure 9c.

# **B.3** Proportion of Positive and Negative Data

We analyzed the distribution of positive and negative data samples across seven aspects, as shown in Figure 9d. All 500 abstracts included aspect A (Research Aims), while 499 covered aspect I (Research Methods or Intervention) and aspect O (Research Results or Outcomes). In contrast, aspect D (Treatment Duration) was less common, appearing in only 174 abstracts. Overall, the ratio

of positive to negative samples was 2862:638.

# **B.4** Length of Traceable Summaries

As shown in Table 9, all 2,862 positive summaries had an average length of 28.06 tokens, with the longest containing 77 tokens and the shortest just 3. On average, each summary cited 1.78 sentences, with the number ranging from 1 to 7. Among all aspects, summaries related to aspect S (Side Effects) had the highest average token count, while those concerning aspect I (Research Methods or Intervention) cited the most sentences. 1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1021

1022

1024

1025

# **C** Generation Pipelines

In this section, we provide a detailed description of the design and training of our three baseline methods: TRACK-THEN-SUM, SUM-THEN-TRACK, and END-TO-END.

# C.1 TRACK-THEN-SUM

As illustrated in Figure 10, the TRACK-THEN-SUM 1026 generation pipeline consists of two phases: tracking 1027 and summarization. In the first phase, the tracker 1028 module  $\mathcal{T}$  retrieves the sentences most relevant to 1029 the given aspect using a default threshold of 0.5. In 1030 the second phase, the summarizer module S gen-1031 erates a concise summary based on the selected 1032 sentences. Finally, the summary and the cited sen-1033

	Summary					Citations								
	Α	Ι	0	Р	Μ	D	S	A	Ι	0	Р	Μ	D	S
Min	13	15	12	4	3	4	4	1	1	1	1	1	1	1
Max	56	73	77	69	77	75	75	5	7	6	5	6	5	4
Avg.	29.33	37.81	34.75	25.64	25.37	17.82	25.67	1.51	2.33	2.58	1.61	1.74	1.25	1.46

Table 9: Length of summaries (in tokens) and number of citations (in sentences) in positive samples.



Figure 10: TRACK-THEN-SUM summarization pipeline.

tences are merged to form the final system output.

## C.2 Tracker $\mathcal{T}$

We implement the sentence tracing task as a binary classification of sentences within the abstract.

**Data Collection:** We applied sentence tokenization to each abstract in the training set. For every sentence, we created (c, a) pairs by combining it with each predefined aspect  $a \in A$ . Each pair was labeled with a binary variable y based on the corresponding *citations* field: if the sentence index appeared in the *citations* associated with aspect a, we assigned y = 1; otherwise, y = 0. In total, we obtained 35.5K sentence-aspect-label pairs, forming the training dataset  $D_T$ .

**Training:** Given the constructed dataset  $\mathcal{D}_{\mathcal{T}}$ , we initialized tracker  $\mathcal{T}$  using a pre-trained language model (LM) as the backbone. The model was subsequently fine-tuned on  $\mathcal{D}_{\mathcal{T}}$  using a standard binary classification objective which maximizes the log-likelihood of the observed labels:

$$\max_{\mathcal{T}} \mathbb{E}_{((c,a),y) \sim \mathcal{D}_{\mathcal{T}}} \log p_{\mathcal{T}}(y \mid (c,a))$$

We fine-tuned the tracker  $\mathcal{T}$  using the QLoRA technique, initializing from the 4-bit quantized version of the LLaMA-3.1-8B-Instruct backbone<sup>8</sup>, on  $\mathcal{D}_{\mathcal{T}}$ . To enable binary classification, we appended a lightweight classification head that maps the model's output to a single scalar representing the predicted probability. Training was conducted on six NVIDIA A6000 GPUs with a batch size of 32, gradient accumulation steps of 2, and a total of 5 epochs. We employed a learning rate of  $1 \times 10^{-5}$ , applied a weight decay of 0.01, set the random seed to 3407 for reproducibility, and used 200 warmup

steps. The full training process took 17 hours and 2 minutes.

