It's Not About the Terms: Structured Topic Descriptions for Scientific Corpora

Anonymous ACL submission

Abstract

Topic models uncover thematic structures in large document collections by assigning documents to topics and representing each topic as a ranked list of terms. However, these lists are often hard to interpret and insufficient for knowledge-intensive exploration, especially in scientific domains. We propose the task of Topic Description for Scientific Corpora, which focuses on generating structured, concise, and informative summaries for topicspecific document sets. To this end, we adapt two LLM-based pipelines: Selective Context Summarisation (SCS), which uses maximum marginal relevance to select representative documents; and Compressed Context Summarisation (CCS), a hierarchical approach based on the RAPTOR framework that recursively abstracts subsets of documents to compress the input. We evaluate both methods using SUPERT and a multi-model LLM-as-a-Judge across three topic modeling backbones (CTM, BERTopic, TopicGPT) and three scientific corpora. SCS consistently outperforms CCS in quality and robustness, while CCS performs better on larger topics despite a higher risk of information loss. Our findings highlight tradeoffs between selective and compressed strategies and provide new benchmarks for topiclevel summarisation. Code and data for two of the three datasets will be released.

1 Introduction

005

007

011

017 018

019

028

Gaining an overview of large scientific corpora
is useful for exploring research areas, identifying
common methodologies, and tracking emerging
developments. A common entry point is topic
modeling, which reveals underlying topics and
presents them as ordered lists of terms. Algorithms such as Latent Dirichlet Allocation (LDA;
Blei (2012)), Contextualised Topic Models (CTM;
Bianchi et al. (2021)) and BERTopic (Grootendorst,
2022) are widely used for this purpose. While ef-

fective for organising unlabelled data, these methods only provide term-based topic representations, making them difficult to interpret (Chang et al., 2009). Most topic modeling pipelines stop at this level, which limits their usefulness for knowledgeintensive tasks, particularly in scientific domains where understanding a research topic requires insight into research goals, methods, and purposes. 043

045

047

049

051

054

055

057

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

Recent work has sought to improve interpretability by enriching topic representations with machinegenerated labels or short contextual snippets (Lau et al., 2011; Popa and Rebedea, 2021; Rosati, 2022; Azarbonyad et al., 2023). However, these approaches often rely on surface-level signals, lack domain-specific grounding and fail to incorporate document-level context. Consequently, they offer limited support for understanding the underlying content of complex domains such as science.

We address this problem by introducing the task of **Topic Description for Scientific Corpora**, which aims to generate structured and informative summaries for topics derived from topic models. These descriptions enrich the topic representation by incorporating document-level context while remaining aligned with the topic terms, offering a clearer view of the underlying research themes.

To this end, we adapt two pipelines based on large language models (LLMs). The first, *Selective Context Summarisation* (SCS), uses Maximum Marginal Relevance (MMR; Carbonell and Goldstein (1998)) to select a representative subset of topic documents prior summarisation. The second, *Compressed Context Summarisation* (CCS), adapts the RAPTOR framework (Sarthi et al., 2024), applying recursive summarisation over a hierarchy constructed from the topic's documents.

We evaluate these pipelines across three topic modeling backbones—CTM (Bianchi et al., 2021), BERTopic (Grootendorst, 2022), and Top-icGPT (Pham et al., 2024)—on three scientific corpora. Evaluation is conducted using SUPERT (Gao

et al., 2020), a reference-free semantic similarity metric, and a multi-model LLM-as-a-Judge framework using open-source models. Results show that the MMR-based pipeline consistently produces more focused and concise topic descriptions than the RAPTOR-based method. We also analyze how topic-level properties, such as size and cohesion, affect effectiveness, and complement our findings with qualitative examples.

Our contributions are:

086

090

100

101

102

105

106

107

108

109

110

127

- We introduce **Topic Description for Scientific Corpora** as the task of enriching topic model outputs with structured, interpretable, document-grounded summaries.
- We adapt and compare two LLM-based pipelines for topic-level summarisation in scientific corpora.
- We propose a robust evaluation strategy combining SUPERT and a multi-model LLM-asa-Judge framework.

2 Related Work

We review prior work on topic modeling, enhanced topic representations, and multi-document scientific summarisation. Our work builds on these areas by combining topic model outputs with LLM-based summarisation to enrich topic representations.

2.1 Topic Modeling

Topic modeling is widely used for uncovering the-111 matic structure in large text collections. Latent 112 Dirichlet Allocation (LDA; (Blei, 2012)) remains 113 a foundational model, assuming documents are 114 mixtures of latent topics and topics are distribu-115 tions over words. Contextualized Topic Models 116 (CTM; (Bianchi et al., 2021)) extend this frame-117 work by incorporating document embeddings from 118 pre-trained language models such as BERT (De-119 vlin et al., 2019) and Sentence-BERT (Reimers 120 and Gurevych, 2019). BERTopic (Grootendorst, 121 2022) clusters BERT embeddings for document 122 topic assignment, while TopicGPT (Pham et al., 123 2024) employs decoder-only LLMs to directly gen-124 erate topics. These models are applied across vari-125 126 ous domains, including scientific literature.

2.2 Enriching Topic Representations

Beyond term lists, several methods aim to create more interpretable topic representations. Early work retrieved candidate labels from external sources such as Wikipedia and ranked them by relevance to topic terms (Lau et al., 2011; Bhatia et al., 2016). Later approaches used generative models to create more descriptive labels from topic terms (Alokaili et al., 2020). BART-TL (Popa and Rebedea, 2021) fine-tunes a BART model using weakly supervised training signals derived from heuristic labels. In the scientific domain, topic interpretation often involves producing richer textual outputs. One method clusters citation statements and summarizes them using Longformer to reflect citation intent (Rosati, 2022). Topic Pages (Azarbonyad et al., 2023) construct structured descriptions by combining definition extraction using SciBERT with contextual snippets and co-occurrence-based linking. LimTopic (Azhar et al., 2025) applies BERTopic and LLMs to generate titles and summaries for topics found in scientific limitation sections. Our work uses LLMs to generate documentgrounded topic descriptions reflecting methods, purposes, and research objects.

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

2.3 Multi-Document Scientific Summarisation

Multi-document scientific summarisation (MDSS) synthesizes coherent summaries from clusters of scientific papers. Transformer-based methods such as KGSum (Wang et al., 2022) encode documents into knowledge graphs and use two-stage decoding for improved coherence. PRIMERA (Xiao et al., 2022) applies entity-level masking during pretraining to improve salience modeling, and its effectiveness extends to domain-specific datasets such as Multi-XScience (Lu et al., 2020). Hybrid pipelines combine extractive and abstractive stages. A biomedical-focused system combines BERT-based extraction with a PEGASUS decoder for summarisation (Shinde et al., 2022), while SKT5SciSumm (To et al., 2024) uses SPECTER (Cohan et al., 2020) embeddings for clustering followed by T5-based generation, outperforming larger models like GPT-4 on some tasks. The 3A-COT framework (Zhang et al., 2024) structures LLM prompting into Attend-Arrange-Abstract stages to improve factuality and reduce redundancy. We adopted this framework in our setting with minor adjustments to generate a unified, structured output appropriate for our context.

We build on recent LLM-based MDSS advances, adapting them to topic modeling settings.

273

3 Task Definition

179

181

182

184

186

187

188

189

190

191

192

195

196

197

198

199

207

208

210

211

212

213

214

215

216

217

218

223

225

226

We define **Topic Description for Scientific Corpora** as the task of generating structured, interpretable summaries for topic model outputs. Each topic T_k is defined by:

- A set of topic-specific documents D_k ⊆ D, where each document is assigned to a single dominant topic,
- A ranked list of topic terms $W_k = \{w_1, \ldots, w_n\}.$

The goal is to generate a topic description S_k that summarises the main content of D_k , remains aligned with W_k , and follows a unified structure across topics. Each S_k contains a brief introduction to the topic, followed by the key research objects, methods, and purposes reflected in the underlying documents. This format supports comparable and structured exploration of scientific corpora.

