STRUCTSUM Generation for Faster Text Comprehension

Anonymous ACL submission

Abstract

We consider the task of generating structured 001 representations of text using large language models (LLMs). We focus on tables and mind 004 maps as representative modalities. Tables are more organized way of representing data, while 006 mind maps provide a visually dynamic and flexible approach, particularly suitable for sparse content. Despite the effectiveness of LLMs on different tasks, we show that current models struggle with generating structured outputs. In 011 response, we present effective prompting strategies for both of these tasks. We introduce a 012 taxonomy of problems around factuality, global and local structure, common to both modalities 014 and propose a set of critiques to tackle these issues resulting in an absolute improvement in accuracy of +37pp (79%) for mind maps and +15pp (78%) for tables. To evaluate semantic coverage of generated structured representations we propose AUTO-QA, and we verify the adequacy of AUTO-QA using SQuAD dataset. We further evaluate the usefulness of structured 023 representations via a text comprehension user study. The results show a significant reduction in comprehension time compared to text when using table (42.9%) and mind map (31.9%), without loss in accuracy.

1 Introduction

028

034

040

The overwhelming amount of information available online poses a significant challenge for users seeking to quickly grasp and process relevant information. Current large language models (LLMs), such as PALM-2 (PaLM2, 2023), Gemini (Gemini Team, 2023) and ChatGPT (OpenAI, 2022), while capable of providing text-based responses to user queries, often fail to adequately structure and organize this information in a way that facilitates comprehension (Tang et al., 2023). This can lead to information processing bottlenecks that hinder users' ability to efficiently extract meaningful insights from text.



Figure 1: Overview of (a) tables and (b) mind map generation prompts. Prompting steps are colored. Figure (a) show divide-and-generate prompt.. Input passages is first segmented into sub-passages, followed by multiple table generation. Figure (b) shows the generation process for mind maps. After the main concept generation, an iterative expansion phase takes place where the mind map is expanded until termination.

To address this issue, we introduce the notion of structured summaries, or STRUCTSUM in short. STRUCTSUMs are derived by hierarchically organizing information and inducing semantic connections from an input text passage. Without loss of generality, we focus on tables (Wu et al., 2022; Li et al., 2023) and mind maps (Buzan, 1996; Huang et al., 2021) as possible STRUCTSUM instantiations:

042

045

051

053

055

057

060

061

062

063

• **Tables** are well-studied in the NLP literature. However the vast majority of the work focused on simpler tasks where tables are inputs – such as QA (Herzig et al., 2020), semantic parsing (Bogin et al., 2019), NLG (Andrejczuk et al., 2022; Puduppully and Lapata, 2021; Laha et al., 2020), etc. – rather than outputs. Indeed, faithfully transforming an arbitrary text passage into a table is a difficult task as the model must deal with different challenges, such as reasoning at multiple levels, dealing with missing information, and visually consistent formatting. Motivated by the limitations above, we propose to generate multi-

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

115

ple tables instead. We argue that this is a simpler task for an LLM, as shown in Figure 2, which compares single-table and multi-table generation side by side. We therefore propose a divide-andgenerate prompting approach (see Figure 1) that first divides the input text into multiple text passages, each representing a sub-topic, followed by an LLM prompt to generate a table-caption pair for each smaller passage. This decomposition allows the model to generate smaller, focused and more informative tables, especially for complex text passages with multiple sub-topics.

065

066

080

081

084

086

090

092

096

101

102

103

104

105

106

107

Mind maps (Hu et al., 2021; Wei et al., 2019) ٠ are less studied in the literature, but are helpful for comprehension and learning (Buzan, 1996; Dhindsa et al., 2011). Mind maps are complementary to tables in their structure, allowing for more flexibility and dynamism than tables, as they are inherently schema-less. However, generating mind maps with LLMs presents several challenges: (i) the model first need to select a central concept, that is the fulcrum of all the successive extractions, as mind maps revolve around a central root node; (ii) being a schemaless abstraction, each connecting branch has its own independent sub-topic, making it difficult to automatically add branches all at once; (iii) to ensure readability and well-structuredness each leaf node should terminates the path in a way that concludes the idea or sub-topic; (iv) depending on the information density, some paths may be shorter than others. Therefore, the model should decide whether or not a branch is worth expanding. Following the structure of these observations, we propose an iterative prompting technique for mind map generation. As show in Figure 1, we initialize the mind map by generating the root concept. At each iteration, we decide either to expand the current mind map further or stop the process. During the expansion step, we prompt the model to add branches to the current leaf nodes. We represent the mind map as a JSON object, as it is easy to parse and verify.

108Through extensive experimentation with PALM-1092 (PaLM2, 2023), we show that LLMs are not al-110ways effective at generating STRUCTSUMs that are111factual and structurally correct. To overcome these112issues we propose a pipeline for structured data gen-113eration. Our pipeline consists of structure-specific114prompts followed by critics to assess output quality

along three different dimensions, that are common both to tables and mind maps: (i) *Factuality*, (ii) *Local Structure* and (iii) *Global Structure*. We found that our proposed critics improved the overall quality of the generated output by +37pp for mind maps and +15pp for tables.

To ensure the usefulness of STRUCTSUM for text-comprehension tasks, we propose Auto-QA as a measure of output coverage. We automatically generate QA pairs from input text and use structured outputs to answer these questions. Furthermore, we verify the appropriateness of using Auto-QA by comparing Auto-QA with human generated QA pairs on SQuAD (Rajpurkar et al., 2016) development set.

Finally, starting from the initial hypothesis that STRUCTSUMs can enhance the effectiveness of information-seeking scenarios, we conducted a user study to evaluate their impact on users' ability to process information, using a text comprehension user study. Results demonstrate how STRUCT-SUMs improve information seeking, specifically on timed text comprehension metrics. We found that by using the structured representation, users can answer questions 42.9% faster for tables and 31.9% for mind maps.

2 Related Work

Structured Output. Generating structured output from text has been explored in the context of information extraction (Li et al., 2023; Pietruszka et al., 2022). Most of the work focus on textto-table (Wu et al., 2022) generation using the model trained on domain specific dataset. Ni and Li (2023) use LLM for information extraction by generating key-value pairs. Tang et al. (2023) evaluate different models on table generation from text by prompting where table structure is provided as format instructions. Mind map generation has been explored in the form of relation graph structure (Hu et al., 2021; Wei et al., 2019) to summarize new articles (Cheng and Lapata, 2016; Hermann et al., 2015). In contrast, we focus on a generation pipeline applicable for multiple structured outputs types by prompting LLM given a text input. We keep the output structure flexible and domain independent by not instructing the model with specific format.