1068

1069

1070

1071

1072

1073

1074

1076

1077

1078

1079

1080

1081

1082

1083

1084

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1099

# C.3 Summarizer S

**Data Collection:** For each summary *sum* in the training set, we extracted related sentences from the abstract based on the *citations* field to form the set C. Each C was paired with its associated aspect a, and combined with the *sum* to form ((C, a), sum). In total, we obtained 2.8K citations-aspect-summary pairs, forming the training dataset  $D_S$ .

**Training:** Similar to the training of  $\mathcal{T}$ , we initialized summarizer S using a pre-trained LM as the backbone. We then fine-tuned summarizer S on  $\mathcal{D}_S$ using a standard next-token prediction objective, which maximizes the likelihood of generating the target summary *sum* given the input (C, a) pair:

$$\max_{\mathcal{S}} \mathbb{E}_{((C,a),sum) \sim \mathcal{D}_{\mathcal{S}}} \log p_{\mathcal{S}}(sum \mid C, a)$$
 10

The input instruction is shown in Table 16. We finetuned Summarizer S using the Unsloth framework, starting from the 4-bit version of the LLaMA-3.1-8B-Instruct base model<sup>9</sup>, on  $\mathcal{D}_S$ . Training was performed on two NVIDIA A6000 GPUs with a batch size of 16, a gradient accumulation step size of 2, and a total of 5 epochs. We used a learning rate of 1e-5, a weight decay of 0.01, a fixed random seed of 3407, and 200 warmup steps. The entire training process took 1 hour and 55 minutes. Additionally, we adopted the train\_on\_responses\_only strategy to focus learning on relevant output segments.

## **C.4 TTS** $\oplus$ *f*.

As mentioned in §5, our TRACK-THEN-SUM 1100 method includes two variants, differing only in 1101

1034

1035

1055

1056

1057

1058

1059

1060

1061

1062

1063

1065

1066

<sup>&</sup>lt;sup>8</sup>Model: meta-llama/Llama-3.1-8B

<sup>&</sup>lt;sup>9</sup>Model: unsloth/Meta-Llama-3.1-8B-Instruct-bnb-4bit



Figure 11: SUM-THEN-TRACK method pipeline.

Algorithm 3: SUM-THEN-TRACK Inference
<b>Require:</b> Tracker $\mathcal{T}$ , Summarizer $\mathcal{S}$
<b>Input:</b> article $d = \{c_1, c_2,, c_n\}$ and aspect $a \in \mathcal{A}$
<b>Output:</b> summary <i>sum</i> and its citations $C'$
1: $sum \leftarrow \mathcal{S}(a, d);$
$2: \mathcal{C}' \leftarrow \emptyset;$
3: foreach $c_i \in \{c_1, c_2,, c_n\}$
4: $\mathcal{T}$ predict <b>relevance</b> given $(sum, c_i)$ ;
5: <b>if relevance</b> == Yes <b>then</b> append $c$ to $C'$ ;

Algorithm 3: SUM-THEN-TRACK inference process.

their input. Specifically, the TTS  $\oplus$  f. variant uses both the set of cited sentences and the full context (i.e., abstract) as input. The input instruction is shown in Table 17. All other settings remain unchanged, except for the batch size, which was set to 8. Under this configuration, training took 8 hours and 36 minutes.

#### C.5 SUM-THEN-TRACK

1102

1103

1104

1105

1106

1107

1108

1109

1110

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

113

#### C.5.1 Inference Overview

As illustrated in Figure 11, the SUM-THEN-TRACK 1111 method consists of two phases: summarization and 1112 tracking. In the first phase, the summarizer S gen-1113 erates an aspect-specific summary sum from an 1114 abstract d based on a given aspect a. In the second 1115 phase, the tracker  $\mathcal{T}$  identifies the sentences most 1116 relevant to this summary using a default similar-1117 ity threshold of 0.5. Finally, the summary and the 1118 corresponding sentences are combined to form the 1119 1120 final output, as shown in Algorithm 3.

