

Structuring the Unstructured: A Systematic Review of Text-to-Structure Generation with a Universal Evaluation Framework

Anonymous ACL submission

Abstract

The evolution of AI systems toward agentic operation and context-aware retrieval necessitates transforming unstructured text into structured formats like tables, knowledge graphs, and charts. While such conversions enable critical applications from summarization to data mining, current research lacks a comprehensive synthesis of methodologies, datasets, and metrics. This systematic review examines text-to-structure techniques and the encountered challenges, evaluates current datasets and assessment criteria, and outlines potential directions for future research. We also introduce a universal evaluation framework for structured outputs, establishing text-to-structure as foundational infrastructure for next-generation AI systems.

1 Introduction

The rapid growth of the agentic AI systems is redefining the paradigm of information processing, where autonomous agents must dynamically acquire, synthesize, and act upon structured knowledge extracted from textual sources (Singh et al., 2025; Buehler, 2025). This agentic revolution creates dual demands for structured knowledge representation: (1) enabling dynamic knowledge grounding during multi-step agentic reasoning, and (2) serving as curated retrieval sources for Retrieval-Augmented Generation (RAG) pipelines (Zhang et al., 2024c, 2025; Singh et al., 2025). Complex structures, such as tables and graphs, play a crucial role in conveying information, as they can intuitively display data relationships and temporal trends (Baek et al., 2024; Huang et al., 2024b; Ghaifarollahi and Buehler, 2024; Li et al., 2024a), while preserving hierarchical dependencies, which are vital for enhancing RAG reliability (Edge et al., 2024; Zhuang et al., 2024; Li et al., 2024c; Wang et al., 2024b). These structures enable both comprehension of complex information and machine-friendly

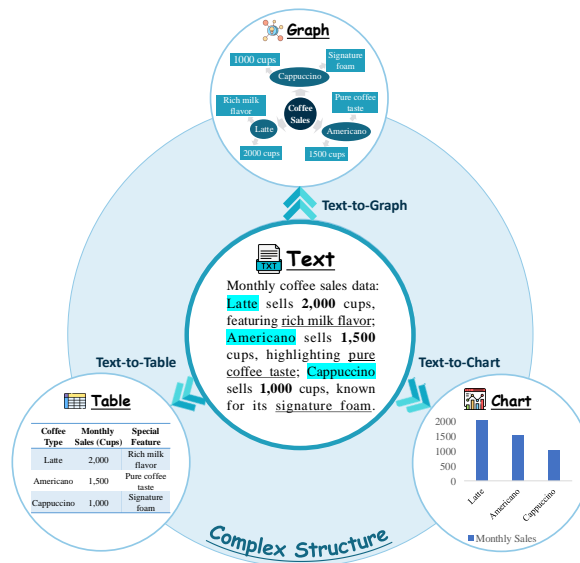


Figure 1: An example of text-to-complex structure conversion. The center represents the text, while the surrounding elements constitute the complex structures. Based on this text, three types of complex structures – table, graph, and chart – can be generated respectively. representations for downstream processing (Jain et al., 2024; Reif et al., 2024).

Traditional NLP methods have demonstrated proficient performance in extracting simple entities and relationships (Sang and Meulder, 2003; Yao et al., 2019), but their efficacy is limited when confronted with complex structured knowledge implicit in the text. The emergence of Large Language Models (LLMs) enhances the capabilities of language models in capturing complex semantics and long-range dependencies, crucial for supporting agentic systems’ dynamic knowledge requirements and RAG’s structured retrieval needs. This advancement drives the growth of research in complex structured information extraction. Currently, the most common forms of structured output include tables, knowledge graphs, mind maps, and charts, which are crucial in downstream tasks such as RAG-enhanced question-answering and agentic data analysis (Caciularu et al., 2024; Xu et al.,

2025). However, this task also poses considerable challenges: complex structured information is represented in diverse forms in the text, and accurately extracting this information requires a deep semantic understanding of the text, including identifying implicit relationships and reasoning. Moreover, there is a lack of task-specific data tailored for training and evaluation. To address these issues, researchers have explored various methods to enhance model performance and constructed several high-quality datasets for benchmarking.

Notwithstanding the significance of this task, the field currently lacks a comprehensive survey to summarize the existing research progress. This gap becomes more critical with the rise of agentic systems that demand structured knowledge representations (Zhang et al., 2024c). Specifically, two key challenges persist: (1) the baseline models and benchmark datasets used for comparison vary significantly across studies; (2) current evaluation systems remain constrained by traditional metrics that systematically misalign with human judgment when assessing structured outputs (as demonstrated in Figure 2). To bridge these gaps, our contributions are threefold:

- We conduct a systematic synthesis and critical analysis of current methods (§3), datasets (§4), and evaluation metrics (§5).
- We establish a novel Text-to-Structure (T2S) benchmark with a universal framework for structured output evaluation (§6).
- Extensive experimental results show that our framework significantly outperforms traditional metrics, providing actionable resources.

2 Task Formulation

As shown in Figure 1, we categorize Text-to-Structure generation into three common types: tables, graphs, and charts. For any task, the input should consist of a textual passage with N tokens, denoted as $\mathcal{T} = [t_1, t_2, \dots, t_N]$, and an instruction text with M tokens, denoted as $\mathcal{I} = [i_1, i_2, \dots, i_M]$ as constituent elements. We then provide formal definitions and detailed explanations for each category to demonstrate their distinct characteristics.

2.1 Table Generation

We define the table T as: $T = (\mathcal{R}, \mathcal{C}, \mathcal{E}, Capt)$, where $\mathcal{R} = \{r_1, r_2, \dots, r_m\}$ is an ordered set of m rows, $\mathcal{C} = \{c_1, c_2, \dots, c_n\}$ is an ordered set of n columns, $\mathcal{E} : \mathcal{R} \times \mathcal{C} \rightarrow \mathcal{D}$ is a function that

Ground Truth	Content Error	Structural Error
Quality		
Human Evaluation		
BERTScore		
ROUGE-L		

Figure 2: Misalignment between traditional metrics and true quality in structured output evaluation.

assigns a data entry to each cell in the table, and $Capt$ denotes the caption of the table. It is allowed to generate multiple tables at the same time (Wu et al., 2022; Tang et al., 2024; Jain et al., 2024).

2.2 Graph Generation

A graph can be considered an extension of a table due to its greater flexibility and schema-less nature. To model more complex structures, we extend the definition of graphs to include attributes and semantics: $G = (\mathcal{V}, \mathcal{E}, \mathcal{A}, \mathcal{L})$, where \mathcal{V} is a set of nodes representing entities or concepts, $\mathcal{E} \in \mathcal{V} \times \mathcal{V}$ is a set of edges representing relations between nodes, $\mathcal{A} : \mathcal{V} \cup \mathcal{E} \rightarrow \mathcal{P}$ is an attribute function assigning properties to nodes and edges, $\mathcal{L} : \mathcal{V} \cup \mathcal{E} \rightarrow \Sigma^*$ is a semantic function mapping nodes and edges to contextual meanings. The framework unifies graph-based structures, such as knowledge graphs and mind maps (Wei et al., 2019; Hu et al., 2021). See Appendix A for more details.