Prompting. Our prompting strategy is rooted in task decomposition techniques. Least-tomost (Zhou et al., 2023), in contrast with chain-ofInput Passage: The Mersey-class cruisers were improved versions of the Leander class with more armour and no sailing rig on a smaller displacement. Like their predecessors, they were intended to protect British shipping. The cruisers had a length between perpendiculars of 300 feet (91.4 m) ... Multiple Table Generation

Single Table Generation				The Mersev-cla	ass cruise	r				
The Mersey-class cruisers were improved versions of the Leander class with more rmour and no sailing rig on a smaller displacement.				Length	Beam	Dra	aught		Displacement	
Displacement	Length	Beam	Draught		300 feet	46 feet	20	feet 2	2 inches	4,050 long tons
4,050 long tons	300 feet	46 feet	20 feet 2 inches		Armament		0			
Speed	Range	Complement	_		weapon		Quantity	y L	ocation	
					BL 8-inch gun	ı	2	F	ore and aft on	pivot mounts
18 knots	8,750 nautical miles	300 to 350 officers and ratings			BL 6-inch gun	ıs	10	F	ive on each b	roadside in sponsons
Armament	_	_	_						•	
				Mersey	-class armour				The Mersey	-class cruiser's mach
Two breech-loading (BL) 8-inch (203 mm) guns, one each		Locati	on	Thick	aness (in)		Attribute	Value		
Armour	-	-	-	Lower	Armoured Dec	k 2 (fla	t) / 3 (sloj	pe)	Engine type	Two-cylinder
A lower armoure	A lower armoured deck that was 2 inches (51 mm) on the flat and 3 inches			Conni	ng Tower	9		_	Shafts	2

Figure 2: Example table generation for the text at top, comparing single table (left) vs multiple table generation (right). Some parts in the table and text were truncated (...) for readability. The full example is reported in Figure 6.

thought (Wei et al., 2022), progresses from easiest 165 166 to hardest questions eventually answering the complete question, while successive prompting (Dua 167 et al., 2022) iteratively generate new questions based on previous answers. Unlike least-to-most, 169 decomposed prompting (Khot et al., 2023) doesn't 170 restrict task decomposition from easiest to the hard-171 est and iteratively generate next steps that can be 172 executed by different systems. Most of the prior work is focused on reasoning for solving QA type 174 problems, in contrast, we are interested in trans-175 forming text to structured formats. Our divide-and-176 generate prompting for multiple table generation 177 (similar to least-to-most) uses an initial prompt 178 to divide the input passage into different topical 179 sub-passages that simplifies the table generation in 180 next step. Different from these tasks our iterative 181 prompting for mind maps requires reasoning over current structured output at each step. 183

Factuality. Attribution is used as a tool for assessing the reliability of LLMs and identifying potential sources of inaccuracy or fabrication in their generated outputs. Current work apply attribution on unstructured text generation settings, such as, question answering (Bohnet et al., 2022) and text generation tasks (Gao et al., 2023a). Diverging from that, our work require verifying the factuality of generated structured outputs.

Evaluation. Due to the cost of human evaluation, LLMs are used to critic the generated outputs (Wang et al., 2023). Recent instructions tuned models, such as, GPT-4 (OpenAI, 2023) and Chat-GPT (OpenAI, 2022) are shown to be strong evaluators. To avoid using external APIs, Kim et al. (2023); Wang et al. (2023) fine-tune a smaller pretrained model to critic model responses. We are interested in evaluating the quality structured outputs using critics and self-correct based on the feedback. As a part of data generation pipeline, our focus is on filtering instances that are incomplete and are not factually grounded. 200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

221

222

223

224

225

226

227

229

230

231

232

3 Generating STRUCTSUMs

We focus on tables and mind maps as a possible STRUCTSUM instantiations.

3.1 Tables: Divide & Generate prompting

Given an input text we would like to transform it into multiple tables. Although generating multiple tables from text may seem unnecessary, singletable generation lead to several issues, as shown in Figure 2 (bottom left). The model often produces complex table structures, resulting in missing cell values or the exclusion of relevant information. Additionally, complex tables are difficult to verify for factual accuracy and can require additional mental effort from the user to understand.

To address these limitations, we propose a divide-and-generate approach that dynamically partitions the passage into smaller subtopic segments. While deterministic rule-based chunking methods (e.g., based on word or sentence count) can be employed, they often produce suboptimal results due to potential under-chunking, over-chunking, and the absence of division for certain instances. Therefore, the chunking must be adaptive and depend on the input text and its sub-topic distribution. We use a one-shot prompt for this step, as shown in Appendix C (Figure 14). After the chunking, we prompt the model to generate a table along with its caption for each sub-passage obtained in the



Figure 3: Example mind map output. The full example along with the input text is reported in Figure 7.

236

240

241

243

244

245

247

248

249

253

previous step.

Alg	gorithm 1 Mind maps Iterative Prompting
Rec	luire:
	input text passage: input
	maximum number of steps: max_steps
1:	step $\leftarrow 0$
2:	$mindmap \leftarrow GENERATE-ROOT(input)$
3:	while step < max_steps do
4:	$step \leftarrow step + 1$
5:	if CONTINUE-PROMPT(input, mindmap) then
6:	expansions $\leftarrow EXPAND(input, mindmap)$
7:	mindmap \leftarrow JSON-CRITIC(expansions)
8:	else
9:	return mindmap
10:	end if
11:	end while
12:	return mindmap

3.2 Mind maps: Iterative Prompting

Contrary to tables, mind maps are more flexible and present a different set of challenges. The first challenge is representation. We desire a representation that is (i) close to a familiar format, and (ii) is easily parsable and verifiable using current tools. JSON meets both of these requirements. The second challenge is that mind maps, unlike tables where each row can be produced linearly, necessitate attaching information in different locations depending on which branch is being expanded. This requires the model to think radially.

We propose an iterative prompting for mind maps generation. Algorithm 1 shows the overall procedure. Details of each prompt is in Appendix C. We start by generating the root concept that becomes the central node for the mind map. This separate step allows the model to independently reason about the theme of the passage. After generating the root, at each step we prompt the model to decide if current mind map can be expanded further. If the model decides to expand (line 5), we prompt the model using the current mind map to add more branches. Otherwise, the procedure terminates and we return the current mind map. At each expansion step we sample multiple mind maps. Utilizing the fact that JSON verification is cheaper we select the topmost JSON that is parsed correctly. In the rare case, when none of the samples are parsable we call a critic prompt to correct the top JSON (line 7).

256

257

258

259

260

261

262

263

265

266

267

268

269

270

271

272

273

274

275

276

278

279

281

282

283

285

286

287

289

290

291

292

293

294

297

299

300

301

302

303

4 Data Generation Pipeline

We now present our STRUCTSUM data generation pipeline. Although each STRUCTSUM is seemingly different, we identify three dimensions that are common to both table and mind map modalities: (i) Factuality, (ii) Local Structure, and (iii) Global Structure. We use a set of critics, implemented via prompts, to ensure sufficient quality across each dimension. Through our initial experiments we find that tweaking each critic according to the structure is more helpful.

4.1 Factuality Critic

We use post-attribution (Gao et al., 2023a) to verify factuality, as we found that jointly generating and attributing (Gao et al., 2023b) results in (i) unnatural text output and (ii) in the model copying verbatim from the input text passage.