#### C.5.2 Summarizer S

**Data Collection:** We extracted abstract, aspect, and summary fields from the training set, resulting in 2.8K ((*d*, *a*), sum) pairs, denoted as  $\mathcal{D}_S$ .

**Training:** We then initialized summarizer S using a pre-trained LM as the backbone. We then finetuned summarizer S on  $\mathcal{D}_S$  using a standard nexttoken prediction objective, which maximizes the likelihood of generating the target summary *sum* given the input (d, a) pair:

$$\max_{\mathcal{S}} \mathbb{E}_{((d,a),sum) \sim \mathcal{D}_{\mathcal{S}}} \log p_{\mathcal{S}}(sum \mid d, a)$$

1132 The input instruction is shown in Table 18. We fine-1133 tuned Summarizer S using the Unsloth framework, starting from the 4-bit version of the LLaMA-3.1-8B-Instruct base model, on  $\mathcal{D}_S$ . Training was performed on two NVIDIA A6000 GPUs with a batch size of 8, a gradient accumulation step size of 2, and a total of 5 epochs. We used a learning rate of 1e-5, a weight decay of 0.01, a fixed random seed of 3407, and 200 warmup steps. The entire training process took 7 hour and 32 minutes. Additionally, we adopted the train\_on\_responses\_only strategy to focus learning on relevant output segments. 1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

#### C.5.3 Tracker $\mathcal{T}$

**Data Collection:** We first applied sentence tokenization to all abstracts in the training set. For each abstract, every sentence c was paired with each summary sum, forming (c, sum) pairs. Each pair was then labeled with y based on the *citations* field. This process resulted in 35.5k ((c, sum), y) pairs, denoted as  $\mathcal{D}_{\mathcal{T}}$ .

**Training:** Given the constructed dataset  $\mathcal{D}_{\mathcal{T}}$ , we initialized tracker  $\mathcal{T}$  using a pre-trained language model (LM) as the backbone. The model was subsequently fine-tuned on  $\mathcal{D}_{\mathcal{T}}$  using a standard binary classification objective which maximizes the log-likelihood of the observed labels:

$$\max_{\mathcal{T}} \mathbb{E}_{((c,sum),y) \sim \mathcal{D}_{\mathcal{T}}} \log p_{\mathcal{T}}(y \mid (c,sum))$$

We fine-tuned the tracker  $\mathcal{T}$  using the QLoRA technique, initializing from the 4-bit quantized version of the LLaMA-3.1-8B-Instruct backbone, on  $\mathcal{D}_{\mathcal{T}}$ . To enable binary classification, we appended a lightweight classification head that maps the model's output to a single scalar representing the predicted probability. Training was conducted on six NVIDIA A6000 GPUs with a batch size of 32, gradient accumulation steps of 2, and a total of 5 epochs. We employed a learning rate of  $1 \times 10^{-5}$ , applied a weight decay of 0.01, set the random seed to 3407 for reproducibility, and used 200 warmup steps. The full training process took 22 hours and 12 minutes.

## C.6 END-TO-END

The END-TO-END approach employs a single 1174 model  $\mathcal{M}$ , to jointly perform summarization and 1175

Model	API Src.	Input Prices	<b>Output Prices</b>	Input Length	Output Length	Costs
Llama-3.1-8B-Inst.	DeepInfra	\$0.03	\$0.05	131K	8K	\$0.030
Llama-3.3-70B-Inst.	DeepInfra	\$0.23	\$0.40	131K	8K	\$0.250
Mistral-7B-Inst (V0.3).	DeepInfra	\$0.029	\$0.055	32K	8K	\$0.040
Mistral-8x7B-Inst.	DeepInfra	\$0.24	\$0.24	131K	4K	\$0.600
Gemma-3-12B-Inst.	DeepInfra	\$0.05	\$0.100	128K	8K	\$0.070
Gemma-3-27B-Inst.	DeepInfra	\$0.10	\$0.20	128K	8K	\$0.110
GPT-40	OpenAI	\$2.50	\$10.0	128K	16K	\$2.838
GPT-4o-mini	OpenAI	\$0.15	\$0.60	128K	16K	\$0.147
					SUM	· \$4 085

Table	10.	Details	on th	e 116e	of	different	model	ΔPI
Table	10.	Details	on ui	e use	01	umerent	moder	ALIS



Figure 12: END-TO-END generation pipeline.