We assess each S_k along four dimensions. Relevance requires that the description accurately reflects the key aspects of the topic by incorporating topic terms into meaningful context. Factuality ensures that the information is grounded in the original documents and does not introduce unsupported claims. Coherence refers to the logical flow and consistency of the description, ensuring it presents a unified explanation of the topic's main ideas. Fluency concerns the linguistic quality of the output; descriptions should use clear, accessible language that balances readability and technical precision.

4 Methodology

We adapt two LLM-based approaches for generating topic descriptions from sets of documents associated with each topic: *Selective Context Summarisation (SCS)*, which uses **Maximum Marginal Relevance** (MMR; (Carbonell and Goldstein, 1998)) to select a small, diverse subset of representative documents, and *Compressed Context Summarisation (CCS)*, which builds a hierarchical structure over all topic documents using recursive clustering and abstraction, following the tree-based indexing strategy of the **RAPTOR** framework (Sarthi et al., 2024). Both methods operate independently of the underlying topic modeling backbone.

In both pipelines, the generation process is guided by the same prompt template, adapted from the 3A-COT framework (Zhang et al., 2024), with topic terms provided as guidance. The full prompt is provided in Appendix A. An overview of the pipelines is shown in Figure 1.

4.1 Selective Context Summarisation (SCS)

SCS builds on an existing integration of LLMs into topic representation, as implemented in the BERTopic library¹. In the original implementation, representative documents for each topic are selected and passed to an LLM alongside topic terms to generate a short label. We extend this idea to generate informative topic descriptions that summarise the core content of each topic.

Given a topic, we select the ten highest-ranked terms and concatenate them to form a single string. This is then embedded using a pre-trained sentence embedding model. All documents within the topic are embedded in the same vector space and those most similar to the topic vector are retrieved.

To ensure the selected documents are both relevant and diverse, we apply Maximum Marginal Relevance (MMR; (Carbonell and Goldstein, 1998)). MMR iteratively selects documents that are similar to the topic vector while penalizing redundancy with respect to previously selected documents. This results in a representative and non-redundant subset of documents that captures the breadth of the topic and fits within the context window of the LLM.

In the generation process, we use the top 10 most representative documents and the top 10 most relevant topic terms for each topic. These are inserted into the shared prompt template (see Appendix A) and passed to the LLM, which generates the description based on this context.

4.2 Compressed Context Summarisation (CCS)

The Compressed Context Summarisation method is a slight adaptation of the tree-based indexing strategy from the RAPTOR framework (Sarthi et al., 2024), which constructs a recursive hierarchy of summaries through iterative clustering and abstraction. While RAPTOR is originally designed for retrieval over long documents, we use its core strategy to organize documents associated with each topic and generate descriptive summaries.

Unlike the original RAPTOR pipeline, which begins by segmenting long documents into smaller chunks, we start directly from the short documents already assigned to each topic (e.g., abstracts),

¹https://maartengr.github.io/BERTopic/getting_started/ representation/Ilm.html



Figure 1: Overview of the two topic description pipelines. *SCS* selects a representative subset of documents using MMR and summarises them with an LLM. *CCS* summarises all topic documents via hierarchical clustering and recursive abstraction.

without additional segmentation. These documents are embedded and projected into a lowerdimensional space using UMAP (McInnes et al., 2020) to improve clustering quality.

The projected embeddings are then clustered using Gaussian Mixture Models (GMMs), which support soft assignment, allowing documents to belong to multiple clusters. Each cluster is summarised using an LLM, with the top 10 topic terms provided at each stage for additional guidance. This produces an abstract summary that captures the main content of the clustered documents. These summaries are recursively re-embedded and re-clustered, forming a tree structure in which each internal node summarises its child nodes.

This recursive summarisation continues until only one cluster remains or no further abstraction is necessary. A final root node is added at the top of the tree, which is not part of the original RAP-TOR design, but is introduced in our adaptation. It serves as the output of the method: a topic description generated by the LLM that summarises the top-level content in the tree.

By including all topic documents and organizing them hierarchically, CCS sidesteps LLM context length limitations and produces descriptions grounded in the complete topic context. The prompt template used for all summarisation steps is the same one used in SCS.

5 Experimental Setup

To assess the effectiveness and generalisability of the proposed topic description pipelines, we conducted experiments across diverse scientific domains and topic modeling backbones. This section describes the datasets, modeling configurations, and models used for generation and embedding. 307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

329

330

331

332

333

334

5.1 Datasets

We evaluate our method on three domain-specific scientific corpora, using the abstracts of Englishlanguage research papers. Each dataset covers a distinct field to assess generalizability.

ACL Anthology The ACL Anthology² contains publications in computational linguistics and NLP from conferences such as ACL, EMNLP, and NAACL. We use the official GitHub version, extracting metadata and abstracts. Non-English entries and missing abstracts are removed, resulting in 52,126 clean abstracts.

NIPS Papers The NIPS Papers Dataset³ includes papers from the Neural Information Processing Systems (NIPS) between 1987 and 2016. We retain only English abstracts, removing missing entries and performing basic preprocessing. The final dataset contains 3,916 abstracts.

Quantum Computing Domain experts curated this dataset using a Boolean query on Scopus to retrieve recent papers (2010–2024) on quantum computing hardware. We retain only unique English abstracts, yielding 45,830 documents. Due to licensing restrictions, the dataset cannot be released; the full query is provided in Appendix B.

305

306

²https://github.com/acl-org/acl-anthology/tree/master/ python

³https://www.kaggle.com/datasets/benhamner/nips-papers

5.2 Topic Modeling

337

338

339

340

341

342

343

345

347

351

364

367

370

371

372

374

378

In order to provide a fair comparison testbed among different topic modeling approaches, we select three backbones: CTM, BERTopic and TopicGPT. Each variant builds upon quite distinct topic modeling methods, from classical bag-of-words statistical estimation (CTM), to plain clustering of vector representations of texts (BERTopic), up to straightforward multi-step zero-shot topic generation (TopicGPT). This choice composes a diverse set of setups with the goal of posing distinct levels and kinds of difficulty for creating topic descriptions. We apply each topic modeling method to each dataset, which leads to 9 topic models.

For all approaches involving training, we perform hyper-parameter optimization to find the best coherence and diversity metrics for each combination of topic model and dataset. For coherence, we use the Gensim implementation of the Coherence Model (Řehůřek and Sojka, 2010), specifically its default C_V metric (Röder et al., 2015). For diversity, we calculate the Inverted Rank-Biased Overlap (Webber et al., 2010; Terragni et al., 2021) of the top 10 keywords per topic.

In Appendix C, we show (Table 3) the scores and number of topics for each topic model, together with a full overview of their implementation and optimization details. Broadly, CTM has shown the best coherence, followed by BERTopic and TopicGPT. Conversely, TopicGPT has modeled in general the greatest number of topics, followed by BERTopic and then CTM. For which, following previous art (Grootendorst, 2022; Pham et al., 2024), we choose to attribute to each document just its most pronounced topic. This allows for better comparability among the two other backbones, which work with singleton topic inference.

5.3 LLM and Embedding Models

We use the DeepSeek-V3 (DeepSeek-AI et al., 2024) model to generate topic descriptions across all pipelines. For embedding-based retrieval, we use ModernBERT (Warner et al., 2024a), a competitive model for sentence-level semantic similarity.

6 Evaluation Strategy

Evaluating topic descriptions is inherently challenging due to the lack of gold-standard references
and the wide variation in topics across different
domains. We rely on reference-free evaluation metrics that assess quality without requiring human-

written summaries. We adopt two complementary strategies: **SUPERT**, a semantic similarity metric designed for multi-document summarisation, and an **LLM-as-a-Judge** framework, which uses prompting-based evaluation with LLMs. 384

385

386

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

6.1 SUPERT

SUPERT (Gao et al., 2020) is a reference-free evaluation metric developed for multi-document summarisation tasks. It creates a pseudo-reference by selecting key sentences from the input documents and compares generated summaries based on their semantic similarity to this reference. The similarity is computed using contextualized embeddings and soft token alignment. SUPERT has been shown to align well with human judgments of relevance, making it well-suited for assessing how much essential content is preserved in a topic description.