Critic cost is one aspect that requires consideration. For example, for tables, verifying each cell could be more robust, however, it increases the number of LLMs calls (listed in Table 1), from $\mathcal{O}(1)$ to $\mathcal{O}(\#$ number of cells).

For simplicity, we choose a single prompt per STRUCTSUM: for tables we ask the LLM to attribute each row, while for mind maps we ask to attribute each path from root to leaf. We convert the input text passage to a list of sentences and ask the model to cite, following the [x,y] format for attribution, the source sentence(s) where the information can be found or [NA] in case this is not possible. The prompts are reported in Figure 15.

4.2 Local Structure Critic

For tables, a common issue arises from the model misplacing values in incorrect columns. For example, placing "66 years" in the *Birth date* column or an address in the *Company Name* column. To detect such errors, we leverage each column header as a category and verify whether all cell values within

Critic	# LLM calls			
Chuc	Tables	Mindmaps		
Factuality	$\mathcal{O}(1)$	$\mathcal{O}(1)$		
Local Structure	$\mathcal{O}(\#cols)$	$\mathcal{O}(\text{#paths})$		
Global Structure	NA	$\mathcal{O}(1)$		

Table 1: Cost for each critic in terms of #LLM calls as proxy. #cols is number of columns in output table. #paths is the number of paths in a mind map from root node to a terminal node.

that column belong to the same category. For mind maps, we observed that a well-defined terminal node can often represent the entire path leading to it. We use this fact and prompt the model to verify whether the terminal node is a specific value, rather than a general concept. The prompts are reported in Figure 16.

4.3 Global Structure Critic

304

305

306

307

311

330

331

332

335

336

Global critic allows us to verify the overall structure of the output. This means understanding
whether all the information contained in a STRUCTSUM makes sense globally.

For tables, we simply verify whether the table is well formatted: e.g. we verify equal number columns in the header and subsequent rows, therefore ignoring semantic content of the table and only focusing on form rather than the content. This is realized via simple heuristics.

For mind maps, we used a stricter approach, to ensure that information were semantically valid on a global level. Specifically, we convert the mind map into a familiar format like table of contents (ToC), which we hypothesize is more likely to be seen during the pre-training phase of existing LLMs, and ask the model to check if the ToC is at right level of abstraction. The prompts are reported in Figure 17.

5 Semantic Coverage using AUTO-QA

In this section we propose an automatic way to assess the quality and the general usefulness of STRUCTSUMs introducing AUTO-QA coverage as proxy metric. This metric measures the semantic coverage or percentage of questions that are answerable when using a STRUCTSUM s, instead of the full text passage t. Formally it is defined as:

$$COV(s) = \frac{1}{|\text{GenQA}(t)|} \sum_{i=1}^{|\text{GenQA}(t)|} \mathbb{1}_{E_{a_i}} \left[Q(s, q_i) \right]$$

where GenQA(x) is a function that generates (q, a) pairs given the input text passage $t, Q(s, q_i)$ is a function that generates an answer given in input a STRUCTSUM s and the question q_i , whereas the indicator function $\mathbb{1}_{E_{a_i}}(x)$ asses the answer equivalence between a_i and x. Figure 18 in Appendix C, show all the prompts associated with AUTOQA module (Deutsch et al., 2021; Fabbri et al., 2022). 338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

354

355

356

357

359

360

361

362

363

364

365

366

367

368

369

370

371

373

374

375

376

377

378

379

Independently of perceived quality, it is worth noting that this simple metric can be thought as an abstractiveness measure or compression quality for a given STRUCTSUM s. A value of 1 indicates no information loss at the expense of no compression/abstraction, whereas a value of 0 indicates theoretically maximum compression at the expense of not providing any useful information. A target value is therefore application specific and must be adjusted accordingly ¹.

QA pairs generation GenQA(t) is implemented by prompting the LLM to generate a list of question-answer (QA) pairs conditioned on the input text t. To ensure that the quality of QA pairs is sufficient, after generation, we we apply a threestep procedure. First, we removed duplicate questions via string match. Second, we removed answers if none of the words appeared in the input text, thereby ensuring with reasonable certainty that the answer is grounded in the text without being overly stringent. Third, we performed a cyclic consistency check, where we prompted the model to answer the generated question based on input text.

Question answering We use a simple prompt for function $Q(s, q_i)$. For tables, we convert the table representation to a markdown table format, whereas for mind maps we simply serialize the information as a JSON object.

Answer Equivalence As the model might generate verbose answers, verifying whether two answers are the same is a problem of semantic similarity. Instead of using lexical matching, that is $\mathbb{1}_{E_{a_i}}(x) := a_i = x$, we prompt the model to check if two answers are equivalent.

¹It is possible to include coverage as a critic. But we opted not to do so, as the threshold for coverage depends on the specific use case. This also allowed us to analyze coverage independently, without being influenced by other factors.

```
is cleaned up by page filtering to remove disam-
biguation pages, redirect pages, deleted pages, and
```

7.1 Filtering for Tables Generation

Not all input paragraphs are well suited for table generation. As a proxy for selecting adequate passages, we used regex-based filters to only include passages with more that 20 numeric values and removed passages with less than three sentences. In a real world setting, we would like a systematic way of deciding which modality is adequate for a given text. We leave this exploration as future work.

For all the experiments, we use the Unicorn (PaLM-2-unicorn, 2023) variant of PALM-2, a fine-tuned transformer-based model with UL2 (Tay et al., 2022) like objectives. PALM-2 improves on

PaLM (Chowdhery et al., 2023) through optimized

scaling, richer training data and instructing tun-

To test our pipeline on a diverse set of input passages, we selected Wikipedia text as the source. Specifically, we started with the English section of WIKI40B (Guo et al., 2020) dataset. The dataset

non-entity pages. We then iterated over the dataset

to extract passages that are inputs to our prompt.

ing (Wei et al., 2021; Chung et al., 2022).

8 Results

6

7

381

386

387

397

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

Model

Dataset

In this section, we present the results of our experiments using PALM-2.

8.1 Quality impact of prompting style and automated critics

We assessed the quality of generated structured data through manual human ratings. The study was conducted on 100 instances for mind maps. For multi-table generation, we choose 100 individual table-text pair for annotation ². Input passages were obtained via data filtering strategy described in Section 7.

Guidelines Annotators were asked to rate each instance as "Good" or "Bad" by checking the overall quality of the output. For both modalities, annotators were asked to check for factuality as well

Tables			Mindma	aps
Single Table	54		СоТ	39
Multi Table	63		Iterative	42

Table 2: Table / mind map accuracy per prompt style. Outputting multiple tables provides higher quality for the table modality. For mindmaps, an iterative approach is to be preferred to a CoT approach. Full prompts are reported in the Appendix C.