Algorithm 4: END-TO-END Inference
<b>Require:</b> Model $\mathcal{M}$
<b>Input:</b> article $d = \{c_1, c_2,, c_n\}$ and aspect $a \in \mathcal{A}$
<b>Output:</b> summary $sum$ and its citations $C'$
1: $\mathcal{C}' \leftarrow \emptyset$ ;
2: $(sum, \mathcal{C}') \leftarrow \mathcal{M}(a, d);$

Algorithm 4: END-TO-END inference process.

sentence tracking, as shown in Figure 12.

#### C.6.1 Inference Phase

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

Given a abstract d and an aspect  $a \in A$ ,  $\mathcal{M}$  generates a summary focused on a and  $\mathcal{C}'$  on which the summary relies, as illustrated in Algorithm 4.

## C.6.2 Training Phase

Data Collection. We extracted abstract, aspect, summary, and citations fields from the training set and then combined them into ((d, a), (sum, C)) pairs. As a result, we obtained 2.8K training instances, denoted by D<sub>M</sub>.

**Training.** We then initialized  $\mathcal{M}$  with a pre-trained LM and trained it on  $\mathcal{D}_{\mathcal{M}}$  using a standard conditional language modeling objective, maximizing the likelihood:

$$\max_{\mathcal{M}} \mathbb{E}_{((d,a),(sum,\mathcal{C}))\sim \mathcal{D}_{\mathcal{M}}} \log p_{\mathcal{M}}(sum,\mathcal{C} \mid d,a)$$

The input instruction is shown in Table 19. We 1192 fine-tuned  $\mathcal{M}$  using the Unsloth framework, start-1193 ing from the 4-bit version of the LLaMA-3.1-8B-1194 Instruct base model, on  $\mathcal{D}_{\mathcal{M}}$ . Training was per-1195 1196 formed on two NVIDIA A6000 GPUs with a batch size of 8, a gradient accumulation step size of 2, 1197 and a total of 5 epochs. We used a learning rate of 1198 1e-5, a weight decay of 0.01, a fixed random seed 1199 of 3407, and 200 warmup steps. The entire training 1200

process took 8 hours and 16 minutes. Additionally,1201we adopted the train\_on\_responses\_only strat-<br/>egy to focus learning on relevant output segments.1202

1204

1205

1206

1207

1208

1209

1218

1219

#### **D** API Cost

#### **D.1** Dataset Collection Costs

We initially generated our dataset with the free credits provided by the Mistral-Large API, so the cost for this part is \$0.

#### **D.2** Evaluation Costs

We incurred approximately \$4.085 in API costs 1210 to obtain results from eight different models on 1211 the test set, as detailed in Table 10. The test set 1212 comprises 700 data samples, each formatted into 1213 prompts, resulting in approximately 100K input 1214 tokens in total. The number of output tokens varies 1215 across LLMs; standard text generation models typ-1216 ically produce around 50K output tokens. 1217

#### **E** Experiment Analysis

#### **E.1** Full Context $\oplus C$ vs. C only

In this section, we present an example to illustrate 1220 how incorporating full context impacts summary 1221 generation and, in turn, affects claim recall. When 1222 the cited sentences (i.e., the tracker  $\mathcal{T}$  output) re-1223 main fixed, providing the full document as addi-1224 tional input enables the summarizer S to better 1225 resolve abbreviations and domain-specific termi-1226 nology, thereby enhancing claim recall. As shown 1227 in Table 11, TTS  $\oplus f$  resolves the abbreviation "RT" as "radiation therapy", which leads the NLI model 1229 (TRUE) to determine that the subclaim is entailed 1230 by the reference text during entailment evaluation. 1231 This results in an increase in the overall claim recall 1232 score from 2/4 to 3/4. 1233