6.2 LLM-as-a-Judge

We build on recent work from the Eval4NLP 2023 Shared Task (Leiter et al., 2023), which explored prompting LLMs as explainable and referencefree evaluation metrics. Our setup is inspired by the best-performing system (Kim et al., 2023), which demonstrated that zero-shot prompting, finegrained scoring, and deterministic decoding lead to better alignment with human preferences.

To align evaluation with our task definition, we assess topic descriptions along four dimensions: **Relevance, Factuality, Coherence**, and **Fluency**. These criteria correspond to the aspects outlined in Section 3, and reflect the qualities expected from a high-quality topic description. We compute the **Mean Aspect Score (MAS)**, as the average across these four evaluation dimensions.

When selecting an LLM-as-a-judge model, we prioritized open-source models with strong alignment to human judgment. To this end, we chose Qwen2.5-7B-Instruct (Yang et al., 2024), which achieved the highest alignment among open-source models in the LLMEval benchmark (Gu et al., 2024). To account for variability in model outputs, we included two additional models. Our first choice was the Orca family, as both Orca-13B and OpenOrca-Platypus2-13B have shown promising alignment in prior studies (Kim et al., 2023; Leiter and Eger, 2024). However, due to their 4k context window limitations, we selected Mistral-7B-OpenOrca⁴, which maintains similar

⁴https://huggingface.co/Open-Orca/Mistral-7B-OpenOrca

Method	Metric	ACL			NIPS			Quantum		
		CTM	BERTopic	TopicGPT	CTM	BERTopic	TopicGPT	CTM	BERTopic	TopicGPT
SCS	SUPERT	0.475	0.477	0.508	0.459	0.465	0.519	0.489	0.486	0.557
	MAS-Qwen	50.645	62.561	72.004	56.703	71.962	81.612	57.088	64.319	78.650
CCS	SUPERT	0.467	0.474	0.501	0.453	0.469	0.515	0.458	0.484	0.552
	MAS-Qwen	49.612	62.400	65.299	58.766	66.295	78.108	56.647	62.839	78.829

Table 1: SUPERT & MAS-Qwen scores across methods, datasets, and topic modeling backbones

alignment while supporting longer contexts (32k). As a third model from a different architecture line, we added Gemma-3-27B (Kamath et al., 2025) to ensure diversity across the various model families.

As it is not possible to evaluate a generated description against all documents associated with a topic at once due to the limited context window of LLMs, we instead sample 5 random draws of 10 documents each from the full topic set. Each batch is evaluated independently, and we report the mean score across the five runs. This approach reflects a more realistic human evaluation scenario, where annotators are unlikely to read all the documents in a large collection. Moreover, it aligns with a key assumption in topic-level summarisation, where a strong topic representation should capture the central content of the topic and remain consistent and relevant across different subsets of its documents. Appendix D lists the evaluation prompts.

7 Results

432

433

434

435

436

437

438

439

440

441

442 443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466 467

468

469

470

471

In this section, we present SUPERT scores and Mean Aspect Score (MAS) from the LLM-as-a-Judge evaluation, using Qwen-2.5-7B-Instruct as our primary model. We confirm that trends hold across Mistral-7B-OpenOrca and Gemma-3-27B, showing robustness across model families. We first examine overall pipeline effectiveness, then analyze how topic size affects description quality.

7.1 Performance Across Domains and Backbones

We compare the two topic description pipelines across datasets and topic modeling backbones. The results, shown in Table 1, demonstrate a clear and consistent advantage for SCS. It achieves the highest SUPERT and MAS scores in almost all configurations, highlighting its robustness across domains and backbone models. CCS performs competitively, achieving strong SUPERT scores in several configurations, but slightly falls behind SCS on MAS in most settings. To validate the consistency of the evaluation results across LLMs, we measured the correlation between the MAS scores produced by the three judge models using Kendall's tau-b. The results demonstrate a strong agreement between Qwen-2.5-7B-Instruct and Gemma-3-27B, and moderate agreement across the other model pairs, as shown in Table 2. Moreover, the MASs show 472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

Judge Models	$ au_b$	p
Qwen & Gemma	0.7255	$4.304\cdot10^{-6}$
Qwen & Mistral	0.5033	$2.99\cdot 10^{-3}$
Mistral & Gemma	0.5163	$2.24 \cdot 10^{-3}$

Table 2: Kendall's τ_b correlations between MAS scores of judge models.

consistent behavior across document draws. For transparency, we include a detailed presentation of MAS on each document set draw in Appendix F.

7.2 Effect of Topic Size

To better understand how topic characteristics impact description quality, we analyze the effect of topic size on MAS distributions for SCS and CCS, cross-validated with topic cohesion (mean cosine distance among topic documents).

Figure 2 shows the distribution of winners among the probed pipelines by topic size quartile. It is noteworthy that SCS achieves the highest ranking among the first, second, and third topic size quartiles. The only exception is the Large category, where CCS matches SCS with an equal number of wins. Additionally, while the number of SCS wins tends to decline as topic size increases, CCS shows an upward trend from Small to Large categories, matching SCS in the largest quartile.

Cross-validation against topic cohesion confirms that description quality remains remarkably consistent across all topic cohesion quartiles for both SCS and CCS, indicating that these approaches are robust to variation in topical coherence and that



Figure 2: Winner count on LLM-Eval MAS per topic size quartile over all topic models.

the observed size effects above are not confounded by cohesion variations. We provide a thorough presentation against topic cohesion in Appendix G. Appendix H shows SUPERT-based results by topic size & cohesion, showing a similar trend to MAS. Appendix E reports Kendall's τ_b correlations between LLMs on topic size preferences.

8 Discussion

504

505

507

510

511

514

516

518

519

520

521

527

529

This section discusses both pipeline effectivenessand highlights trends by topic size and structure.

8.1 Selective vs. Compressed Approaches to Topic Description

Effectiveness Advantage of Selective Sampling Our results demonstrate a consistent effectiveness advantage for the Selective Context Summarisation (SCS) pipeline across multiple datasets and topic modeling backbones. The MMR selection process in SCS provides a balanced set of relevant and diverse documents, creating focused yet comprehensive input for the LLM. This selective approach seems to reduce noise from peripheral documents while ensuring core topic terms remain prominent throughout the description generation process. In scientific corpora, we hypothesize that this advantage may be amplified, since documents on the same research topic often share similar objects of study, purposes, and methodologies.

Limitations of Hierarchical Compression
CCS's hierarchical structure, despite its theoretical
capacity to process entire document sets, suffers
from what we term "error propagation" and
"keyword attrition." As abstractions build upward
through the tree, inaccuracies at lower levels can

amplify in subsequent steps, while important terminology may become diluted during recursive summarisation. These phenomena likely contribute to CCS's generally lower effectiveness across our evaluation metrics. From an efficiency standpoint, SCS demonstrates a superior compute-to-quality ratio, requiring only a single document set pass compared to CCS's multiple rounds of embedding, clustering, and LLM calls. The stability of SCS effectiveness across different topic modeling backbones (CTM, BERTopic, and TopicGPT) further highlights its robustness as a general-purpose topic description method that can integrate with existing topic modeling workflows regardless of their underlying approach. 537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

8.2 Scalability and Topic Size Effects

Size-Dependent Effectiveness Patterns Analysis of topic size effects reveals an intriguing pattern: while SCS dominates for small to medium-sized topics, CCS becomes competitive and even outperforms SCS for the largest topics (4th quartile), as shown in Figure 2. This finding highlights important scalability considerations for topic description applications. For smaller topics, SCS effectively identifies a representative subset that captures the topic's essence. However, as topics grow larger, the fixed selection size (10 documents in our setup) becomes limiting. When topics contain hundreds of documents, even carefully selected subsets may miss important sub-themes or variations. CCS shows a valuable property for larger topics: its hierarchical summarisation approach scales with topic size, preserving coverage of diverse sub-themes that fixed-size selection may miss.