2.3 Chart Generation

A chart can be considered a visualization that goes beyond the structure of a graph. The input may include data-specific contexts, such as tables or data attributes (e.g., data types, ranges, etc.) (Tian et al., 2023). The goal is to generate visualization charts that feature the appropriate specifications, including the type of the chart, the visual encodings, and other relevant details (Zhang et al., 2024b).

3 Methods

3.1 Fine Tuning

Supervised Fine Tuning (SFT) is still the default approach to text-to-structure generation since it only necessitates paired $text \rightarrow target$ instances. Although intuitive, its performance differs across tasks because of variances in data needs, structural intricacy, and computational expenses. In text-to-table generation, sequence-to-sequence models

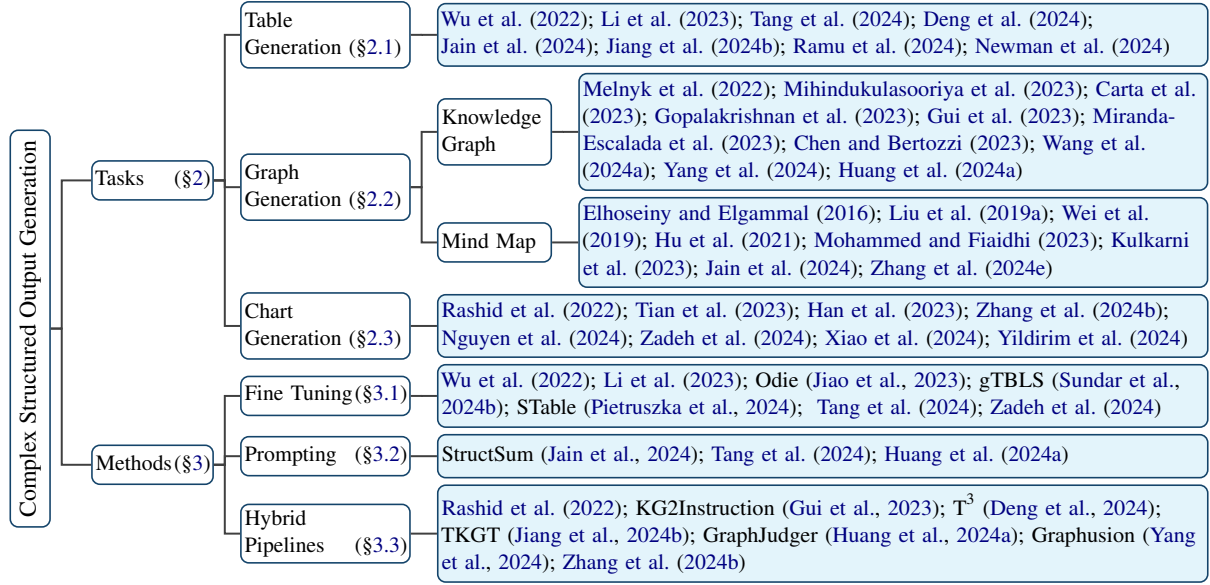


Figure 3: Taxonomy of representative works in text-to-structured generation categorized by tasks and methods.

with header-aware serialization (Wu et al., 2022; Li et al., 2023) produce validated tables but need enormous annotated datasets, which restricts low-resource applicability. QA-style cell queries (Sundar et al., 2024a,b) decrease data requirements, while being challenged by intricate layouts. Instruction tuning with Chain-of-Thought prompting (Wei et al., 2022) and LLM-synthesized pairs increases robustness but includes computational expense (Jiao et al., 2023; Tang et al., 2024). For text-to-KG generation, SFT tends to follow a two-step pipeline, entity recognition followed by relation extraction (Gui et al., 2023; Wang et al., 2024a; Yang et al., 2024), which enhances structure but has a risk of error propagation. In text-to-chart tasks, instruction tuning with RL-based feedback makes it more adaptable, although performance relies on the quality of the feedback, particularly for uncommon chart types (Xia et al., 2024; Zadeh et al., 2024).

3.2 Prompting

The prosperity of LLMs (Dubey et al., 2024; OpenAI, 2024; DeepSeek-AI et al., 2025) has triggered the evolution of diverse prompting methods for text-to-structure tasks (Wei et al., 2022; Dua et al., 2022; Khot et al., 2023), with each presenting trade-offs in accuracy, efficiency, and simplicity. Basic prompts such as “generate tables/graphs/charts” are quick but have poor results (Jain et al., 2024; Huang et al., 2024a). Incorporating schema constraints enhances structural fidelity at the expense of potentially reducing flexibility in tasks with implicit relations (Tang et al., 2024; Deng et al., 2024; Jain et al., 2024). Chain-of-Thought prompts are

also used to enhance reasoning and accuracy (Wang et al., 2023b; Tang et al., 2024). Decomposed prompting improves accuracy by breaking up input (Khot et al., 2023; Jain et al., 2024), but can lose global context, hurting coherence in applications such as cross-document structure construction.

3.3 Hybrid Pipeline

3.3.1 Table Generation

Harnessing table’s natural triple structure (row, column, cell) (Deng et al., 2024; Sundar et al., 2024b), it is common to let LLMs perform either tuple extraction from text (Wang et al., 2023a), or decomposition of text into atomic propositions (Hoyle et al., 2023), which subsequently serves as input for three approaches: *Code Aggregation* (Rozière et al., 2023; Deng et al., 2024) exploits LLM’s coding versatility for efficient data integration but falters on complicated structures and scalability. *KG with RAG* builds knowledge graphs for noise-robust RAG-based generation at the cost of flexibility and scalability (Jiang et al., 2024b). *Schema-Based Inference* models both qualitative nuances and quantitative facts despite similar scalability limitations (Ahuja et al., 2025).