However, providing additional context beyond1234the cited sentences may cause the summarizer S to1235incorporate irrelevant or unsupported information1236(i.e., content not present in the cited sentences),1237

**Reference:** Summary  $\neg$  Citations  $\rightarrow 0, 1, 7$ 

- 1'. A total of 50 participants were involved in the study.
- 2'. Participants with cutaneous neurotropic melanoma of

the head and neck.

1238

1239

1240

1241

1242

1243

1244

1245

1246

1947

1249

1250

1251

1252

1253

1254

1255

1256

1257

1258

1259

1261 1262

1264

1266

- 3'. 23 participants were assigned to the observation group. 4'. 27 participants were assigned to the radiation therapy group.
- Citation 0: BACKGROUND: Cutaneous neurotropic melanoma (NM) of the head and neck (H&N) is prone to local relapse, possibly due to difficulties widely excising the tumor. Citation 1: This trial assessed radiation therapy (RT) to the primary site after local excision. Citation 7: During 2009-2020, 50 participants were randomized: 23 to observation and 27 to RT.

27 10 KI.
<b>TTS Output:</b> Subclaims $\rightarrow$ 7 Citations $\rightarrow$ 7
1'. A total of 50 participants were randomized in the study.
2'. 23 participants were assigned to the observation group.
3'. 27 participants were assigned to the RT group.
(TRUE) Claim Recall: 2/4. 1': ✓, 2': ✓, 3': ✗
<b>TTS</b> $\oplus$ <i>f</i> . <b>Output:</b> Subclaims $\neg$ Citations $\rightarrow$ 7
1'. A total of 50 participants were randomized in the study.
2'. 23 participants were assigned to the observation group.
3'. 27 participants were assigned to the radiation therapy
(RT) group.

(TRUE) Claim Recall: 3/4. 1′: ✓, 2′: ✓, 3′: ✓

Table 11: An example of summaries generated by TTS and TTS  $\oplus f$ , along with their claim recall comparison (PMID: 38851639, Aspect: Patients).

which could reduce claim precision or citationbased metrics. Nonetheless, our evaluation results do not show a noticeable drop in other metrics. This may be attributed to the instruction explicitly directing the summarizer S to generate summaries strictly based on the cited sentences, with the additional context serving only as reference.

#### E.2 Agreement with Human Evaluation

To evaluate the relationship between the system outputs and task-level evaluation scores, we employ both Spearman's correlation coefficient ( $\rho$ ) (Kendall and Gibbons, 1990) and Pearson's correlation coefficient (r) (Sheskin, 2003). Pearson's r measures the strength of a *linear relationship* between two continuous variables, which is appropriate when assuming interval-scaled outputs and normally distributed scores (Benesty et al., 2009). In contrast, Spearman's  $\rho$  captures monotonic relationships based on rank order, making it more robust to non-linear patterns and outliers (Hauke and Kossowski, 2011). Using both metrics provides a comprehensive view of how well the automatic system outputs align with human-centric evaluation criteria, accounting for both linear trends and ordinal consistency.

Specifically, we randomly sampled ten abstracts from the test set, and asked the annotator to follow the procedure in Algorithm 2 to assess outputs from the best-performing method (TTS  $\oplus$  *f*.) using four

	Comp	leteness	Conc	iseness	F1 Score		
Evaluator	CLR	CIR	CLP	CIP	$F_1^{\text{cl.}}$	$F_1^{\text{ci.}}$	
Human	81.1↑	74.3↑	68.6↑	78.1↓	74.3↑	76.2↓	
TRUE	78.2	73.4	65.7	79.5	71.4	76.3	

Table 12: Comparison of evaluation results between human annotator and the TRUE model on 10 sampled abstracts.