Effectiveness Nuances Across Size Deciles The relationship between topic size and method effectiveness shows additional nuance when examined at finer granularity. Figure 3 displays MAS per topic size decile. Notably, SCS demonstrates a consistent dominance in quality across the initial six deciles. CCS then assumes the lead for the seventh and eighth deciles, before SCS regains dominance for the largest topics. This suggests that while CCS outperforms SCS for some larger topics, it also has a saturation point, likely due to a bottleneck in the hierarchical compression of information. This scale-dependent effectiveness suggests that practical applications might benefit from a hybrid approach that adaptively selects between methods based on topic size. Our analysis confirms that



Figure 3: Mean Aspect Score per topic size decile.

these patterns persist when controlling for topic cohesion, indicating that the observed effects are genuinely related to scale. This highlights topic size as a key factor in designing and evaluating topic description pipelines for scientific corpora.

8.3 Qualitative Analysis

588

591

593

594

595

601

To complement our quantitative results, we conducted a targeted qualitative analysis of 45 topic descriptions. We examined 15 top-scoring, 15 lowscoring and 15 descriptions with diverging SU-PERT and LLM-as-Judge scores. This enabled us to examine the behaviours of the methods beyond aggregate metrics. Examples illustrating content quality across different models and methods are provided in Appendix I.

Characteristics of Selective Context Summarisation Our analysis reveals that SCS consistently generates clear and coherent summaries in highscoring cases, with strong alignment to the provided topic terms and good coverage of central concepts (see Example 1). It demonstrates notable resilience to incoherence in topic terms (Example 2), as any inconsistencies in the topic terms do not compound through multiple summarisation lay-610 ers. SCS descriptions maintain coherence across different datasets and topic modeling backbones, indicating robust transferability. However, in low-613 performing cases, particularly when topic terms are 614 overly general or lossy, the method tends to pro-615 duce generic or shallow outputs. This limitation is 616 617 exacerbated when the selected representative documents contain primarily general knowledge rather 618 than specific insights. We observe that SCS cap-619 italizes on well-selected topic terms from the underlying model, creating a synergistic effect where 621

strong topic models yield better descriptions.

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

Characteristics of Compressed Context Summarisation CCS exhibits distinctive strengths in handling complex or technical topics, often producing more detailed descriptions than SCS. However, this method shows lower alignment with the original topic terms, in several cases, generating dense and nuanced content that only partially connects to the provided terms. This misalignment creates challenges in verifying how faithfully the description represents the intended topic (see Example 3). The hierarchical summarisation approach in CCS appears to struggle with effectively prioritizing the most important content, often resulting in information overflow manifested as lengthy lists or excessive detail. This limitation stems from "document grounding distance" in effect in the hierarchical summarisation process, which may not optimally distinguish central from peripheral information. Finally, CCS demonstrates greater sensitivity to topic term quality, with more frequent failures when topic terms are incoherent (see Example 4 comparatively to Example 2).

9 Conclusion & Future Work

We introduced the task of Topic Description for Scientific Corpora, which aims to create structured, document-based summaries that go beyond term lists. To address this, we adapted two LLMbased pipelines: Selective Context Summarisation (SCS) and Compressed Context Summarisation (CCS). SCS consistently achieved better performance across datasets and topic modeling backbones. CCS showed advantages for large topics due to its scalable, recursive structure. Our findings highlight a trade-off between selective and compressed strategies. SCS excels in precision and stability, while CCS offers broader coverage for large-scale topics. Together, they provide practical foundations and insights for developing interpretable topic representations in scientific domains.

This work suggests directions for further exploration, including methodological improvements and practical applications. Instead of single-vector retrieval of SCS, future work could examine more fine-grained retrieval strategies to improve coverage and adaptability for complex or broad topics. Our evaluation strategy combines SUPERT and LLM-as-a-Judge; future research could investigate alternative setups based on automatic factuality and coverage check-ups with retrieval techniques.

701

703

710

Limitations

Despite our multi-faceted evaluation strategy, several limitations remain. First, we do not include hu-674 man assessment. Although we combine SUPERT 675 and LLM-as-a-Judge to approximate quality, expert feedback is essential, especially in scientific domains where interpretability and factual accuracy require domain knowledge. The use of both 679 SUPERT and LLM-based evaluation offers complementary strengths: SUPERT captures content 681 relevance via semantic similarity, while LLM-as-a-Judge enables structured, fine-grained evaluation. This dual setup mitigates some metric-specific bi-684 ases, though it cannot fully substitute for human judgment.

> This challenge is compounded by limitations in the topic modeling stage itself. The quality of topic descriptions is directly tied to the coherence and relevance of the underlying topics and their terms. Despite optimization, CTM often produced noisy or domain-unspecific topics. Similarly, TopicGPT occasionally generated topics that were overly broad or narrowly scoped. These issues degraded the resulting descriptions, even with grounded generation. This reliance on topic model quality is another central limitation in this present study. Still, such limitations are inherent to real-world applications (academic and industrial alike) when attempting to gain an overview of large-scale (scientific) corpora. Our analysis reflects these challenges rather than avoids them.

In addition, while the chosen LLMs are among the strongest models, their outputs remain sensitive to prompt design and can hallucinate content. Our pipelines use a fixed 3A-COT-derived prompting strategy, but prompt wording significantly affects LLM output. No ablation or robustness analysis was conducted to assess this sensitivity. Also, even strong LLMs are prone to hallucination, especially when context is sparse or ambiguous. This is only partially mitigated by the factuality criterion in our LLM-as-a-Judge evaluation.

Finally, our evaluation is confined to abstracts
in English-language scientific corpora. This raises
questions about the generalisability of the approach
to full-text documents, other genres such as patents,
or non-English data.

References

Areej Alokaili, Nikolaos Aletras, and Mark Stevenson. 2020. Automatic generation of topic labels. *CoRR*, abs/2006.00127. 719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

- Hosein Azarbonyad, Zubair Afzal, and George Tsatsaronis. 2023. Generating topic pages for scientific concepts using scientific publications. In Advances in Information Retrieval - 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2-6, 2023, Proceedings, Part II, volume 13981 of Lecture Notes in Computer Science, pages 341–349. Springer.
- Ibrahim Al Azhar, Venkata Devesh Reddy, Hamed Alhoori, and Akhil Pandey Akella. 2025. Limtopic: Llm-based topic modeling and text summarization for analyzing scientific articles limitations. *CoRR*, abs/2503.10658.
- Shraey Bhatia, Jey Han Lau, and Timothy Baldwin. 2016. Automatic labelling of topics with neural embeddings. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 953–963, Osaka, Japan. The COLING 2016 Organizing Committee.
- Federico Bianchi, Silvia Terragni, and Dirk Hovy. 2021. Pre-training is a hot topic: Contextualized document embeddings improve topic coherence. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pages 759–766, Online. Association for Computational Linguistics.
- David M Blei. 2012. Probabilistic topic models. Communications of the ACM, 55(4):77-84.
- Jaime G. Carbonell and Jade Goldstein. 1998. The use of mmr, diversity-based reranking for reordering documents and producing summaries. In SIGIR '98: Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, August 24-28 1998, Melbourne, Australia, pages 335–336. ACM.
- Jonathan D. Chang, Jordan L. Boyd-Graber, Sean Gerrish, Chong Wang, and David M. Blei. 2009. Reading tea leaves: How humans interpret topic models. In Advances in Neural Information Processing Systems 22: 23rd Annual Conference on Neural Information Processing Systems 2009. Proceedings of a meeting held 7-10 December 2009, Vancouver, British Columbia, Canada, pages 288–296. Curran Associates, Inc.
- Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel S. Weld. 2020. SPECTER: document-level representation learning using citationinformed transformers. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, pages 2270–2282. Association for Computational Linguistics.