3.3.2 Graph Generation

KG construction begins with Named Entity Recognition (NER), Relation Extraction (RE), and disambiguation (He and Choi, 2021; Gui et al., 2023), but is susceptible to noise and hallucinations. Therefore, structure-based denoising aids reliability (Xie et al., 2018; Deng et al., 2023), and LLM-guided verification boosts accuracy (Liu et al., 2023a;

Dataset	Size	Domain	T2S	Annot.	Sch.	Text Unit			Complexity		GTS
						Sent.	Para.	Doc.	Reas.	Struct.	
Text-Table Datasets											
WikiBio (Lebret et al., 2016)	728K	Open-Domain	X	X	∞	✓			★	✱	X
E2E (Novikova et al., 2017)	50.6K	Restaurant	X	✓	□		✓		★	✱	✓
Rotowire (Wiseman et al., 2017)	4.9K	Sports	X	✓	□			✓	★★	✱✱	✓
WikiTableText (Bao et al., 2018)	5.0K	Open-Domain	X	X	∞		✓		★	✱	✓
MLB (Padupully et al., 2019)	26.3K	Sports	X	✓	□		✓		★★	✱✱	✓
WikiTablePara (Laha et al., 2019)	171	Open-Domain	X	✓	∞		✓		★	✱✱	✓
WikiTableT (Chen et al., 2021)	1.5M	Open-Domain	X	✓	∞		✓		★	✱✱	X
InstructIE (Jiao et al., 2023)	14.7K	Open-Domain	✓	✓	∞		✓		★★	✱✱	✓
arxivDIGESTables (Newman et al., 2024)	2.2K	arXiv	✓	✓	∞			✓	★★★	✱✱	✓
CPL (Jiang et al., 2024b)	850	Law [‡]	✓	✓	∞			✓	★★	✱✱	✓
CT-Eval (Shi et al., 2024)	86.6K	Open-Domain [‡]	✓	✓	∞		✓		★	✱✱	✓
DescToTTo (Ramu et al., 2024)	1.3K	Open-Domain	✓	✓	∞		✓		★	✱✱✱	✓
LiveSum (Deng et al., 2024)	3.8K	Sports	✓	✓	□			✓	★★★	✱	✓
Struc-Bench Table (Tang et al., 2024)	4.1K	Sports	X	✓	□			✓	★★	✱✱	✓
Text-Graph Datasets											
WebNLG (Gardent et al., 2017)	25.3K	Open-Domain	X	✓	□		✓		★	✱	✓
SCIERC (Luan et al., 2018)	500	AI Papers	X	✓	□		✓		★★	✱✱	✓
AGENDA (Koncel-Kedziorski et al., 2019)	40K	AI Papers	X	X	□		✓		★★	✱✱	X
GenWiki (Jin et al., 2020)	1.3M	Open-Domain	X	X	□		✓		★	✱	✓
DART (Nan et al., 2021)	82.2K	Open-Domain	X	✓	□		✓		★★	✱	✓
EventNarrative (Colas et al., 2021)	224K	EventKG	X	X	□		✓		★★★	✱✱✱	X
KeLM (Agarwal et al., 2021)	18M	Open-Domain	X	X	□	✓			★★	✱	X
REBEL (Cabot and Navigli, 2021)	879K	Open-Domain	X	X	□	✓			★★	✱	X
Text2MindMap (Hu et al., 2021)	44.6K	News	✓	✓	□			✓	★★	✱✱	✓
InstructIE (Gui et al., 2023)	362K	Open-Domain [‡]	✓	X	□		✓		★★	✱✱	✓
MINE (Mo et al., 2025)	100	Open-Domain	✓	X	∞			✓	★★★	✱✱✱	X
Text-Chart Datasets											
AutoChart (Zhu et al., 2021)	10.2K	Data Analysis	X	✓	□		✓		★★	✱	X
Pew (Kantharaj et al., 2022)	9.3K	Society & Policy	X	✓	□		✓		★★★	✱✱✱	✓
Statista (Kantharaj et al., 2022)	34.8K	Market & Stats	X	✓	□		✓		★★	✱✱	✓
Text2Chart (Rashid et al., 2022)	717	Data Analysis	✓	✓	□		✓		★★	✱✱	✓
ChartX (Xia et al., 2024)	48K	Open-Domain	X	✓	□		✓		★★★	✱✱✱	✓
Text2Chart31 (Zadeh et al., 2024)	11.1K	Data Analysis	✓	X	□		✓		★★★	✱✱✱	✓

Table 1: Comparison and classification of existing text-to-structure generation benchmarks across key dimensions. **Size** denotes the number of (text, data) pairs in the dataset, including the training set. **Domain** shows the text source domain. **T2S** (Text-to-Structure) specifies whether the dataset is specifically designed for text-to-structure tasks. **Annot.** (Annotation) indicates human annotation involvement. **Sch.** (Schema) shows schema limitedness (‘□’ for limited, ‘∞’ for unlimited). **Text Unit** shows the granularity level of text (sentence / paragraph / document). **Reas.** (Reasoning) indicates the complexity of reasoning required for text-to-structure conversion, and **Struct.** (Structure) shows the structural complexity of the data, where a higher number of ‘★’ or ‘*’ indicates greater complexity in the corresponding dimension. **GTS** (Gold Test Set) denotes manual verification of the test set. ‘‡’ indicates the presence of Chinese text in the dataset. See Appendix B.1 for more details.

Huang et al., 2024a). Clustering-based node merging alleviates sparsity but threatens semantic oversimplification (Mo et al., 2025). Mindmaps (Jain et al., 2024), on the other hand, take a hierarchical approach: rooted iterative prompting guarantees structural clearness but at the expense of the dense interlinking that is characteristic of KGs.

3.3.3 Charts Generation

Text-to-chart generation generally employs a two-phase pipeline: *Data point identification*, which has inherent trade-offs: table-based extraction is structured but inflexible (Rashid et al., 2022), while text-based approaches are more adaptable but may be inconsistent (Zhang et al., 2024b); *Vision encoding*, which suffers compounded constraints: unclear data hinders chart-type prediction, while visualization synthesis amplifies earlier mistakes (Rashid et al., 2022; Zhang et al., 2024b; Zadeh et al., 2024). Both phases necessitate more resilient approaches.

4 Datasets

Table 1 provides a comprehensive list of all currently available benchmarks and contrasts them along eight axes. Most structure-to-text datasets can be adapted for text-to-structure tasks (Wiseman et al., 2017; Obeid and Hoque, 2020; Zhu et al., 2021; Kantharaj et al., 2022; Zhao et al., 2023; Lin et al., 2024), although information loss sometimes hinders suitability (Nie et al., 2018; Chen et al., 2020; Parikh et al., 2020), prompting dedicated text-to-structure dataset creation. Datasets are categorized by human supervision, schema constraints, text length, and exclusive human-annotated test sets. Complexity is quantified via structural complexity of outputs and reasoning complexity in generation (see Appendix B for further details). Analysis reveals significant gaps: *Text-Table* datasets dominate, yet few demand high complexity; *Text-Graph* datasets generally lack crucial human supervision; and *Text-Chart* datasets are scarce. Future

work should prioritize annotated, high-complexity datasets, especially for underrepresented graph and chart domains, to ensure balanced progress.