Critic	Tables	Mind maps
Baseline [†]	63	42
\hookrightarrow Structure	70	71
\hookrightarrow Factuality	78	79

Table 3: Human annotation accuracy at different pipeline stages. The use of critics is a critical step to improve perceived quality. Local and Global Structure critic provides a significant lift for mind maps. The increase in performance for Factuality, is similar for both Tables and mind map.

as the structural quality of the output. To help the annotators measure the structural quality we asked the annotators to check "table structure", "table header", "column header-value match" for tables. For mind maps, they were asked to check "incomplete branches", "not a good main concept", "too dense / too sparse" and "wrong edge connections". We also encouraged the annotators to mark the instance as bad if they find any other issues.

423

424

425

426

497

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

Prompt style Table 2 show the results for both the modalities. For table generation task we find that annotators prefer multiple tables generation outputs compared to single table generation. For mind map, we compare chain-of-thought (Wei et al., 2022) with our proposed iterative generation strategy described in Algorithm 1. We find that iterative generation were preferred over simpler prompt outputs.

8.2 Do Critics Align with Human Ratings?

Through our human annotations results in Section 8.1, we find that many generated outputs are not of acceptable quality. To improve the quality of the generated data and to avoid costly human annotations, we propose to use a combination of critics as a measure of data quality. To verify the efficacy of our critics, we first filtered the generated dataset with our critics. Specifically, we performed a logical AND of individual critics and filtered the

 $^{^{2}}$ We made sure that the input passages are the same for the different ablations within modalities. For multi-table generation, we choose 100 text-table pairs generated using 52 input passages.



Figure 4: AUTO-QA based coverage. A point $\langle X, Y \rangle$ in each line show that X% of data has at least Y% of coverage measured using AUTO-QA.

instances that do not pass the criterion. We then conducted the same evaluation of Section 8.1.

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

Results in Table 3 show that using the proposed critics the overall quality is improved by a significant margin. We observe that data filtered using *Structure* (Global and Local) and *Factuality* critics improve the percentage of acceptable instances generated using the pipeline. We find that the quality of mind maps improve by absolute +37pp. Similarly, for tables quality improves by absolute +15pp. These results indicate that the critics were able to retain good examples and that the selection criterion is in agreement with human judgement.

8.3 Measuring Coverage via Auto-QA

Results in Figure 4 show AUTO-QA coverage for mind maps and tables. The curve shows for a particular coverage threshold what percentage of data meets that threshold. Overall, we observe that tables have better coverage compared to mind maps, meaning that they have an higher abstractiveness or information retention capacity. Interestingly, even though both modalities are perceptually different, we notice that both of them follow similar trends.

8.4 Is Auto-QA a reasonable metric?

We investigate the feasibility of using AUTO-QA 475 as a surrogate for manually written QA pairs. We 476 aim to determine whether AUTO-QA can generate 477 QA pairs of comparable quality to those written 478 by humans, and leading to a similar evaluation of 479 semantic coverage. To verify the same, we use ran-480 domly selected 1000 < passage, question, answer> 481 triples from the SQuAD (Rajpurkar et al., 2016) 482 validation set (common for all the experiments). 483

	QA Type		
	Auto	Human	
Mind map	55.6	61.4	
Multi-Table (Divide-and-generate)	66.8	69.3	
Single Table	57.1	58.8	
Query Focused (Single Table)	81	85.5	

Table 4: QA accuracy on different modalities as context, generated using SQuAD validation set. AUTO-QA is automatic question-answer pair generation. Human QA are original SQuAD questions curated by humans.

Using the text *passage* as input we generate different STRUCTSUMS. Next, we generate a QA pair corresponding to each text *passage*. This QA pairs acts as a substitute for human written QA pair for AUTO-QA study. The goal is to check whether, keeping the passage and output STRUCTSUM the same, there is a correlation in performance between human generated QA pairs and automatically generated QA pairs. 484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

Table 4 shows the overall results. Second (Mind maps) and third (Multi-Table) row show the comparison between Human QA and Auto QA for our proposed divide-and-generate prompt for tables and iterative prompt for Mindmap generation. We can see that AUTO-QA has comparable results and is a reasonable substitute for human generated questions as a measure of semantic coverage.

8.5 Multiple Tables vs Single Table

To check whether generating multiple tables is better at covering more information, we perform a comparison between the ability to answer question by generating single or multiple tables. On comparing Multi-Table and Single Table row in Table 4, we observe that for both AUTO-QA and Human QA generating multiple table provides more coverage. So in addition to the benefits such as comparatively better verifiability and robust generation, multiple table generation are also better at covering more semantic information.

8.6 Query Focused Generation

In many cases user intent is known in advance,514for example, a user query to search or LLM-based515Assistant interface (e.g., ChatGPT, Gemini, etc.).516We explore the possibility of generating structured517data in the presence of a query. We perform a518*preliminary* analysis by adding the query in single519table generation prompt. As we can see in last row520



Figure 5: Results for timed text comprehension based user study. Plots show 95% confidence interval over time taken in seconds to answer question with different structure combinations as context. For both tables (left) and mind map (right), compared to text only, we observe significant reduction (42.9% and 31.9% resp.) in average time taken by annotators to answer the question.

in Table 4, query focused generation improve the performance by more than 20 for AUTO-QA and 25 points for Human-QA. Since this need further investigation in terms of prompting and more importantly detailed analysis of output quality. We leave a comprehensive exploration for query focused structured data generation as future work.

8.7 Are STRUCTSUMs useful?

We evaluate whether STRUCTSUM are useful abstractions for the users. For this we design a *timed* text comprehension based user study. We assume that user is looking to answer a specific query, i.e. has a specific intent. We measure time taken to satisfy the user intent as a proxy of usefulness.

We create an intent in the form of a question along with different context combinations. For example, for a question q, we create $\langle q, s \rangle$, $\langle q, t \rangle$, $\langle q, s + t \rangle$ as possible combinations, where s is a STRUCTSUM and t is the input text passage. Each of these combinations are presented to different annotators while ensuring that no annotator see the same question twice. We then measure how long it takes to answer the question in each scenario.

STRUCTSUMs for the study were generated using our data generation pipeline and critic-based filtering, as discussed in Section 4. In total 600 instances were annotated, equally divided into different context combinations for mind maps and table generation. Annotators consistently answered correctly across all context combinations (Appendix B), suggesting that the level of context did not significantly impact their accuracy. Figure 5 shows the overall results. The plots show 95% confidence interval of time taken by the annotators when using different modality:

552

553

554

556

557

558

559

560

561

562

563

564

566

567

568

569

570

571

572

573

574

575

576

578

579

580

582

- **Tables.** Figure 5a shows that on average annotators with access to tables were able to answer almost 42.9% time faster on average compared to annotators with only text. Furthermore, we observe that presenting both table and text is also useful to the annotators.
- Mind maps. Figure 5b shows the results for the study with mind map. A similar trend can be observed, with a reduction of approximately 31.9% in average time between annotators with mind maps compared to annotators that only used text to answer the question.

We note that $\langle q, s + t \rangle$ performs worse than $\langle q, s \rangle$, we believe this is due to the fact that the annotators cross-checked the answer from both the modalities. Leading to increase in time to answer the question.