Figure 13: Spearman's correlation coefficient ( $\rho$ ) and Pearson's correlation coefficient (r) between TRUE and human evaluation scores across four evaluation metrics.

evaluation metrics. As indicated in Table 12, human evaluations score higher than the TRUE model on most metrics, achieving an F1 score of 74.3 for claims and 76.2 for citations quality. For each of the four evaluation metrics, we computed the Spearman correlation coefficient ( $\rho$ ) and Pearson correlation coefficient (r) between the automatic evaluation results and human judgments. As shown in Figure 13, the Spearman correlation coefficient between human and automatic evaluation results is  $\rho = 0.612$ , and the Pearson correlation coefficient is r = 0.577. The agreement is relatively lower for claim-related metrics, whereas citation-related metrics demonstrate stronger consistency with human judgments.

#### E.3 Comparison of Entailment Evaluators

We experiment with two additional instructionfollowing LLMs as entailment evaluators: the proprietary GPT-40 (Hurst et al., 2024) and the opensource Mistral-Large (Mistral AI, 2025). Building on the experimental setup described in §6.3, we replace the TRUE model with each of these evaluators to assess the outputs generated by the TTS  $\oplus$  *f*. variant. The evaluation results are presented in Table 13. Among the models, GPT-40 produces scores that most closely align with human judgments, followed by Mistral.

	Compl	eteness	Conci	nciseness F1 Sco		
Evaluator	CLR	CIR	CLP	CIP	$F_1^{\text{cl.}}$	$F_1^{\text{ci.}}$
Human	81.1	74.3	68.6	78.1	74.3	76.2
TRUE	78.2	73.4	65.7	79.5	71.4	76.3
GPT-40	80.2	77.1	67.0	76.2	73.0	76.7
Mistral	75.6	76.8	70.1	74.5	72.8	75.6

Table 13: Comparison of evaluation results between human annotator and three entailment evaluators on 10 sampled abstracts. 1267

1268

1269

1270

1271

1272

1273

1274

1275

1276

1277

1278

1279

1280

1282

1283

1285

1286

1287

1288

1289

1291

# F Data Samples of TRACSUM Dataset

PMID	abstract	aspect	summary	citations	
31638282	The multinational phase 3 CheckMate 238 trial compared adjuvant therapy with nivolumab versus ipilimumab among patients with resected stage III or IV melanoma (N = 906)	d	Unknown.	[]	
33294860	In this study, we incorporate anal- yses of genome-wide sequence and structural alterations with pre- and on- therapy transcriptomic and T cell reper- toire features in immunotherapy-naive melanoma patients treated with	a	The study aims to predict response to immune checkpoint blockade by integrating genomic, transcriptomic, and immune repertoire data.	[4]	
34650833	Combination immunotherapy with sequential administration may en- hance metastatic melanoma (MM) pa- tients with long-term disease control. High Dose Aldesleukin/Recombinant Interleukin-2 (HD rIL-2) and ipili- mumab (IPI) offer	m	The study used High Dose Aldesleukin/Recombinant Interleukin-2 (HD rIL-2) at 600,000 IU/kg and ipilimumab (IPI) at 3 mg/kg.	[1,3]	
37479483	BACKGROUND: Continuous combi- nation of MAPK pathway inhibition (MAPKi) and anti-programmed death- (ligand) 1 (PD-(L)1) showed high re- sponse rates, but only limited im- provement in progression-free survival (PFS) at the cost of a high frequency	р	The study involved 33 patients with treatment-naïve BRAFV600E/K- mutant advanced melanoma, with 32 randomized into four cohorts.	[3,8]	
33593880	PURPOSE: Triple-negative breast can- cer (TNBC) is an aggressive disease with limited therapeutic options. An- tibodies targeting programmed cell death protein 1 (PD-1)/PD-1 ligand 1 (PD-L1) have entered the therapeutic landscape in TNBC, but only a minor- ity of patients benefit. A way to reli- ably enhance immunogenicity, T-cell infiltration, and predict responsiveness is critically needed. PATIENTS AND METHODS: Using mouse models of TNBC	i	This study used mouse models of TNBC to evaluate immune activa- tion and tumor targeting of intra- tumoral IL12 plasmid followed by electroporation (Tavo), conducted a single-arm prospective clinical trial of Tavo monotherapy in patients with treatment-refractory advanced TNBC, and expanded findings using publicly available breast cancer and melanoma datasets.	[3,4,5]	
38870745	BACKGROUND: Treatment op- tions for immunotherapy-refractory melanoma are an unmet need. The MASTERKEY-115 phase II, open- label, multicenter trial evaluated talimogene	s	Treatment-related adverse events (TRAEs), including grade 3 TRAEs, serious AEs, and fatal AEs, oc- curred in 76.1%, 12.7%, 33.8%, and 14.1% of patients, respectively.	[11]	
33127652	PURPOSE: Increased -adrenergic re- ceptor (-AR) signaling has been shown to promote the creation of an immuno- suppressive tumor microenvironment (TME)	0	The combination of propranolol with pembrolizumab in treatment- naïve metastatic melanoma is safe and shows very promising activity with an objective response rate of 78%.	[ 12,14 ]	