5 Evaluation Metrics

Evaluating complex structured outputs poses multifaceted challenges distinct from traditional IE tasks. While the latter often benefit from well-defined formats and objective standards (Xu et al., 2024), the inherent diversity and complexity of structured outputs necessitate a more comprehensive evaluation framework. As shown in Figure 4, after generating *structured output* from the *original text*, prevailing methods bifurcate into: (1) **Direct Evaluation** jointly assesses the *ground truth*, *original text*, and *structured output* (§5.1, §5.2, §5.3); (2) **Indirect Evaluation** first generates intermediate *content* (e.g., propositions, question-answer pairs) from the *original text*, then uses this *content* and the *structured output* to produce further *content* (e.g., answers), and finally compares the two *generated content* sets (§5.3). Through analysis of this established framework, we identify limitations that motivate our novel evaluation paradigm.

5.1 Human-based Evaluation

Human-based evaluation remains the gold standard, assessing content coverage, structural validity, and factual consistency (Jiao et al., 2023; Tang et al., 2024; Jain et al., 2024). It can also be evaluated by manual pairwise comparison (Zadeh et al., 2024), processing pipeline stages (Rashid et al., 2022; Deng et al., 2024), or indirectly via comprehension tasks (Jain et al., 2024). While human evaluation provides accurate results, it is costly in terms of time and resources, highlighting the complexity of quality that automated metrics must capture.

5.2 Rule-based Evaluation

Metrics like Exact Match, ROUGE-L (Lin, 2004), Levenshtein Distance (Haldar and Mukhopadhyay, 2011), and chrF (Popovic, 2015) rely on predefined rules focusing on surface-level features: token overlap, sequence edits, or n-gram alignment (Post, 2018; Nalawati and Dini Yuntari, 2021; Wu et al., 2022; Tang et al., 2024). Although ROUGE-L shows the strongest correlation with human evaluation among these metrics (Li et al., 2023; Wang et al., 2023c), it still fails to capture semantic and structural aspects, leading to misalignment in structured output assessment despite its convenience.

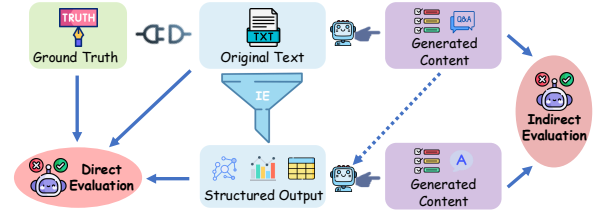


Figure 4: Comparison of direct evaluation and indirect evaluation methods for structured outputs.

5.3 Model-based Evaluation

Direct Scoring Learning-based methods use pre-trained models like SentenceBERT (Reimers and Gurevych, 2019), BERTScore (Zhang et al., 2020), BLEURT (Sellam et al., 2020), and BARTScore (Yuan et al., 2021) to evaluate text by semantic and contextual similarity, demonstrating better alignment with human judgements than token-matching rules on semantics (Jiao et al., 2023; Tang et al., 2024; Deng et al., 2024). However, these models remain inadequate for assessing structural correctness, as they lack explicit schema understanding and relational reasoning, and are insensitive to fine-grained content variations critical for quality assessment. On the other hand, direct prompting of LLMs partially addresses this gap by evaluating content-structure similarity via Chain-of-Thought prompting (Chiang and Lee, 2023; Liu et al., 2023b; Tang et al., 2024) although its reliability depends on prompt design and model biases.

Indirect Scoring Complex structured outputs often lack unique ground truth due to valid variations, like row/column order, equivalent values, or complete absence (Guo et al., 2020; Li et al., 2023; Tang et al., 2024; Jain et al., 2024), necessitating automated quality assessment methods: NLI-based alignment provides superior robustness to structural diversity than rigid rule-matching (Liu et al., 2019b; Ramu et al., 2024); QA-based evaluation (Jain et al., 2024; Deng et al., 2024) better assesses how well the structured outputs retain original information. Although these indirect evaluation approaches are more robust and generalizable, they demand stronger reasoning capabilities from models, such as the techniques from TableQA (Zhang et al., 2024d; Wang et al., 2024c), KGQA (Jiang et al., 2023b; Luo et al., 2024; Jin et al., 2024), and chart reasoning (Masry et al., 2023; Akhtar et al., 2023; Masry et al., 2024; Zhang et al., 2024a) to avoid error propagation while assessing outputs.

6 Text-to-Structure (T2S) Benchmark

To address the limitations of existing evaluation frameworks summarized in Section 5, we propose a novel framework capable of accurately and comprehensively evaluating structured outputs. Qualitatively, we address the previous vague criteria by implementing concrete metrics: an F1-style semantic metric using LLM-as-judge for **content** evaluation, overcoming limitations of string-based metrics like ROUGE and BERTScore (Tang et al., 2024), and defining two new orthogonal dimensions for **structure**, *alignment* for format consistency and *clarity* for unambiguous representations, inspired by header evaluation of prior works (Wu et al., 2022; Jiao et al., 2023) and Jain et al. (2024). To validate our framework’s superiority, we introduce a novel Text-to-Structure (T2S) benchmark, detailing: criteria (§6.1), datasets and models (§6.2), results (§6.3), and validation of effectiveness (§6.4).

6.1 Scoring Criteria

We evaluate structured outputs across two dimensions: *content faithfulness* and *structure coherence*.

Content Faithfulness Measures the accuracy and completeness of the generated information, ensuring it reflects the source text without errors or omissions. It uses an F1-style metric combining: (1) **precision**, which identifies conflicts between output and reference data (e.g., factual errors or semantic mismatches) and (2) **recall**, which measures how thoroughly the output captures key details from the source. These are combined into a single score via harmonic mean:

$$\text{FAITHFULNESS} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

Structure Coherence Evaluates logical organization and interoperability of the output. It examines (1) **alignment** with expected formats (e.g., schema consistency, data types) and (2) **clarity** of presentation (e.g., unambiguous representation, redundancy-free content). Again, we employ the harmonic mean for these two orthogonal metrics, as it penalizes imbalance and more accurately reflects joint effectiveness:

$$\text{COHERENCE} = 2 \cdot \frac{\text{alignment} \cdot \text{clarity}}{\text{alignment} + \text{clarity}}$$

We employ LLM-as-judge with detailed scoring rubrics (0-100 scale) for all four sub-metrics: *precision*, *recall*, *alignment*, and *clarity*. Detailed explanations and implementation are in Appendix C.1.

6.2 Datasets and Models

We carefully select six datasets from the openly available sources listed in Table 1, aiming to encompass a diverse range of complexities, cover multiple languages, and prioritize open-domain datasets. Specifically, three datasets for text-to-table generation: InstructIE, Struc-Bench, and LiveSum; two for text-to-graph generation: DART and InstructIE; and Text2Chart31 for text-to-chart generation. We select several mainstream LLMs for evaluation, which can be divided into two groups: non-thinking models (e.g., Llama, GPT) and thinking models (e.g., R1, O4-mini). For the former group, we also evaluate multiple models under SFT. Other models are evaluated in a zero-shot setting. For non-thinking models, we perform inference at temperature 0; for thinking models, we use the documented recommended hyperparameters, averaging results across five runs. When employing DeepSeek-R1 and GPT-4o as LLM-as-judge, we fix the temperature parameter at 0 to ensure that every evaluation remains deterministic. Further details are provided in Appendix C.2 and C.3.