9 Conclusion

In this work we study the potential of structured representations like tables and mind maps to enhance information comprehension. Utilizing our divideand-generate prompting and iterative expansion, we achieved significant improvements in output quality (+37pp for mind maps, +15pp for tables) using structure-specific prompts and critics. We proposed AUTO-QA based coverage metric that automatically generates QA pairs from the input text and uses STRUCTSUM outputs to answer them.

551

521

523

686

687

689

690

691

10 Limitations

583

We outline the limitations of our work to ensure transparency and inspire future research. First, the 585 structured output representations we experimented 586 with are limited to tables and mind maps. How-587 ever, to comprehensively evaluate the effectiveness of our critics and pipeline, it is desirable to also 589 evaluate other input and output modalities, e.g. im-590 age and video, considering the recent advances in 591 VLMs. Our structsum generation is performed us-592 ing a LLM, however, primary aim of this study is validating faster text comprehension. Furthermore, 594 prior work have shown a reasonable portability of prompts across similar models (Zhou et al., 2023; Khot et al., 2023). Secondly, our work and experimental findings are limited to only English sources. 598 We plan to also explore multilingual structured summaries in future works. Third, we would to 600 warn against the risk of blindly trusting models to generate structured summaries from an input accurately. Although we take extra care to increase the factuality of the outputs via the use of critics, and experimentally validate QA coverage, we believe 605 that special care should be taken to verify outputs in accuracy-sensitive applications.

Despite these limitations, our work serves as an
initial step in constructing reliable structured summarization evaluations, models and applications.
We hope future research can greatly benefit from
this starting point.

References

614

615

616

618

619

620

621

623

624

625

627

628

629

633

- Ewa Andrejczuk, Julian Eisenschlos, Francesco Piccinno, Syrine Krichene, and Yasemin Altun. 2022.
 Table-to-text generation and pre-training with TabT5.
 In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 6758–6766, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
 - Ben Bogin, Matt Gardner, and Jonathan Berant. 2019. Global reasoning over database structures for textto-SQL parsing. 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 3659–3664, Hong Kong, China. Association for Computational Linguistics.
- Bernd Bohnet, Vinh Q. Tran, Pat Verga, Roee Aharoni, Daniel Andor, Livio Baldini Soares, Massimiliano Ciaramita, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, Kai Hui, Tom Kwiatkowski, Ji Ma, Jianmo Ni, Lierni Sestorain Saralegui, Tal Schuster, William W. Cohen, Michael Collins, Dipanjan

Das, Donald Metzler, Slav Petrov, and Kellie Webster. 2022. Attributed question answering: Evaluation and modeling for attributed large language models.

- Tony Buzan. 1996. The Mind Map Book: How to Use Radiant Thinking to Maximize Your Brain's Untapped Potential. Plume.
- 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.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Daniel Deutsch, Tania Bedrax-Weiss, and Dan Roth. 2021. Towards Question-Answering as an Automatic Metric for Evaluating the Content Quality of a Summary. *Transactions of the Association for Computational Linguistics*, 9:774–789.
- Harkirat S. Dhindsa, Makarimi-Kasim, and O. Roger Anderson. 2011. Constructivist-visual mind map teaching approach and the quality of students' cognitive structures. *Journal of Science Education and Technology*, 20(2):186–200.
- Dheeru Dua, Shivanshu Gupta, Sameer Singh, and Matt Gardner. 2022. Successive prompting for decomposing complex questions. In *Proceedings of the* 2022 Conference on Empirical Methods in Natural Language Processing, pages 1251–1265, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Alexander Fabbri, Chien-Sheng Wu, Wenhao Liu, and Caiming Xiong. 2022. QAFactEval: Improved QAbased factual consistency evaluation for summarization. In Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 2587–2601, Seattle, United States. Association for Computational Linguistics.
- Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and Kelvin Guu. 2023a. RARR: Researching and revising what language models say, using language models. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16477–16508, Toronto, Canada. Association for Computational Linguistics.

Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023b. 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.

695

700

703

704

705

706

707

710

711

712

713

715

716

717

718

719

721

723

725

726

727

728

729

730

731

733

734

740

741

742

743

744

745

746

747

- Gemini Team. 2023. Gemini: A family of highly capable multimodal models. Technical report, Google.
- Mandy Guo, Zihang Dai, Denny Vrandecic, and Rami Al-Rfou. 2020. Wiki-40b: Multilingual language model dataset. In *LREC 2020*.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. *Advances in neural information processing systems*, 28.
- Jonathan Herzig, Pawel Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Eisenschlos. 2020. TaPas: Weakly supervised table parsing via pre-training. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4320–4333, Online. Association for Computational Linguistics.
- Mengting Hu, Honglei Guo, Shiwan Zhao, Hang Gao, and Zhong Su. 2021. Efficient mind-map generation via sequence-to-graph and reinforced graph refinement. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8130–8141, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Luyang Huang, Shuyang Cao, Nikolaus Parulian, Heng Ji, and Lu Wang. 2021. Efficient attentions for long document summarization. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1419–1436, Online. Association for Computational Linguistics.
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2023. Decomposed prompting: A modular approach for solving complex tasks. In *The Eleventh International Conference on Learning Representations*.
- Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoo Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2023. Prometheus: Inducing finegrained evaluation capability in language models.
- Anirban Laha, Parag Jain, Abhijit Mishra, and Karthik Sankaranarayanan. 2020. Scalable Micro-planned Generation of Discourse from Structured Data. Computational Linguistics, 45(4):737–763.
- Tong Li, Zhihao Wang, Liangying Shao, Xuling Zheng, Xiaoli Wang, and Jinsong Su. 2023. A sequenceto-sequence&set model for text-to-table generation.

In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5358–5370, Toronto, Canada. Association for Computational Linguistics. 748

749

750

751

752

753

754

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

799

800

- Xuanfan Ni and Piji Li. 2023. Unified text structuralization with instruction-tuned language models. *arXiv preprint arXiv:2303.14956*.
- OpenAI. 2022. Chatgpt: Optimizing language models for dialogue.
- OpenAI. 2023. Gpt-4 technical report.
- PaLM-2-unicorn. 2023. PaLM-2 google ai blog. https://blog.google/technology/ai/ google-palm-2-ai-large-language-model/. Accessed: 2023-05-10.
- PaLM2. 2023. PaLM2 technical report. https://ai.google/static/documents/ palm2techreport.pdf. Accessed: 2023-05-10.
- Michał Pietruszka, Michał Turski, Łukasz Borchmann, Tomasz Dwojak, Gabriela Pałka, Karolina Szyndler, Dawid Jurkiewicz, and Łukasz Garncarek. 2022. Stable: Table generation framework for encoder-decoder models. *arXiv preprint arXiv:2206.04045*.
- Ratish Puduppully and Mirella Lapata. 2021. Datato-text Generation with Macro Planning. *Transactions of the Association for Computational Linguistics*, 9:510–527.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Xiangru Tang, Yiming Zong, Jason Phang, Yilun Zhao, Wangchunshu Zhou, Arman Cohan, and Mark Gerstein. 2023. Struc-bench: Are large language models really good at generating complex structured data?
- Yi Tay, Mostafa Dehghani, Vinh Q Tran, Xavier Garcia, Jason Wei, Xuezhi Wang, Hyung Won Chung, Dara Bahri, Tal Schuster, Steven Zheng, et al. 2022. Ul2: Unifying language learning paradigms. In *The Eleventh International Conference on Learning Representations*.
- Tianlu Wang, Ping Yu, Xiaoqing Ellen Tan, Sean O'Brien, Ramakanth Pasunuru, Jane Dwivedi-Yu, Olga Golovneva, Luke Zettlemoyer, Maryam Fazel-Zarandi, and Asli Celikyilmaz. 2023. Shepherd: A critic for language model generation. *arXiv preprint arXiv:2308.04592*.
- Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. 2021. Finetuned language models are zero-shot learners. *arXiv preprint arXiv:2109.01652*.