776

- 786 787
- 790 791
- 795 796
- 799

805 806

809

810 811

813

815 816 817

818 819

822

823 824 825

826 827

829

832

833

- DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, and 81 others. 2024. Deepseek-v3 technical report. CoRR, abs/2412.19437.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Yang Gao, Wei Zhao, and Steffen Eger. 2020. SUPERT: Towards new frontiers in unsupervised evaluation metrics for multi-document summarization. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 1347-1354, Online. Association for Computational Linguistics.
- Maarten Grootendorst. 2022. Bertopic: Neural topic modeling with a class-based tf-idf procedure. Preprint, arXiv:2203.05794.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Yuanzhuo Wang, and Jian Guo. 2024. A survey on llm-as-a-judge. CoRR, abs/2411.15594.
- Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean-Bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, Gaël Liu, and 79 others. 2025. Gemma 3 technical report. CoRR, abs/2503.19786.
- JoongHoon Kim, Sangmin Lee, Seung Hun Han, Saeran Park, Jiyoon Lee, Kiyoon Jeong, and Pilsung Kang. 2023. Which is better? exploring prompting strategy for llm-based metrics. In Proceedings of the 4th Workshop on Evaluation and Comparison of NLP Systems, pages 164-183.
- Jey Han Lau, Karl Grieser, David Newman, and Timothy Baldwin. 2011. Automatic labelling of topic models. In Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, pages 1536–1545, Portland, Oregon, USA. Association for Computational Linguistics.
- Christoph Leiter and Steffen Eger. 2024. Prexme! large scale prompt exploration of open source llms for machine translation and summarization evaluation. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 11481-11506.

Christoph Leiter, Juri Opitz, Daniel Deutsch, Yang Gao, Rotem Dror, and Steffen Eger. 2023. The eval4nlp 2023 shared task on prompting large language models as explainable metrics. In Proceedings of the 4th Workshop on Evaluation and Comparison of NLP Systems, pages 117–138.

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

- Yao Lu, Yue Dong, and Laurent Charlin. 2020. Multi-XScience: A large-scale dataset for extreme multidocument summarization of scientific articles. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 8068-8074, Online. Association for Computational Linguistics.
- Leland McInnes, John Healy, and Steve Astels. 2017. hdbscan: Hierarchical density based clustering. The Journal of Open Source Software, 2(11):205.
- Leland McInnes, John Healy, and James Melville. 2018. Umap: Uniform manifold approximation and projection for dimension reduction. arXiv preprint arXiv:1802.03426.
- Leland McInnes, John Healy, and James Melville. 2020. Umap: Uniform manifold approximation and projection for dimension reduction. Preprint, arXiv:1802.03426.
- Zach Nussbaum, John X. Morris, Brandon Duderstadt, and Andriy Mulyar. 2024. Nomic embed: Training a reproducible long context text embedder.
- Chau Pham, Alexander Hoyle, Simeng Sun, Philip Resnik, and Mohit Iyyer. 2024. TopicGPT: A promptbased topic modeling framework. 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 2956-2984, Mexico City, Mexico. Association for Computational Linguistics.
- Cristian Popa and Traian Rebedea. 2021. BART-TL: weakly-supervised topic label generation. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021, pages 1418-1425. Association for Computational Linguistics.
- Radim Řehůřek and Petr Sojka. 2010. Software framework for topic modelling with large corpora.
- Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERTnetworks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3982-3992, Hong Kong, China. Association for Computational Linguistics.
- Michael Röder, Andreas Both, and Alexander Hinneburg. 2015. Exploring the space of topic coherence measures. In Proceedings of the eighth ACM international conference on Web search and data mining, pages 399-408.

948

949

950

951

- 969 970 971 972 973
- 973 974

- Domenic Rosati. 2022. Moving beyond word lists: towards abstractive topic labels for human-like topics of scientific documents. *CoRR*, abs/2211.05599.
- Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D. Manning.
 2024. RAPTOR: recursive abstractive processing for tree-organized retrieval. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.

895

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

925

926

927

928

930

931

934

935

936

937

938

939

942

943

947

- Kartik Shinde, Trinita Roy, and Tirthankar Ghosal. 2022. An extractive-abstractive approach for multidocument summarization of scientific articles for literature review. In *Proceedings of the Third Workshop on Scholarly Document Processing*, pages 204–209, Gyeongju, Republic of Korea. Association for Computational Linguistics.
 - Silvia Terragni, Elisabetta Fersini, and Enza Messina. 2021. Word embedding-based topic similarity measures. In *International conference on applications* of *Natural Language to information systems*, pages 33–45. Springer.
- Huy Quoc To, Ming Liu, Guangyan Huang, Hung-Nghiep Tran, Andr'e Greiner-Petter, Felix Beierle, and Akiko Aizawa. 2024. Skt5scisumm – revisiting extractive-generative approach for multidocument scientific summarization. *Preprint*, arXiv:2402.17311.
- Pancheng Wang, Shasha Li, Kunyuan Pang, Liangliang He, Dong Li, Jintao Tang, and Ting Wang. 2022.
 Multi-document scientific summarization from a knowledge graph-centric view. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 6222–6233, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, Nathan Cooper, Griffin Adams, Jeremy Howard, and Iacopo Poli. 2024a. Smarter, better, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetuning and inference. *CoRR*, abs/2412.13663.
- Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, Nathan Cooper, Griffin Adams, Jeremy Howard, and Iacopo Poli. 2024b. Smarter, better, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetuning and inference. *Preprint*, arXiv:2412.13663.
- William Webber, Alistair Moffat, and Justin Zobel. 2010. A similarity measure for indefinite rankings. ACM Transactions on Information Systems (TOIS), 28(4):1– 38.
- Wen Xiao, Iz Beltagy, Giuseppe Carenini, and Arman Cohan. 2022. PRIMERA: Pyramid-based masked

sentence pre-training for multi-document summarization. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5245–5263, Dublin, Ireland. Association for Computational Linguistics.

- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, and 22 others. 2024. Qwen2.5 technical report. *CoRR*, abs/2412.15115.
- Yongbing Zhang, Shengxiang Gao, Yuxin Huang, Zhengtao Yu, and Kaiwen Tan. 2024. 3a-cot: An attend-arrange-abstract chain-of-thought for multidocument summarization. *International Journal of Machine Learning and Cybernetics*.

A Summarisation Prompts

We used the deepseek-v3 model to generate topic descriptions across all methods. To ensure consistency and structure in the outputs, we define a fixed system message and adopt a 3-step prompting framework inspired by the 3A-COT method (Zhang et al., 2024). This includes *attending* to key aspects, *arranging* extracted information, and generating the final *abstract*. The exact prompt templates used are provided below.

System Prompt

You are a scientific research assistant who organizes information into structured markdown documents. Your writing style sounds natural and professional. Avoid using Marketing and HR language.

Prompt 1: System prompt used for topic description generation

Attend Prompt

[DOCUMENTS]

What are the research purposes in this document?

What are the research object in this document? What are the research methods in this document? What are the research result in this document? What are the main findings in this document?

Please answer the above questions:

Prompt 2: Attend prompt for extracting key information

Arrange Prompt

[ATTEND_OUTPUT]

Organize the above important information. Arrange this information in a logical order or relevance to build a coherent narrative, and consider how information from different articles can be combined to complement and connect with each other.

Prompt 3: Arrange prompt for structuring extracted content

Abstract Prompt

[ABSTRACTS]

[ARRANGE_OUTPUT]

Based on the above abstracts, key information and the keywords: {topic_words}, write a summary.

Make sure to include key information, research objectives and ideas. The summary should be structured as clean MARKDOWN with ONLY the following Headings:

Brief Introduction into the Topic, Key Research Objects, Key Research Methods, Key Research Purpose.

Each Heading should have only Keypoints listed. Avoid the use of additional MARKDOWN subsections. Avoid adding your own opinion, interpretation, or conclusions or Future Work. Use the information provided in the text only.

Prompt 4: Abstract prompt for final topic description generation

975

976

977

980

B **Quantum Dataset Query**

In Query 1, we present the full boolean query used for collecting the source documents for the Quantum Computing dataset. Specially, the query is specialized in the hardware part of this scientific field.

Topic Models С

This section presents implementation details and results of three topic modeling approaches used in our comparative analysis.