6.3 Results

In the Text-to-Table task, LiveSum emerges as the most challenging benchmark in terms of content faithfulness, with all evaluated models (except O4-mini) scoring no higher than 70; InstructIE is intermediate, while Struc-Bench is the simplest, achieving nearly 90. Regarding structure coherence, both LiveSum and Struc-Bench achieve higher scores due to their predefined headers, reflecting strong instruction-following capabilities. In contrast, the absence of fixed headers in InstructIE prompts models to generate more diverse table formats, resulting in lower coherence scores. In the Text-to-Graph task, the DART dataset demonstrates higher faithfulness compared to InstructIE because both datasets, originally designed for KG-to-text, contain triples often irretrievable from text and this limitation is more pronounced in InstructIE and leads to lower recall. Meanwhile, both datasets show moderate performance in structure coherence, primarily due to challenges in entity and relation naming. In the Text-to-Chart task, larger models generally exhibit superior performance, consistently achieving scores above 80, while smaller models trained with SFT lag significantly behind. For further analysis, see Appendix D.

Model	Table						Graph				Chart	
	InstructIE (Jiao et al., 2023)		Struc-Bench (Tang et al., 2024)		LiveSum (Deng et al., 2024)		DART (Nan et al., 2021)		InstructIE (Gui et al., 2023)		Text2Chart31 (Zadeh et al., 2024)	
	FAITH.	COH.	FAITH.	COH.	FAITH.	COH.	FAITH.	COH.	FAITH.	COH.	FAITH.	COH.
◆ <i>SFT</i>												
Mistral-7B	58.08	68.85	87.65	90.05	53.31	98.54	75.07	83.17	44.58	69.67	64.54	63.41
Llama-3.1-8B	59.31	66.44	85.01	85.95	61.34	90.30	38.94	55.13	34.93	64.99	77.06	77.25
R1-Distill-Llama-8B	62.86	68.08	54.56	58.97	43.27	86.22	59.03	69.43	37.27	69.73	41.66	40.45
◆ <i>Zero-Shot</i>												
Phi-3-medium	29.15	39.89	74.04	83.99	41.61	91.73	64.60	72.59	43.40	71.44	60.24	60.45
Phi-3-medium (CoT)	30.36	37.23	73.42	81.06	46.36	93.19	64.88	72.32	41.28	70.44	57.75	58.27
Mistral-7B	53.00	59.73	56.19	67.05	30.14	87.26	70.27	73.06	42.74	66.54	52.25	52.03
Mistral-7B (CoT)	54.02	60.37	56.25	70.34	41.27	85.97	64.68	70.96	33.35	66.01	45.13	46.51
Mixtral-8x22B	56.35	68.25	84.16	85.20	49.12	87.91	65.70	71.07	40.07	68.89	61.62	62.15
Mixtral-8x22B (CoT)	58.27	68.29	80.55	74.94	54.87	87.37	71.27	74.53	47.47	71.29	58.92	58.89
Llama-3.1-405B	63.88	69.63	89.43	90.43	62.29	99.11	68.74	75.61	46.72	75.82	86.49	86.62
Llama-3.1-405B (CoT)	64.54	72.80	89.38	90.74	66.25	98.64	72.31	78.09	47.94	78.48	87.23	86.40
Llama-3.3-70B	62.88	69.19	80.58	80.00	58.48	99.19	71.09	77.54	47.11	76.30	87.44	87.12
Llama-3.3-70B (CoT)	65.34	69.46	86.47	85.84	62.79	98.68	77.11	78.99	46.56	78.51	88.78	87.83
GPT-3.5-turbo	55.09	69.01	78.44	85.37	45.38	79.18	69.75	76.72	41.89	73.34	85.17	84.79
GPT-3.5-turbo (CoT)	58.44	68.80	71.88	81.54	45.09	86.37	70.17	75.99	41.78	73.51	83.06	83.62
GPT-4o	70.72	76.36	87.35	92.22	63.21	91.99	74.15	77.18	46.77	77.00	85.48	86.01
GPT-4o (CoT)	68.11	74.40	86.05	90.03	62.16	98.50	79.89	79.33	53.28	77.65	85.99	86.63
QwQ-32B	68.23	68.79	81.21	87.12	58.06	96.20	69.42	78.02	56.30	82.15	86.62	86.71
DeepSeek-R1	76.86	74.75	87.77	91.05	64.89	96.85	78.25	79.89	52.15	79.23	86.07	86.90
Grok-3-mini	75.57	75.40	85.60	92.28	66.35	95.26	73.32	80.34	56.03	82.68	87.63	88.65
O4-mini	75.71	74.85	83.45	90.02	84.48	99.16	74.82	80.19	56.46	82.32	85.21	85.93

Table 2: The performance of major LLMs across multiple datasets. The highest results are **bolded**, with numerical scores categorized into five color-coded intervals: [90, 100] - excellent, [80, 90) - strong, [70, 80) - moderate, [60, 70) - limited, and [0, 60) - insufficient, using distinct chromatic gradients for visual differentiation.

Task	Metric	Faithfulness			Coherence		
		P(r)	S(ρ)	K(τ)	P(r)	S(ρ)	K(τ)
Table	ROUGE-L	0.267	0.227	0.166	0.607	0.617	0.452
	BERTScore	0.238	0.198	0.142	0.562	0.576	0.415
Graph	ROUGE-L	0.237	0.118	0.093	0.318	0.165	0.127
	BERTScore	0.100	0.043	0.036	0.211	0.105	0.078
Chart	CodeBLEU	0.249	0.186	0.144	0.266	0.297	0.220
	METEOR	0.281	0.241	0.186	0.317	0.376	0.274

Table 3: Correlation coefficients (Pearson/Spearman/Kendall) between traditional metrics and human evaluation across two dimensions on three different tasks.

6.4 Evaluation of Scoring Criteria

To demonstrate how our proposed scoring criteria address evaluation limitations, we aim to answer the following questions:

RQ1. Do traditional metrics align with human judgment criteria? (§ 6.4.1)

RQ2. Can LLM-based evaluation ensure consistent and nuanced assessment? (§ 6.4.2)

Furthermore, we provide a case study analyzing traditional metric failures and our framework’s strengths (§6.4.3).