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.

806 807

809

810

811 812

813

814

815

816

817 818

819

821

823

824 825

- Yang Wei, Honglei Guo, Jinmao Wei, and Zhong Su. 2019. Revealing semantic structures of texts: Multigrained framework for automatic mind-map generation. In Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19, pages 5247–5254. International Joint Conferences on Artificial Intelligence Organization.
- Xueqing Wu, Jiacheng Zhang, and Hang Li. 2022. Textto-table: A new way of information extraction. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2518–2533, Dublin, Ireland. Association for Computational Linguistics.
 - Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. 2023. Least-to-most prompting enables complex reasoning in large language models. In *The Eleventh International Conference on Learning Representations.*

829 830 831

832 833

834

835

837

838

841

843

847

850

851

A Data Generation statistics

Table 5 shows different statistics of data generated using our prompts. For tables generation we observe that our methods generate almost two (~ 1.9) tables per instance and the tables have 7.1 rows and 3.3 columns on average. Mind maps have an average of 11.8 nodes with a depth of 2.2. We show example mind map and table generation in Figure 7 and Figure 6 respectively.

Tables				
Avg #words per chunk	114.8			
Avg #sentences per chunk	3.9			
Avg #words per input	240.6			
Avg #sentences per input	8.1			
Avg #rows	7.1			
Avg #cols	3.3			
Avg #tables	1.9			
Max #tables	11			
Mind map				
Avg #words	194.6			
Avg #sentences	7.9			
Avg #nodes	11.8			
Avg depth	2.2			

Table 5: Table / mind map text input and output statistics. On average two (~ 1.9) tables (top) are generated per input text instance. Mind maps (bottom) contains 11.8 nodes on average.

B User Study

Usually, mind maps are represented as a graph as shown in Figure 7. However, for the text comprehension user study described in Section 8.7, to avoid bias due to color or orientation, we simplify the representation as a tree (Figure 10). To establish the known query intent, annotators' are first shown with input question, e.g., Figure 8. Next, on clicking Show content button, annotators are shown context in the form of either text (Figure 9), structure (Figure 10), or structure + text (Figure 11). The question-answer pairs were generated automatically conditioned on input text (Section 5). Annotators were also allowed to mark an instance un-answerable. The user study for tables is performed in a similar manner. We annotated 100 question-answer pairs for both mind maps and tables. Each input instance is annotated with three different context combinations, leading to 600 total

Tables				
Table	95.6			
Text	94.1			
Table+Text	94.1			
Mindmaps				
Mind map	97.7			
Mind map Text	97.7 94.3			

Table 6: Answer accuracy (as percentage) for different context combinations. Structure context performs on par/better compared to text.

annotations. We filtered instances that were marked un-answerable by the annotators (32% and 22% for tables and mind map study resp.). To avoid penalizing for spelling errors or other typing mistakes, the answers were evaluated via human evaluation. Table 6 shows the overall accuracy as percentage of questions answered correctly in different context. Irrespective of context combinations, annotators were able to answer the questions correctly with a high accuracy.

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

C Prompts

In this Section we report the different prompts used in this study. In our implementation we use Jinja (https://jinja.palletsprojects.com/) to specify the prompt template. The Mersey-class cruisers were improved versions of the Leander class with more armour and no sailing rig on a smaller displacement. Like their predecessors, they were intended to protect British shipping. The cruisers had a length between perpendiculars of 300 feet (91.4 m), a beam of 46 feet (14.0 m) and a draught of 20 feet 2 inches (6.1 m). They displaced 4,050 long tons (4,110 t). The ships were powered by a pair of two-cylinder horizontal, direct-acting, compound-expansion steam engines, each driving one shaft, which were designed to produce a total of 6,000 indicated horsepower (4,500 kW) and a maximum speed of 18 knots (33 km/h; 21 mph) using steam provided by a dozen cylindrical boilers with forced draught. The Mersey class carried enough coal to give them a range of 8,750 nautical miles (16,200 km; 10,070 mi) at a speed of 10 knots (19 km/h; 12 mph). The ships' complement was 300 to 350 officers and ratings. Their main armament consisted of two breech-loading (BL) 8-inch (203 mm) guns, one each fore and aft on pivot mounts. Their secondary armament was ten BL 6-inch (152 mm) guns, five on each broadside in sponsons. Protection against torpedo boats was provided by three quick-firing (QF) 6-pounder Hotchkiss guns and three QF 3-pounder Hotchkiss guns. The ship was also armed with a pair of submerged 14-inch (356 mm) torpedo tubes and carried a pair of 14-inch torpedo carriages. The Mersey-class ships were protected by a lower armoured deck that was 2 inches (51 mm) on the flat and 3 inches (76 mm) on the slope. It sloped down at the bow to reinforce the ram. The armoured sides of the conning tower were 9 inches (229 mm) thick.

	(a)	Input	text for	table	genetation
--	-----	-------	----------	-------	------------

Caption: The Mersey-class cruisers were improved versions of the Leander class with more armour and no sailing rig on a smaller displacement.

Displacement	Length	Beam	Draught	
4,050 long tons (4,110 t)	300 feet (91.4m)	46 feet (14.0 m)	20 feet 2 inches (6.1 m)	
Speed	Range	Complement	-	
18 knots (33km/h; 21 mph)	8,750 nautical miles (16,200 km; 10,070 mi) at a speed of 10 knots (19 km/h; 12 mph)	300 to 350 officers and ratings		
Armament	-	-	-	
Two breech-loading (BL) 8-inch (203 mm) guns, one each fore and aft on pivot mounts. Their secondary armament was ten BL 6-inch (152 mm) guns, five on each broadside in sponsons. Protection against torpedo boats was provided by three quick-firing (QF) 6-pounder Hotchkiss guns and three QF 3-pounder Hotchkiss guns. The ship was also armed with a pair of submerged 14-inch (356 mm) torpedo tubes and carried a pair of 14-inch torpedo carriages.				
Armour – A				
A lower armoured deck that was 2 inches (51 mm) on the flat and 3 inches (76 mm) on the slope. It sloped down at the bow to reinforce the ram. The armoured sides of the conning tower were 9 inches (229 mm) thick.				

(b) Single table generation output.