Boolean Query

TITLE-ABS-KEY ("quantum comput*" OR "quantum processor" OR "quantum circuit" OR "quantum logic gate" OR "quantum gate" OR "logical qubit" OR qubit OR "quantum system" OR "quantum information processing" OR "quantum control" OR "quantum electronics" OR "quantum hardware" OR "noisy intermediate-scale quantum era" OR "NISQ" OR "multiqubit circuit" OR "quantum simulation" OR "quantum simulator") AND TITLE-ABS-KEY ((cryogen* OR "magnetic field" OR laser OR photoluminescence OR silicon OR "electric fields" OR magnetism OR fluorescence) OR ("neutral atom" OR "cold atom" OR "trap*atom" OR "atom trap" OR "rydberg" OR atoms OR "optical lattice*" OR magic OR "optical tweezer*" OR strontium OR ytterbium OR "photonic crystal fibre") OR ("ion traps" OR "trapped ions" OR "ions" OR "integrated waveguide" OR "laser induced deep etching" OR "on-chip coupling") OR (superconduct* OR "SQUIDs" OR "Josephson junction device*" OR "indium bump" OR "NbN films" OR "single flux quantum" OR "quantum flux" OR "SQUID") OR (center OR diamond OR "NV center" OR "NV centre" OR "color centre" OR "colour center" OR "silicon vacancy centre" OR "silicon vacany center") OR (photon* OR "gaussian boson sampl*" OR "squeezed light source" OR niobate OR "superconducting nanowire singlephoton detector" OR "SNSPD") OR (topology OR "topological quantum computing" OR "topological insulator*") OR (semiconductor OR "molecular beam epitaxy" OR "semiconducting*" OR "crystal lattice*" OR phonons)) AND PUBYEAR > 2009 AND PUBYEAR < 2026

Query 1: Boolean query used for collecting the source documents for the Quantum Computing dataset.

C.1 Implementation Details

12

CTM In the CTM backbone, we use the Gihub implementation⁵ of the original contribution (Bianchi et al., 2021). Here, we choose to optimize over four hyper-parameters: number of topics (40-100), number of epochs (10-50), activation function ({sigmoid, relu, softplus}), number of neurons (100-500). All other hyper-parameters use

985

991

⁵https://github.com/MilaNLProc/contextualized-topicmodels

standard values from the implementation.

BERTopic For this approach, we use the known BERTopic package⁶. This standard pipeline con-995 sists of mainly three stages: Embedding (Em.) stage, Dimensionality Reduction (DR) stage, the 997 Clustering (Cl.) stage, and Topic Representation (TR) stage. For the Em stage, we use the nomic-ai/modernbert-embed-base⁷ model (Nuss-1000 baum et al., 2024), which is an embedding model 1001 trained on the ModernBERT (Warner et al., 2024b) 1002 encoder. For the DR and Cl. stages, we opt for the 1003 standard pairing with UMAP (McInnes et al., 2018) and HDBSCAN (McInnes et al., 2017). Finally, in 1006 the TR stage, we use class-TFIDF, which is firstly 1007 introduced in (Grootendorst, 2022). Overall, in this approach, we have also four hyper-parameters: 1008 UMAP - number of neighbors (5-50), number 1009 of components (2-15) and min. distance (0.0-1010 0.5); HDBSCAN - min. cluster size (10-50). For 1011 UMAP, we fix the metric to cosine, and euclidean 1012 for HDSCAN. All other hyper-parameters use stan-1013 dard values from their implementations. 1014

TopicGPT We follow the original TopicGPT 1015 pipeline (Pham et al., 2024), using the open-source 1016 implementation available at GitHub⁸ and altering only the document-assignment stage to align 1018 with BERTopic and CTM. For topic generation, 1019 we randomly sample 1,000 documents from each 1020 dataset and leverage GPT-4 to propose an initial set of top-level topics, which we then iteratively 1022 refine into subtopics to build a complete hierar-1023 chical structure. In the subsequent assignment phase—applied to the full datasets—we replace 1025 1026 TopicGPT's default routine (which, for each document, prompts GPT-3.5-turbo with the finalized 1027 hierarchy and returns the best-matching topic with 1028 a supporting quote) with a two-part prompt to GPT-3.5-turbo: (i) assign each document to its best-matching topic in our hierarchy; and (ii) ex-1031 tract ten representative keywords per document. 1032 Finally, we post-process all extracted keywords for each topic by tokenizing them on whitespace, converting to lowercase, stripping punctuation, ag-1035 gregating token frequencies, and selecting the ten 1036 most frequent tokens per topic-thereby exactly 1037 matching the output format of our BERTopic and 1038 CTM backbones. 1039

Own Assignment Prompt Template

You will receive a document and a topic hierarchy. Assign the document to the most relevant topic of the hierarchy. Then, output the topic label, and supporting keywords from the document. DO NOT make up new topics or keywords.

[Topic Hierarchy]

{tree}

[Instructions]

1. Topic label must be present in the provided topic hierarchy. You MUST NOT make up new topics.

2. The keywords must be taken from the document. You MUST NOT make up keywords or quotes. All keywords MUST NOT contain stop words.

[Document]

{Document}

Double check that your assignment exists in the hierarchy! Your response should be in the following format:

[Topic Level] Topic Label: keyword1, keyword2, etc

Your response:

Prompt 5: Prompt template used for document-to-topic assignment in the TopicGPT adaptation.

C.2 Results

Table 3 presents topic modeling evaluation re-1041 sults across three datasets (ACL, NIPS, and Quan-1042 tum) for three different topic modeling approaches: 1043 CTM, BERTopic, and TopicGPT. The evaluation 1044 metrics used in the comparison are Coherence, Di-1045 versity, and Number of Topics (N.Topics). CTM 1046 consistently achieves the highest coherence scores across all three datasets (0.664 for ACL, 0.601 for 1048 NIPS, and 0.692 for Quantum). It also maintains 1049 high diversity scores above 0.94 for all datasets. 1050 BERTopic shows moderate coherence performance 1051 (0.504 for ACL, 0.458 for NIPS, and 0.546 for 1052 Quantum), with somewhat lower diversity met-1053 rics, particularly for the Quantum dataset (0.799). 1054 TopicGPT demonstrates coherence scores between 0.458 and 0.526 across datasets, with strong diver-1056

⁶https://maartengr.github.io/BERTopic/index.html

⁷https://huggingface.co/nomic-ai/modernbert-embed-base

⁸https://github.com/chtmp223/topicGPT

sity in the NIPS dataset (0.963) but lower diver-1057 sity for ACL (0.881) and Quantum (0.809). Re-1058 garding the number of topics identified, TopicGPT 1059 produces substantially more topics than the other 1060 approaches, particularly for NIPS (276). BERTopic 1061 identifies the fewest topics overall with just 24 for 1062 the NIPS dataset. For the ACL dataset, the number 1063 of topics is more consistent across models (CTM: 1064 72, BERTopic: 70, TopicGPT: 66). The Quan-1065 tum dataset shows moderate variation, with CTM 1066 identifying 59 topics, BERTopic 72, and TopicGPT 1067 significantly more at 169. 1068

D Evaluation Prompts

We evaluate summaries along four dimensions: relevance, coherence, factuality, and fluency. Each is scored independently using a dedicated prompt, detailed below.

1074 Aspect Definitions

1069

1071

1072

1073

1075

1076

1077

1078

1079

1080

1082

1083

1084

1085

1086

1087

1088

1089

1090

1091

1092

1094

1095

1097

1098

1099

1100

1101

- **Relevance:** The rating measures how well the summary captures the key points of the documents. Consider whether all and only the important aspects are contained in the summary.
- **Coherence:** This rating evaluates how seamlessly the sentences of the summary flow together, creating a unified whole. Assess how smoothly the content transitions from one point to the next, ensuring it reads as a cohesive unit.

• Factuality: This rating gauges the accuracy and truthfulness of the information presented in the summary compared to the original documents. Scrutinize the summary to ensure it presents facts without distortion or misrepresentation, staying true to the source content's details and intent.

• Fluency: This rating evaluates the clarity and grammatical integrity of each sentence in the summary. Examine each sentence for its structural soundness and linguistic clarity.

E Gemma-3-27B & Mistral-7B-OpenOrca Results

To complement the main results, we report the MAS obtained using Gemma-3-27B and Mistral-7B-OpenOrca in Table 4. These models provide

Evaluation Prompt Template

Instruction:

In this task you will evaluate the quality of a summary written for multiple documents.