6.4.1 Validity Gap (RQ1)

For table generation and graph generation tasks, prior works predominantly adopt ROUGE-L and BERTScore by extracting textual content from structured outputs and comparing them with

ground truth references (Wu et al., 2022; Jiao et al., 2023; Huang et al., 2024a). For chart generation tasks, existing approaches typically evaluate code similarity between generated and reference codes using metrics like METEOR and CodeBLEU (Zadeh et al., 2024). To investigate whether these metrics align with human evaluation, we employ three annotators to score outputs based on our two orthogonal criteria, *Faithfulness* and *Coherence*, with final scores averaged across annotators. We conduct correlation analyses between each automated metric and human-rated *Faithfulness* and *Coherence* following established methodologies (Jiao et al., 2023), with results summarized in Table 3. The table reveals that only in table generation tasks do ROUGE-L and BERTScore exhibit moderate correlations with Coherence ($r, \rho > 0.5$), whereas all other metrics exhibit weak alignment with human evaluation. This demonstrates traditional metrics’ limitations and underscores the urgent need for specialized metrics capable of reliable quality assessment in this domain.

6.4.2 Reliability and Discriminability (RQ2)

To investigate the alignment between LLM-based scoring and human evaluation under our proposed criteria, three annotators score 30 randomly sampled instances per dataset following the rubric out-

Task	Model	Prec.	Rec.	Align.	Clr.
Table	GPT-4o	14.90	15.79	23.76	15.14
	DeepSeek-R1	9.37	8.59	8.32	4.55
Graph	GPT-4o	14.76	11.08	7.16	7.42
	DeepSeek-R1	4.91	8.91	7.09	3.08
Chart	GPT-4o	15.32	14.77	6.12	7.91

Table 4: Root Mean Squared Error (RMSE) between human and LLM scores across three task categories.

Metric	Faithfulness		Coherence	
	Precision	Recall	Alignment	Clarity
Pearson (r)	0.745	0.771	0.761	0.728
Spearman (ρ)	0.747	0.766	0.712	0.679
Kendall (τ)	0.664	0.684	0.670	0.621

Table 5: Correlation coefficients (Person / Spearman / Kendall) of human and LLM scores for four evaluation metrics.

lined in the evaluation prompt. The Root Mean Square Error (RMSE) between human and LLM scores is averaged across three task categories and aggregated in Table 4, which reveals that DeepSeek-R1 demonstrates superior accuracy compared to GPT-4o, with subsequent case studies confirming that its enhanced reasoning fidelity directly contributes to reduced prediction errors. Given DeepSeek-R1’s current lack of multimodal input capabilities, GPT-4o remains our evaluator for charts to maintain assessment consistency. We further analyze the correlation between human ratings and DeepSeek-R1 scores, computing Pearson’s r , Spearman’s ρ , and Kendall’s τ coefficients, with results summarized in Table 5. All four evaluation metrics exhibit strong positive correlations with human judgments. To compare our proposed metrics with conventional metrics, we evaluated five models on the InstructIE (Jiao et al., 2023) dataset using both traditional metrics (ROUGE and BERTScore) and LLM-based assessments (FAITH. and COH. scores from GPT-4o and DeepSeek-R1). The results in Table 6 demonstrate that DeepSeek-R1 scores show a high positive correlation with GPT-4o, eliminating concerns about LLM evaluators favoring their own outputs. While GPT-3.5-turbo achieves the second-highest scores in traditional metrics, it ranks last across all LLM-based evaluations, validating our metrics’ capability to detect reliability flaws obscured by conventional measurements.

6.4.3 Case Study

We provide five carefully examined case studies tackling important shortcomings in automatic evaluation metrics in Appendix E. Our examination uncovers systematic misalignment situations where

Model	ROUGE	BERT.	FAITH _{R1}	FAITH _{4o}	COH _{R1}	COH _{4o}
Llama-3.1-405B	0.561	0.742	63.88	69.53	69.63	76.66
Llama-3.3-70B	0.517	0.738	62.88	67.23	69.19	75.08
GPT-3.5-turbo	0.586	0.770	55.09	62.86	69.01	75.82
GPT-4o	0.613	0.772	<u>70.72</u>	<u>73.07</u>	76.36	81.44
DeepSeek-R1	0.541	0.738	76.86	77.02	<u>74.75</u>	<u>79.81</u>

Table 6: Performance comparison of five LLMs on InstructIE (Jiao et al., 2023) dataset through conventional metrics and LLM-based evaluation, and we **bold** the best performance for each metric.

conventional metrics (BERTScore, ROUGE, etc.) incorrectly prefer outputs with: structural incompleteness (Case 1), factual hallucinations (Cases 2, 3), semantic disarray (Case 4), and numerical or logical errors (Cases 3, 5), whereas our model accurately predicts outputs, maintaining content faithfulness and structural coherence.

7 Future Opportunities

Future work may enhance the text-to-structure generation of agentic AI from multiple directions. First, optimizing real-time conversion with efficient algorithms and lightweight models can support dynamic environments, which is critical for tasks such as question answering and summarization. Second, improving accuracy and consistency through multimodal learning (e.g., integrating text with images) and self-reflection mechanisms can reduce errors in complex texts, thereby benefiting the information state conversion of agentic information retrieval. Third, reinforcement learning with sophisticated reward functions – balancing utility, efficiency, and ethics – can combine structured outputs with external goals to enhance autonomous decision-making capabilities. Finally, multi-task learning frameworks can simultaneously generate different structures, such as tables and charts, with consistency and flexibility. These efforts aim to improve the autonomy, accuracy, and applicability of agentic AI in complex real-world scenarios.

8 Conclusion

In conclusion, this work aims to provide a comprehensive review of currently existing methods, benchmark datasets, and evaluation metrics for text-to-structure extraction by designing a universal evaluation framework specifically tailored for structured outputs. Through these contributions, we not only synthesize current research progress but also elucidate how these technologies can empower agentic AI systems to achieve enhanced autonomy and efficiency in dynamic real-world environments.

Limitations

The limitations of this study mainly stem from two aspects. First, since we strategically focus the most representative studies in this field, the scope of the literature review may not comprehensively cover all research advances related to complex structures. Second, given the rapid development of LLMs, the Text-to-Structure benchmark evaluation may not cover the latest published LLMs. However, despite our best efforts to survey a comprehensive selection of important papers and mainstream LLMs, the practical constraints of a single submission inherently limit our ability to cover every aspect or all LLM variants. Future research directions should aim for a broader literature review while integrating datasets from cross-disciplinary studies to enhance their generalizability.

Ethics Statement

Our paper presents a comprehensive survey of text-to-complex structure extraction, with a specific focus on tables, graphs, and charts. The datasets and models employed in this study are all open-source. All datasets are used in their pre-anonymized forms as released by the original creators, complying with established licensing agreements. Our evaluation uses strictly unaltered benchmark data, applying only random sampling to select representative subsets. Crucially, no additions, modifications, or external data curation occur at any processing stage, preserving original dataset integrity while enabling efficient LLM assessment. The annotation was strictly limited to validating sampled structural outputs, conducted voluntarily by three doctoral researchers from the authors' institution. Therefore, to the best of the authors' knowledge, we believe that this study introduces no additional risk.