Contion: The M

Lower Armoured Deck

Conning Tower

Caption: The Mersey-class cruiser

Length Beam I		Draught		Displacement	
300 feet (91.4m) 46 feet (14.0 m)		20 feet 2 in	ches (6.1 m)	4,050 long tons (4,110 t)	
Caption: Armame	nt				
Weapon		Quantity	Location		
BL 8-inch (203)	nm) gun	2	Fore and aft on pivot mounts		
BL 6-inch (152	mm) guns	10	Five on eac	ch broadside in sponsons	
(QF) 6-pounder guns	r Hotchkiss	3	-		
QF 3-pounder	Hotchkiss guns	3	-		
14-inch (356 m	m) torpedo tube	2	Submergeo	t	
14-inch torpedo	o carriages	2	-		

cuption. The worsey class cruiser's machinery			
Attribute	Value		
Engine type	Two-cylinder horizontal,		
Shafts	2		
Horsepower	6,000 horsepower (4,500 kW)		
Speed	18 knots (33 km/h; 21 mph)		
Biolers	Twelve cylindrical boilers with forced draught		
Range	8,750 nautical miles (16,200 km; 10,070 mi) at a speed 10 knots (19 km/h; 12 mph)		
Complement	300 to 350 officers and ratings		
Caption: Mersey-cla	ss armour		
Location	Thickness (in)		

2 (flat) / 3 (slope)

9

(c) Multiple table generation output.

Figure 6: Example outputs for single and multiple table generation approach. Text in (a) shows the input. (b) and (c) show the outputs for single and multiple table generation respectively.

Kathleen "Kay" Daly (January 8, 1919 – October 16, 1975) was an Irish-born American advertising executive and one of the four "celebrated Daly sisters". At Norman, Craig & Kümmel she was the creative force behind the famous Maidenform "I Dreamed" campaign and Revlon's legendary 1952 Fire And Ice campaign, working with photographer Richard Avedon. She also was responsible for the line "Every woman alive loves Chanel Number Five". She went on to join Revlon in 1961 as vice president and creative director. Kathleen Daly was born in Castlecaufield, County Tyrone, Ulster, Ireland, in 1919. Northern Ireland was created two years later with Tyrone one of its six counties. The family emigrated early in the 1920s. She grew up as one of four sisters, Maggie, Kay, Maureen, and American-born Sheila. They became known for their writing and work in journalism, fashion, and advertising, and were called "the celebrated Daly sisters" by Time magazine in 1966. Life magazine ran a feature story on them in 1949 and a follow-up in 1959. All four were at least once employed by the Chicago Tribune. When she moved to San Francisco after World War II, Kay Daly famously rented space on a billboard to advertise for an apartment. It not only netted her an apartment, but netted her nationwide fame and countless marriage proposals. She had a brief marriage to BMW executive and film producer Richard Bradford (part of the famous Bradford family of Plymouth Colony), who fathered her sons John (Kelly), Richard, and Peter. She then was married to journalist and executive Warren Leslie, who adopted and raised her sons, until her death on October 16, 1975, of pancreatic cancer. She was survived by husband Warren, sons Kelly, Peter, and Richard Bradford, and stepsons Warren and Michael Leslie.



(a) Input text for mind map generation.

(b) Mind map output.

Figure 7: Example mind map (below) generation for the input text (above). We use mermaid.js (https://mermaid.js.org/) to visualize the output.

Question: What was Kathleen Daly's birth place?



Figure 8: Example UI frame that is shown at the beginning of each annotation instance.

Question: What was Kathleen Daly's birth place?

Passage

Text: Kathleen "Kay" Daly (January 8, 1919 – October 16, 1975) was an Irish-born American advertising executive and one of the four "celebrated Daly sisters". At Norman, Craig & Kümmel she was the creative force behind the famous Maidenform "I Dreamed ..." campaign and Revlon's legendary 1952 Fire And Ice campaign, working with photographer Richard Avedon. She also was responsible for the line "Every woman alive loves Chanel Number Five". She went on to join Revlon in 1961 as vice president and creative director. Kathleen Daly was born in Castlecaufield, County Tyrone, Ulster, Ireland, in 1919. Northern Ireland was created two years later with Tyrone one of its six counties. The family emigrated early in the 1920s. She grew up as one of four sisters, Maggie, Kay, Maureen, and American-born Sheila. They became known for their writing and work in journalism, fashion, and advertising, and were called "the celebrated Daly sisters" by Time magazine in 1966. Life magazine ran a feature story on them in 1949 and a follow-up in 1959. All four were at least once employed by the Chicago Tribune. When she moved to San Francisco after World War II, Kay Daly famously rented space on a billboard to advertise for an apartment. It not only netted her an apartment, but netted her nationwide fame and countless marriage proposals. She had a brief marriage to BMW executive and film producer Richard Bradford (part of the famous Bradford family of Plymouth Colony), who fathered her sons John (Kelly), Richard, and Peter. She then was married to journalist and executive Warren Leslie, who adopted and raised her sons, until her death on October 16, 1975, of pancreatic cancer. She was survived by husband Warren, sons Kelly, Peter, and Richard Bradford, and stepsons Warren and Michael Leslie.

Unanswerable

Enter your answer:



Figure 9: A followup frame shown after Figure 8 with text as context.

Question: What was Kathleen Daly's birth place? Caption: Kathleen "kay" daly, an irish-born american advertising executive who was one of the four "celebrated daly sisters".

⊢Kay Daly

- Persona	l life
- Bi	rth place
	- Castlecaufield
	- County Tyrone
	- Ulster
	L Ireland
- Bi	rth date
	∟ January 8, 1919
- D	eath place
	L United States
	eath date
	L October 16, 1975
- Spouse	
	- Richard Bradford
	L Warren Leslie
LC	hildren
	– John <mark>(Kelly)</mark> Bradford
	- Richard Bradford
	L Peter Bradford
∟ Career	

- Employer - Norman, Craig & Kümmel L Revion - Occupation L Advertising executive L Notable works - Maidenform 'I Dreamed' campaign

L Revion's Fire And Ice campaign

Unanswerable Enter your answer:

Figure 10: A followup frame shown after Figure 8 with structure (mind map) output as context.