Table 14: Seven traceable aspect-based summary samples from TRACSUM dataset.

## **G** Instructions

#### G.1 LLM Prompt Template

#### Instructions

Given a document consisting of a set of sentences with a marker attached to the head of each sentence. Based on the demonstrations, please summarize the research questions or aims of this study in one sentence and output the sentence markers involved. If there is no relevant information in the document, answer "Unknown".

#### Document

'[ "0: The EORTC-STBSG coordinated two large trials of adjuvant chemotherapy (CT) in localized high-grade soft tissue sarcoma (STS).", "1: Both studies failed to demonstrate any benefit on overall survival (OS).", "2: The aim of the analysis of these two trials was to identify subgroups of patients who may benefit from adjuvant CT." "3: Individual patient data from two EORTC trials comparing doxorubicin-based CT to observation only in completely resected STS (large resection, R0/marginal resection, R1) were pooled.", ... ]'

Summary: . Citations: .

# Demonstrations

#### Document

'[ "0: Giant cell tumor of bone (GCTB) is an aggressive primary osteolytic tumor.", "1: GCTB often involves the epiphysis, usually causing substantial pain and functional disability.", "2: Denosumab, a fully human monoclonal antibody against receptor activator of nuclear factor ligand (RANKL), is an effective treatment option for patients with advanced GCTB.", "3: This analysis of data from an ongoing, open-label study describes denosumab's effects on pain and analgesic use in patients with GCTB. " "4: Patients with unresectable disease (e.g. sacral or spinal GCTB, or multiple lesions including pulmonary metastases) were enrolled into Cohort 1 (N = 170), and patients with resectable disease whose planned surgery was associated with severe morbidity (e.g. joint resection, limb amputation, or hemipelvectomy) were enrolled into Cohort 2 (N = 101).", ... ]' **Summary**: The study aims to evaluate the effects of denosumab on pain and analgesic use in patients with giant cell tumor of bone (GCTB).

Citations: [3]

#### Document

'["0: Common adverse events associated with nivolumab included fatigue, pruritus, and nausea.", "1: Drug-related adverse events of grade 3 or 4 occurred in 11.7% of the patients treated with nivolumab and 17.6% of those treated with dacarbazine." "2: Nivolumab was associated with significant improvements in overall survival and progression-free survival, as compared with dacarbazine, among previously untreated patients who had metastatic melanoma without a BRAF mutation.", "3: (Funded by Bristol-Myers Squibb; CheckMate 066 ClinicalTrials.gov number, NCT01721772.)." ]' **Summary**: Unknown.

Citations: Null.

Table 15: Instructions and demonstrations for generating summaries on aspect A (research aims). The text denotes placeholders to be replaced with aspect-specific descriptions.