References

Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Björck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Weizhu Chen, Yen-Chun Chen, Yi-Ling Chen, Hao Cheng, Parul Chopra, Xiyang Dai, Matthew Dixon, Ronen Eldan, Victor Fragoso, Jianfeng Gao, Mei Gao, Min Gao, Amit Garg, Allie Del Giorno, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Wenxiang Hu, Jamie Huynh, Dan Iter, Sam Ade Jacobs, Mojan Javaheripi,

Xin Jin, Nikos Karampatziakis, Piero Kauffmann, Mahoud Khademi, Dongwoo Kim, Young Jin Kim, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Xihui Lin, Zeqi Lin, Ce Liu, Liyuan Liu, Mengchen Liu, Weishung Liu, Xiaodong Liu, Chong Luo, Piyush Madan, Ali Mahmoudzadeh, David Majercak, Matt Mazzola, Caio César Teodoro Mendes, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Liliang Ren, Gustavo de Rosa, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Yelong Shen, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Praneetha Vaddamanu, Chunyu Wang, Guanhua Wang, Lijuan Wang, Shuohang Wang, Xin Wang, Yu Wang, Rachel Ward, Wen Wen, Philipp Witte, Haiping Wu, Xiaoxia Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Jilong Xue, Sonali Yadav, Fan Yang, Jianwei Yang, Yifan Yang, Ziyi Yang, Donghan Yu, Lu Yuan, Chenruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. [Phi-3 technical report: A highly capable language model locally on your phone](#). *Preprint*, arXiv:2404.14219.

Oshin Agarwal, Heming Ge, Siamak Shakeri, and Rami Al-Rfou. 2021. [Knowledge graph based synthetic corpus generation for knowledge-enhanced language model pre-training](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 3554–3565. Association for Computational Linguistics.

Naman Ahuja, Fenil Denish Bardoliya, Chitta Baral, and Vivek Gupta. 2025. [Map&make: Schema guided text to table generation](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2025, Vienna, Austria, July 27 - August 1, 2025*, pages 30249–30262. Association for Computational Linguistics.

Mubashara Akhtar, Oana Cocarascu, and Elena Simperl. 2023. [Reading and reasoning over chart images for evidence-based automated fact-checking](#). In *Findings of the Association for Computational Linguistics: EACL 2023, Dubrovnik, Croatia, May 2-6, 2023*, pages 399–414. Association for Computational Linguistics.

Jinheon Baek, Sujay Kumar Jauhar, Silviu Cucerzan, and Sung Ju Hwang. 2024. [Researchagent: Iterative research idea generation over scientific literature with large language models](#). *CoRR*, abs/2404.07738.

Junwei Bao, Duyu Tang, Nan Duan, Zhao Yan, Yuanhua Lv, Ming Zhou, and Tiejun Zhao. 2018. [Table-to-text: Describing table region with natural language](#). In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th Innovative Applications of Artificial Intelligence (IAAI-*

661	18), and the 8th AAAI Symposium on Educational	DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang,	718
662	Advances in Artificial Intelligence (EAAI-18), New	Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu,	719
663	Orleans, Louisiana, USA, February 2-7, 2018, pages	Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang,	720
664	5020–5027. AAAI Press.	Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong	721
665	Markus J. Buehler. 2025. Agentic deep graph reasoning	Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue,	722
666	yields self-organizing knowledge networks . <i>CoRR</i> ,	Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu,	723
667	abs/2502.13025.	Chenggang Zhao, Chengqi Deng, Chenyu Zhang,	724
668	Pere-Lluís Huguet Cabot and Roberto Navigli. 2021.	Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji,	725
669	REBEL: relation extraction by end-to-end language	Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo,	726
670	generation . In <i>Findings of the Association for Com-</i>	Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang,	727
671	<i>putational Linguistics: EMNLP 2021, Virtual Event /</i>	Han Bao, Hanwei Xu, Haocheng Wang, Honghui	728
672	<i>Punta Cana, Dominican Republic, 16-20 November,</i>	Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li,	729
673	<i>2021</i> , pages 2370–2381. Association for Computa-	Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang	730
674	<i>tional Linguistics.</i>	Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L.	731
675	Avi Caciularu, Alon Jacovi, Eyal Ben-David, Sasha	Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai	732
676	Goldshtein, Tal Schuster, Jonathan Herzig, Gal Eli-	Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai	733
677	dan, and Amir Globerson. 2024. TACT: advancing	Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong	734
678	complex aggregative reasoning with information ex-	Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan	735
679	traction tools . <i>CoRR</i> , abs/2406.03618.	Zhang, Minghua Zhang, Minghui Tang, Meng Li,	736
680	Salvatore Carta, Alessandro Giuliani, Leonardo Piano,	Miaojun Wang, Mingming Li, Ning Tian, Panpan	737
681	Alessandro Sebastian Podda, Livio Pompianu, and	Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen,	738
682	Sandro Gabriele Tiddia. 2023. Iterative zero-shot	Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan,	739
683	LLM prompting for knowledge graph construction .	Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen,	740
684	<i>CoRR</i> , abs/2307.01128.	Shanghao Lu, Shangyan Zhou, Shanhuang Chen,	741
685	Bohan Chen and Andrea L. Bertozzi. 2023. Autokg:	Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng	742
686	Efficient automated knowledge graph generation for	Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing	743
687	language models . In <i>IEEE International Conference</i>	Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun,	744
688	<i>on Big Data, BigData 2023, Sorrento, Italy, Decem-</i>	T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu,	745
689	<i>ber 15-18, 2023</i> , pages 3117–3126. IEEE.	Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao	746
690	Mingda Chen, Sam Wiseman, and Kevin Gimpel. 2021.	Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan	747
691	Wikitablet: A large-scale data-to-text dataset for gen-	Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin	748
692	erating wikipedia article sections . In <i>Findings of the</i>	Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li,	749
693	<i>Association for Computational Linguistics: ACL/IJC-</i>	Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin,	750
694	<i>NLP 2021, Online Event, August 1-6, 2021</i> , volume	Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxi-	751
695	<i>ACL/IJCNLP 2021 of Findings of ACL</i> , pages 193–	ang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang,	752
696	209. Association for Computational Linguistics.	Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang	753
697	Wenhu Chen, Jianshu Chen, Yu Su, Zhiyu Chen, and	Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng	754
698	William Yang Wang. 2020. Logical natural language	Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi,	755
699	generation from open-domain tables . In <i>Proceed-</i>	Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang,	756
700	<i>ings of the 58th Annual Meeting of the Association</i>	Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo,	757
701	<i>for Computational Linguistics, ACL 2020, Online,</i>	Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yu-	758
702	<i>July 5-10, 2020</i> , pages 7929–7942. Association for	jia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You,	759
703	<i>Computational Linguistics.</i>	Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu,	760
704	David Cheng-Han Chiang and Hung-yi Lee. 2023. Can	Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu,	761
705	large language models be an alternative to human	Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan,	762
706	evaluations? In <i>Proceedings of the 61st Annual</i>	Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean	763
707	<i>Meeting of the Association for Computational Lin-</i>	Xu, Zhenda Xie, Zhengyan Zhang, Zhewen Hao,	764
708	<i>guistics (Volume 1: Long Papers), ACL 2023, Toronto,</i>	Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zi-	765
709	<i>Canada, July 9-14, 2023</i> , pages 15607–15631. Asso-	jia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song,	766
710	<i>ciation for Computational Linguistics.</i>	Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu	767
711	Anthony M. Colas, Ali Sadeghian, Yue Wang, and	Zhang, and Zhen Zhang. 2025. Deepseek-r1: Incent-	768
712	Daisy Zhe Wang. 2021. Eventnarrative: A large-	ivizing reasoning capability in llms via reinforce-	769
713	scale event-centric dataset for knowledge graph-to-	ment learning . <i>Preprint</i> , arXiv:2501.12948.	770
714	text generation . In <i>Proceedings of the Neural Infor-</i>	Zheye Deng, Chunkit Chan, Weiqi Wang, Yuxi Sun,	771
715	<i>mation Processing Systems Track on Datasets and</i>	Wei Fan, Tianshi Zheng, Yauwai Yim, and Yangqiu	772
716	<i>Benchmarks 1, NeurIPS Datasets and Benchmarks</i>	Song. 2024. Text-tuple-table: Towards information	773
717	<i>2021, December 2021, virtual.</i>	integration in text-to-table generation via global tuple	774
		extraction . In <i>Proceedings of the 2024 Conference on</i>	775
		<i>Empirical Methods in Natural Language Processing,</i>	776
		<i>EMNLP 2024, Miami, FL, USA, November 12-16,</i>	777
		<i>2024</i> , pages 9300–9322. Association for Computa-	778
		<i>tional Linguistics.</i>	779