To correctly solve this task, follow these steps:

1. Carefully read the document, be aware of the information it contains.

2. Read the proposed summary.

3. Rate each summary on a scale from 0 (worst) to 100 (best) by its {aspect}. Decimals are allowed.

Definition:

{definition}

Source documents:

{source}

Summary: {summary}

Score:

Prompt 6: Evaluation prompt template used for scoring topic descriptions across relevance, factuality, coherence, and fluency

additional perspectives on the quality of the generated descriptions and help verify the consistency of trends observed with Qwen-2.5-7B-Instruct. 1102

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

To further assess inter-model agreement, we compute Kendall's τ_b between the rankings of method-size combinations (i.e., CCS/SCS across the four topic size categories: *Small*, *Medium-Small*, *Medium-Large*, and *Large*) for each pair of judge models. We evaluate agreement across the full 8-item ranking. This provides a single τ_b score per pair, reflecting overall alignment in method preferences across topic sizes. As shown in Table 5, Qwen2.5-7B-Instruct aligns moderately to strongly with both Gemma-3-27B and Mistral-7B-OpenOrca, while Gemma-3-27B and Mistral-7B-OpenOrca exhibit weaker agreement.

F Impact of drawn documents in LLM-Eval

In order to analyze the impact of using subsets of
documents of topics as reference documents in the
LLM-Eval strategies, we present a detailed visual-
ization of the Quantum dataset results in Figure 4
across all five document draws for each TM and1120
1121

ТМ	Dataset	Coherence	Diversity	N.Topics
	ACL	0.664	0.994	72
CTM	NIPS	0.601	0.949	38
	Quantum	0.692	0.996	59
	ACL	0.504	0.972	70
BERTopic	NIPS	0.458	0.930	24
	Quantum	0.546	0.799	72
	ACL	0.458	0.881	66
TopicGPT	NIPS	0.472	0.963	276
	Quantum	0.526	0.809	169

Table 3: Topic modeling evaluation results across three scientific datasets.

Method	Metric	ACL			NIPS			Quantum		
		CTM	BERTopic	TopicGPT	СТМ	BERTopic	TopicGPT	CTM	BERTopic	TopicGPT
SCS	MAS-Mistral	74.913	74.306	76.307	75.497	75.109	77.240	72.077	73.127	76.120
	MAS-Gemma	81.537	85.163	87.411	82.691	83.865	89.400	82.006	85.255	89.395
CCS	MAS-Mistral	73.092	71.207	74.471	66.384	70.771	75.201	65.031	72.758	74.528
	MAS-Gemma	81.483	85.155	86.799	81.194	83.370	88.549	80.204	85.120	89.125

Table 4: MAS scores across methods, datasets, and topic modeling backbones using Mistral and Gemma as judge models.

Model Pair	$ au_b$	p
Qwen & Gemma	0.6183	0.0340
Qwen & Mistral	0.6910	0.0178
Gemma & Mistral	0.2857	0.3988

Table 5: Kendall's τ_b between full method–size rankings of each model pair.

TD approach. As for the other datasets, a similar scenario holds.

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137 1138

1139

1140

1141

1142

From visual inspection of Figure 4, we observe that scores remain relatively stable across different document draws for the same TM and TD method. When fixing a topic modeling approach and a topic description pipeline, the fluctuations in LLM-Eval MAS are generally small, with most variations remaining within 5 points to the mean on our 100point scale.

While a comprehensive variance analysis across all datasets would provide further statistical rigor, the consistency observed in the Quantum dataset suggests that our sampling approach produces reliable evaluations. The observed stability indicates that randomly sampling 10 documents five times provides a reasonable approximation of how a topic description would be evaluated against the full document collection.

The observed consistency across document draws supports our decision to use this sampling approach as a practical solution to the context window limitations of LLMs. While a more exhaustive analysis would be valuable for future work, the current evidence suggests that our methodology provides reliable evaluations of topic descriptions despite using only subsets of the complete document collections. 1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

G Effect of Topic Cohesion on Mean Aspect Score (MAS)

To study the impact of topic cohesion on the quality 1155 of topic descriptions, We compute the mean cosine 1156 embedding distance among all documents for each 1157 topic. We call this indicator "Topic Cohesion." Fig-1158 ure 5 shows the MAS distributions for all topics 1159 grouped by their topic cohesion quartile. Interest-1160 ingly, topic cohesion plays an almost negligible 1161 role in the MAS distributions across all quartiles. 1162 There is a clear downward trend indicating the an-1163 ticipated TD quality degradation towards topics 1164 of low cohesion. However, this effect is minor 1165 among all TD approaches, only becoming more 1166 pronounced in the low cohesion quartile. Even 1167 there, the best topic descriptions of the two best 1168





approaches, SCS and CCS, are competitive with TD's best scores in the more cohesive quartiles.

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

H Effect of Topic Size and Cohesion on SUPERT

Figure 6 shows the distribution of winners per topic size category based on the SUPERT metric. SCS leads in the first and second quartiles, with CCS gaining a slight edge in the third. In contrast to MAS-Qwen, which shows CCS catching up in the largest category, SUPERT continues to favor SCS in the fourth quartile. This suggests that SCS is more aligned with SUPERT's relevance-focused evaluation, even as topic size increases.



Figure 6: Winner count on SUPERT per topic size quartile over all topic models

Figure 7 shows SUPERT distributions grouped1182by topic cohesion (mean pairwise cosine distance).1183Both SCS and CCS remain relatively stable across1184cohesion levels, with only slight performance1185degradation from high to low cohesion.1186





Figure 4: LLM-Eval MAS for every draw of 10 documents per topic.



Figure 7: Box-plots of SUPERT distributions conditioned on topic cohesion quartile



Figure 8: SUPERT Score per topic size decile.

The SUPERT experiments on effectiveness nuances across size deciles, to be seen in Figure 8 reveal patterns consistent with the results of LLM-Eval in Section 8.2: method effectiveness remains tied to topic size, with SCS dominating slightly for smaller topics (first to sixth decile). The effect of CCS slightly outperforming SCS in LLM-Eval for the larger topics (seventh and eigth deciles) is not to be seen with SCS even achieving minially better scores with SUPERT. Still, the overall trend is replicated from LLM-Eval and this replication further confirms that scale, besides intrinsic topic quality, drives performance differences and emphasizes the viability of adaptive hybrid approaches.

1187

1188

1189

1190

1191

1192 1193

1194

1195

1196

1197

1198

1199

1200

1201

1202

I Examples of generated Topic Descriptions

1203To demonstrate the results of our topic description1204pipeline, we present a selection of examples cov-1205ering different datasets and modelling approaches.1206Each example provides metadata, the generated

summary, and evaluation scores. Together, these1207elements help contextualise the content and quality1208of the generated descriptions.1209

acl-53

Dataset: ACL Model: TopicGPT Method: SCS Category: Best Topic Words: data, LOD, digital, linguistics, studies, linked open, humanities, multilingual, modeling SUPERT: 0.7043 LLMEval: 87.5000

Topic Summary: Linked Open Data (LOD) in the Humanities

Brief Introduction into the Topic

Linked Open Data (LOD) is increasingly being integrated into the Humanities to enhance the accessibility, interoperability, and analysis of machine-readable information. Two notable initiatives, LODinG and DigItAnt, exemplify this trend. LODinG explores broad applications across various Humanities disciplines, while DigItAnt focuses on creating, linking, and exploiting LOD lexica for historical linguistics and epigraphy data.

Key Research Objects

LODinG: Covers digital humanities, computational linguistics, digital lexicography, philology, translatology, cultural and literary studies, media studies, information science, and law.
 DigItAnt: Manages LOD lexica and heterogeneous resources in historical linguistics and epigraphy, supporting the H2IOSC Italian infrastructures cluster project.

Key Research Methods

– LODinG: Collects, models, links, releases, and analyzes multilingual data in LOD format, developing a modular, cross-domain Humanities data model.

- DigItAnt: Designs architecture for LOD lexica and links Linguistic Linked Open Data (LLOD) with non-LLOD datasets, following FAIR principles.

Key Research Purpose

- LODinG: Improves extraction and modeling of multilingual LOD data, demonstrates LLOD applications in interdisciplinary contexts, and develops a cross-domain data model.