#### G.2 Instruction for summarizer S in TRACK-THEN-SUM

#### Instructions

Summarize the research aims or questions of the study in one clear sentence that includes all key details from the input sentences without omitting important information.

#### Sentences

'[ "The EORTC-STBSG coordinated two large trials of adjuvant chemotherapy (CT) in localized high-grade soft tissue sarcoma (STS).", "Both studies failed to demonstrate any benefit on overall survival (OS).", "The aim of the analysis of these two trials was to identify subgroups of patients who may benefit from adjuvant CT." "Individual patient data from two EORTC trials comparing doxorubicin-based CT to observation only in completely resected STS (large resection, R0/marginal resection, R1) were pooled." ]'

#### Summary:

Table 16: Instruction used to generate summaries for aspect A (research aims) in the summarization component of TRACK-THEN-SUM. The text denotes placeholders to be replaced with aspect-specific descriptions.

1300 1301

1297 1298

## G.3 Instruction for summarizer $S (\oplus full \ context)$ in TRACK-THEN-SUM

#### Instructions

Summarize the research aims or questions of the study in one clear sentence that includes all key details from the input sentences without omitting important information. The summary must be based solely on the provided sentences. The full text is for reference only and must not be used to introduce any new information not present in the sentences.

#### Sentences

'[ "The aim of the analysis of these two trials was to identify subgroups of patients who may benefit from adjuvant CT." "Individual patient data from two EORTC trials comparing doxorubicin-based CT to observation only in completely resected STS (large resection, R0/marginal resection, R1) were pooled." ]'

#### Full Context

'[ "The EORTC-STBSG coordinated two large trials of adjuvant chemotherapy (CT) in localized high-grade soft tissue sarcoma (STS).", "Both studies failed to demonstrate any benefit on overall survival (OS).", "The aim of the analysis of these two trials was to identify subgroups of patients who may benefit from adjuvant CT." "Individual patient data from two EORTC trials comparing doxorubicin-based CT to observation only in completely resected STS (large resection, R0/marginal resection, R1) were pooled.", ... ]'

#### Summary:

Table 17: Instruction used to generate summaries for aspect A (research aims) in the summarization component of TRACK-THEN-SUM ( $\oplus f$ .). The text denotes placeholders to be replaced with aspect-specific descriptions.

#### G.4 Instruction for summarizer S in SUM-THEN-TRACK

#### Instructions

Summarize the research aims or questions of the study in one clear sentence based on the given article.

#### Article

'[ "The EORTC-STBSG coordinated two large trials of adjuvant chemotherapy (CT) in localized high-grade soft tissue sarcoma (STS).", "Both studies failed to demonstrate any benefit on overall survival (OS).", "The aim of the analysis of these two trials was to identify subgroups of patients who may benefit from adjuvant CT." "Individual patient data from two EORTC trials comparing doxorubicin-based CT to observation only in completely resected STS (large resection, R0/marginal resection, R1) were pooled.", ... ]'

#### Summary:

Table 18: Instruction used to generate summaries for aspect A (research aims) in the summarization component of SUM-THEN-TRACK. The text denotes placeholders to be replaced with aspect-specific descriptions.

#### G.5 Instruction for model $\mathcal{M}$ in END-TO-END

#### Instructions

Given an article, summarize the research aims or questions of the study in one clear sentence and output the index of the cited sentences.

#### Sentences

'["0: The EORTC-STBSG coordinated two large trials of adjuvant chemotherapy (CT) in localized high-grade soft tissue sarcoma (STS).", "1: Both studies failed to demonstrate any benefit on overall survival (OS).", "2: The aim of the analysis of these two trials was to identify subgroups of patients who may benefit from adjuvant CT." "3: Individual patient data from two EORTC trials comparing doxorubicin-based CT to observation only in completely resected STS (large resection, R0/marginal resection, R1) were pooled.", ... ]'

#### Summary: Citations:

Table 19: Instruction used to generate summaries for aspect A (research aims) in the END-TO-END. The text denotes placeholders to be replaced with aspect-specific descriptions.

1303 1304

1308