Question: What was Kathleen Daly's birth place? Caption: Kathleen "kay" daly, an irish-born american advertising executive who was one of the four "celebrated daly sisters".
L Kay Daly
- Personal life
- Birth place
+ Castlecaufield
+ County Tyrone
- Ulster
L Ireland
Birth date
L January 8, 1919
- Death place
L United States
- Death date
L October 16, 1975
+ Spouse
+ Richard Bradford
L Warren Leslie
L Children
- John (Kelly) Bradford
- Richard Bradford
L Peter Bradford
L Career
- Employer
+ Norman, Craig & Kümmel
L Revion
- Occupation
L Advertising executive
L Notable works
- Maidenform 'I Dreamed' campaign
L Revion's Fire And Ice campaign
Passage
Kathleen "Kay" Daly (January 8, 1919 – October 16, 1975) was an Irish-born American advertising executive and one of the four "celebrated Daly sisters". At Norman, Craig & Kümmel she was the creative force behind the famous Maidenform "I Dreamed" campaign and Revlor's legendary 1952 Fire And Ice campaign, working with photographer Richard Avedon. She also was responsible for the line "Every woman alive loves Chanel Number Five". She went on to join Revloin 1961 as vice precident and creative director. Kathleen Daly was born in Castlecautifield, County Tyrone, Ulster, Iteland, 11919. Northem Ireland was created two years later with Tyrone one of its six counties. The family emigrated early in the 1920s. She grew up as one of four sisters, Maggie, Kay, Maureen, and American-born Shela. They became known for their writing and work in journalism, fashion, and advertising, and were called "the celebrated Daly sisters" by Time magazine in 1966. Life magazine ran a faeture story on them in 1949 and a follow-up in 1959. All four were at least once employed by the Chicago Thibume. When she moved to San Francisco anter World Var II, Kay Dia faradord (part of the famous Bradford family of Pymouth Colony), who fathered her sons John (Kelly), Richard, and Peter. She then was marriade proposals. She had a brief marriage to BMW executive and film producer Richard Bradford (part of the famous Bradford family of Pymouth Colony), who fathered her sons. John (Kelly), Richard, and Peter. She then was married to journalist and executive Warren Leslie, who adopted and raised her sons, until her death on October 16, 1975, of pancreatic cancer. She was survived by husband Warren, sons Kelly, Peter, and Richard Bradford, and stepsons Warren and Michael Leslie.

Unanswerable	
Enter your answer:	Submit

Figure 11: A followup frame shown after Figure 8 with structure (mind map) output and input text as context.

```
Coal. Solutions to global Clima

individual efforts and internat

Choose primary concept that is the root

Output:

MindMap

{{ format_json(step1)|safe }}

END_THOUGHT

Add branches:

MindMap

{{ format_json(step2)|safe }}

END_THOUGHT

Add branches:

MindMap

{{ format_json(step2)|safe }}

END_THOUGHT

Add branches:

MindMap

{{ format_json(step3)|safe }}

END_THOUGHT

Can we add branches ?

Output: Yes

END_THOUGHT

Can we add branches ?

Output: Yes

END_THOUGHT

Add branches:

MindMap

{{ format_json(step4)|safe }}

END_THOUGHT

Can we add branches ?

Output: No

END_THOUGHT

New for the text below:

END_THOUGHT

New for the text below:

END_THOUGHT

END_THO
       END_THOUGHT
Now for the text below:
{{input_text}}
Now for the text below:
{{input_text}}
Now for the text below:
{{input_text}}
Now for the text below:
{{root}}
END_THOUGHT
{% if current_mindmap -%}
Can we add branches?
Output: Yes
Add branches:
MindMap
{{current_mindmap}}
END_THOUGHT
{%- endif %}
Can we add branches ?
{*- if y_n_current %}
Output: {{y_n_current %}
Output: {{y_n_current }}
END_THOUGHT
Add branches:
MindMap
{*- else %}
Output: {*- endif %}
{*- endif %}
```

Figure 12: Iterative prompt in Jinja template format for mind map generation that is used in Algorithm 1.

<pre>{% set step</pre>	<pre>= '{"node": "Global Climate change", "branches ("node": "effects", "branches": [{"node": " ting ice"}{ffects", "branches": [{"node": " ting ice"}{frects", "branches": [{"node": " causes", "branches": [{"node": "causes", "branches": cllution", "branches": [{"node": "Carbon emission ("node": "Burning coal"}]]]}, {"node": "solutions branches": [{"node": "Individual efforts"}, {" e": "International resolutions"]]]] - "} s a diagram used to visually organize ormation into a hierarchy, showing relationships ong pieces of the whole. It is often created und a single concept. Major ideas are connected ectly to the central concept, and other ideas nch out from those major ideas. Mind maps can be inerated based on the content present in text in tiple steps. following example. llowing text: te change has many effects, including melting , heat waves, and droughts. It is caused by the anced greenhouse effect, which is caused by the anced greenhouse of global climate change include widual efforts end international concentione.</pre>
Thought: Print \hookrightarrow clint \hookrightarrow solution \hookrightarrow effective branches \hookrightarrow branches \hookrightarrow response to the clint \hookrightarrow the clint \hookrightarrow branches \hookrightarrow branches \hookrightarrow response to the clint \hookrightarrow clint \hookrightarrow the clint \hookrightarrow the clint \hookrightarrow clint \hookrightarrow clint	mary concept is Global climate change. Global mate change has branches, effects, causes and utions. Effects have branches that include ects, melting ice and heat waves. Causes have nches enhanced greenhouse effect. Solutions have anches, individual efforts and international olutions.
Output: MindMap {{ format_js	on(step) safe }}
Now summariz {{input_text	e the following text as a mind map. }}

Figure 13: Prompt in Jinja template format for mind map generation without iterative process.

Now divide the following passage into Smaller passages → grouped by similar facts. {{input_text}} Summarize the contents of the text below in a table. {{input text}} Use the following format. Caption: A caption for the table you generate. Can be → multiple lines Table: A table in markdown format. Caption:

Figure 14: Text segmentation prompt (top) for multiple table generation. Zero-shot prompt for text to table and caption generation (bottom).

Figure 15: Factuality critic prompts for Table (top) and Mind maps (bottom).

Your task is to check if all the values in a list falls → under a category. Go over all the values one by one → and check if they belong to the assigned category. → Use the following format to answer. Thought: Reasoning for the answer. Answer: Single final answer yes or no. Category: {{category}}

Category: {{category}} Values: {{values}} Thought:

```
There are some words or sentences that describes concept

→ while other describes values associated with them.

→ Values are defined as ordinals, type of job, degree

→ , education level, location, region, date etc.

Answer No If any words is not a specific value otherwise

For example:

Words:

Delhi

Cat

26 May

Lawyer

Thought: All the words are specific content words.

Answer: yes

Words:

IBM

Trucks

Birth

Family

26 May

Thought: Many words such as Trucks, Family, Birth are

→ concept without specific values.

Answer: no

Words:

{(words)}

Thought:
```

Figure 16: Local structure critic prompt for Tables (top, zero-shot) and Mind maps (bottom, few-shot).



Figure 17: Global structure critic few shot prompt for Mind map.

Your task is to generate a list of fact based questions that → can be answered by the text passage. The format → should be [Question][Answer]. Paragraph: {{text}}

Check if the following two answers are equivalent. Use the following format. Question: question text Answer 1: answer text Answer 2: answer text Conclusion: Yes/No Question: {{context_question}} Answer 1: {{answer1}} Answer 2: {{answer2}} Conclusion:

Answer in concise manner the question using the information → below. Say <unknown> when the questions cannot be → answered. {{data}}

Question: {{question}}

Figure 18: Prompts used by the AutoQA pipeline: QA pair generation prompt (*top*); Conditional answer equivalence (*middle*); Question answering prompt (*bottom*)