- DigItAnt: Aims to build and integrate LOD lexica for historical linguistics and epigraphy, support the H2IOSC project, and contribute to Italian Humanities infrastructure.

Example 1: Example topic description for topic ac1-53.

quantum-49 Dataset: Ouantum Model: CTM Method: SCS Category: Ambiguous Topic Words: foremost, boast, ineffective, craft, exacerbate, maintains, complication, pack, prevail, upon **SUPERT:** 0.4825 LLMEval: 62.5 Summary of Quantum Computing and Related Technologies Brief Introduction into the Topic Quantum computing and related technologies focus on advancing fault-tolerant quantum machines, enhancing quantum communication, and exploring novel materials and methods for quantum applications. Research spans from error correction in quantum computing to the development of high-performance quantum platforms, quantum networking, and advanced optical materials. Key Research Objects Fault-tolerant quantum machines: Large-scale quantum computers relying on error correction using logical and physical qubits. Quantum emitters and waveguides: Dipole quantum emitters embedded in dielectric slab waveguides. Quantum platforms: Diamond with engineered coherent defects for quantum control. Quantum simulations: SU(3) gauge theories in Quantum Chromodynamics (QCD) using discrete subgroup S(1080). Quantum networking: Quantum wrapper networking with classical headers and quantum payloads. Quantum box arrays: Self-assembled porous (metal-)organic networks modified by C60 adsorbates. Quantum optics: High refractive index polymers for infrared (IR) optics. Quantum communication: DWDM transmission using ultra-low loss G.654.D fiber and lumped doped fiber amplifiers. Ce³⁺ doped LiYF<inf>4</inf> crystals for ENDOR Quantum spin manipulations: spectroscopy. Quantum geometry: Desargues and Pappus properties in projective geometry applied to quantum systems. Key Research Methods Error correction: Spreading the value of a single logical qubit over multiple physical qubits. Waveguide analysis: Studying the effective quantum yield of dipole emitters in dielectric slab waveguides. Materials engineering: Controlling coherent defects in diamond for quantum applications. **Digitization schemes**: Practical digitization of SU(3) gauge theories using discrete subgroups. Noise characterization: Analyzing in-fiber scattering processes in quantum wrapper networking. Adsorbate modification: Using C60 to modify electronic states in quantum box arrays. Inverse vulcanization: Synthesizing high refractive index polymers from waste materials. DWDM transmission: Demonstrating high-speed data transmission using advanced fiber and amplifiers. ENDOR spectroscopy: Resolving Mims ENDOR spectra for electron-nuclear interactions. Geometric analysis: Applying Desargues and Pappus properties to quantum systems. Key Research Purpose To craft large, fault-tolerant quantum machines capable of robust computations. To **boast** the performance of quantum emitters and waveguides by optimizing their properties. To **maintain** long spin coherence times in diamond-based quantum platforms. To pack the infinite-dimensional gluon field into a minimal qubit representation for QCD simulations. To exacerbate the understanding of noise impacts in quantum networking. To **prevail** in engineering quantum states in on-surface network architectures. To foremost develop low-cost, high-performance materials for IR optics. To complicate the understanding of high-speed data transmission in optical communication. To **upon** investigate electron-nuclear interactions for spin manipulations.

To ineffective link projective geometry principles to quantum physics and information.

Example 2: Example SCS topic description for topic quantum-49.

acl-20

Dataset: ACL Model: BERTopic Method: CCS Category: Ambiguous Topic Words: word2vec, embeddings, vocabulary, wordnet, semantic, semantics, embedding, contextualized, word, linguistic SUPERT: 0.4274 LLMEval: 85.4250

Brief Introduction into the Topic

The extraction and analysis of health-related information from social media platforms, particularly **Twitter**, have become critical in **pharmacovigilance**, **public health monitoring**, and **mental health studies**. Advances in **Natural Language Processing (NLP)** and **transformer-based models** (e.g., BERT, ROBERTA, BERTweet) have improved the ability to process informal and noisy text, enabling the identification of **Adverse Drug Events (ADEs)**, **disease mentions**, and **psychological conditions** such as **depression**, **anxiety**, and **mood disorders**. Social media data, especially **tweets**, provide valuable insights into public health trends and individual health experiences.

Key Research Objects

- Psychological and Mood-Related Terms: Keywords such as depression, depressed, mood, psychological, and anxiety, relevant to ADE classification and diagnosis.

- Tweets: Informal text from Twitter containing colloquial language related to health, drug effects, and psychological conditions.
- Disease Mentions: Identification of diseases such as depression and anxiety in Spanish tweets and other social media content.

- Adverse Drug Events (ADEs): Mentions of adverse effects from medications and dietary supplements in tweets.

- Social Media Platforms: Twitter, Reddit, and health forums as primary sources of health-related data.

Key Research Methods

1. Transformer-Based Models

- BERT, RoBERTa, BERTweet, and CT-BERT are fine-tuned for tasks such as classification, Named Entity Recognition (NER), and information extraction.

 Ensemble Methods: Combining multiple BERT variants to improve accuracy in classifying health-related tweets.

2. Large Language Models (LLMs)

- **Data Augmentation**: Generating synthetic data to address data imbalance in ADE extraction and disease mention detection.

- Retrieval-Augmented Generation (RAG): Mapping informal tweet language to standardized terms like MedDRA Preferred Terms.

3. Advanced NLP Techniques

- Transfer Learning: Fine-tuning multilingual models like **mBERT** for disease mention extraction in Spanish tweets.

- Knowledge Graph Integration: Enhancing models with external evidence from medical gazetteers and ontologies.

- Flair-NER Framework: Evaluating embeddings for capturing complex disease mentions in social media text.

4. Preprocessing and Optimization

- Domain-Specific Preprocessing: Cleaning and preparing social media text for analysis.

- Hyperparameter Optimization: Enhancing model performance through tuning.

Key Research Purpose

- Enhance ADE Discovery: Improve the accuracy and robustness of ADE extraction and normalization from social media data.

- **Improve Disease Mention Extraction**: Develop advanced NER systems for identifying diseases in informal text, particularly in non-English languages like **Spanish**.

- Standardize Health-Related Terms: Map colloquial language to standardized medical terms (e.g., MedDRA Preferred Terms) for better pharmacovigilance and health monitoring.

- **Promote Open Research**: Provide open-source code and methodologies to foster collaboration and further innovation in the field.

- Monitor Public Health Trends: Leverage social media data for real-time monitoring of public health issues, including COVID-19 and psychological conditions.

Example 3: Example topic description for topic ac1-20.

quantum-49

Dataset: Quantum Model: CTM Method: CCS Category: Worst Topic Words: foremost, boast, ineffective, craft, exacerbate, maintains, complication, pack, prevail, upon SUPERT: 0.1685 LLMEval: 7.1150

Brief Introduction into the Topic

The topic revolves around the analysis of certain systems or strategies that are deemed **foremost** in their field but may **boast** features that are ultimately **ineffective**. These systems often **craft** solutions that inadvertently **exacerbate** existing issues, leading to further **complication**. Despite these challenges, the systems **maintain** their prominence, as they are designed to **pack** significant functionality. The research aims to explore why such systems **prevail** and the factors that drive their continued use **upon** further scrutiny.

Key Research Objects

The primary objects of research include systems or strategies that are considered leading in their domain. These objects are characterized by their advanced features, which are often highlighted as strengths but may contribute to inefficiencies or unintended consequences. The study focuses on understanding the dynamics of these systems and their impact on the broader context in which they operate.

Key Research Methods

The research employs a combination of qualitative and quantitative methods to evaluate the effectiveness of the systems under study. This includes analyzing case studies, conducting surveys, and performing comparative assessments to identify patterns and outcomes. The methods aim to uncover the reasons behind the systems' continued use despite their potential drawbacks.

Key Research Purpose

The purpose of the research is to critically examine the systems that are widely regarded as top-tier in their field. It seeks to identify the factors that contribute to their perceived success, as well as the unintended consequences that may arise from their implementation. The study aims to provide a comprehensive understanding of why these systems **prevail** and how they impact their respective domains.

Example 4: Example CCS topic description for topic quantum-49.