780	Zheyang Deng, Weiqi Wang, Zhaowei Wang, Xin Liu, and	<i>Language Generation, INLG 2017, Santiago de Com-</i>	840
781	Yangqiu Song. 2023. Gold: A global and local-aware	<i>postela, Spain, September 4-7, 2017</i> , pages 124–133.	841
782	denoising framework for commonsense knowledge	Association for Computational Linguistics.	842
783	graph noise detection . In <i>Findings of the Association</i>		
784	<i>for Computational Linguistics: EMNLP 2023,</i>	Alireza Ghafarollahi and Markus J. Buehler. 2024. Sci-	843
785	<i>Singapore, December 6-10, 2023</i> , pages 3591–3608.	agents: Automating scientific discovery through	844
786	Association for Computational Linguistics.	multi-agent intelligent graph reasoning . <i>CoRR</i> ,	845
		abs/2409.05556 .	846
787	Dheeru Dua, Shivanshu Gupta, Sameer Singh, and Matt	Seethalakshmi Gopalakrishnan, Victor Zitian Chen,	847
788	Gardner. 2022. Successive prompting for decompos-	Wenwen Dou, Gus Hahn-Powell, Sreekar Nedunuri,	848
789	ing complex questions . In <i>Proceedings of the 2022</i>	and Wlodek Zadrozny. 2023. Text to causal knowl-	849
790	<i>Conference on Empirical Methods in Natural Lan-</i>	edge graph: A framework to synthesize knowledge	850
791	<i>guage Processing, EMNLP 2022, Abu Dhabi, United</i>	from unstructured business texts into causal graphs .	851
792	<i>Arab Emirates, December 7-11, 2022</i> , pages 1251–	<i>Inf.</i> , 14(7):367.	852
793	1265. Association for Computational Linguistics.		
794	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,	Honghao Gui, Shuofei Qiao, Jintian Zhang, Hong-	853
795	Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,	bin Ye, Mengshu Sun, Lei Liang, Huajun Chen,	854
796	Akhil Mathur, Alan Schelten, Amy Yang, Angela	and Ningyu Zhang. 2023. Instructie: A bilin-	855
797	Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang,	gual instruction-based information extraction dataset .	856
798	Archi Mitra, Archie Sravankumar, Artem Korenev,	<i>CoRR</i> , abs/2305.11527 .	857
799	Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien		
800	Rodriguez, Austen Gregerson, Ava Spataru, Bap-	Mandy Guo, Zihang Dai, Denny Vrandečić, and Rami	858
801	tiste Rozière, Bethany Biron, Binh Tang, Bobbie	Al-Rfou. 2020. Wiki-40b: Multilingual language	859
802	Chern, Charlotte Caucheteux, Chaya Nayak, Chloe	model dataset . In <i>Proceedings of The 12th Language</i>	860
803	Bi, Chris Marra, Chris McConnell, Christian Keller,	<i>Resources and Evaluation Conference, LREC 2020,</i>	861
804	Christophe Touret, Chunyang Wu, Corinne Wong,	<i>Marseille, France, May 11-16, 2020</i> , pages 2440–	862
805	Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-	2452. European Language Resources Association.	863
806	lonsius, Daniel Song, Danielle Pintz, Danny Livshits,		
807	David Esiobu, Dhruv Choudhary, Dhruv Mahajan,	Rishin Haldar and Debajyoti Mukhopadhyay. 2011.	864
808	Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,	Levenshtein distance technique in dictionary lookup	865
809	Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,	methods: An improved approach . <i>CoRR</i> ,	866
810	Emily Dinan, Eric Michael Smith, Filip Radenovic,	abs/1101.1232 .	867
811	Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georg-		
812	ia Lewis Anderson, Graeme Nail, Grégoire Mialon,	Yucheng Han, Chi Zhang, Xin Chen, Xu Yang, Zhibin	868
813	Guan Pang, Guillem Cucurell, Hailey Nguyen, Han-	Wang, Gang Yu, Bin Fu, and Hanwang Zhang. 2023.	869
814	nah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov,	Chartlama: A multimodal LLM for chart understand-	870
815	Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan	ing and generation . <i>CoRR</i> , abs/2311.16483 .	871
816	Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan		
817	Geffert, Jana Vranes, Jason Park, Jay Mahadeokar,	Han He and Jinho D. Choi. 2021. The stem cell hypoth-	872
818	Jeet Shah, Jelmer van der Linde, Jennifer Billock,	esis: Dilemma behind multi-task learning with trans-	873
819	Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi,	former encoders . In <i>Proceedings of the 2021 Confer-</i>	874
820	Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu,	<i>ence on Empirical Methods in Natural Language Pro-</i>	875
821	Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph	<i>cessing, EMNLP 2021, Virtual Event / Punta Cana,</i>	876
822	Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia,	<i>Dominican Republic, 7-11 November, 2021</i> , pages	877
823	Kalyan Vasuden Alwala, Kartikeya Upasani, Kate	5555–5577. Association for Computational Linguis-	878
824	Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and	<i>tics</i> .	879
825	et al. 2024. The llama 3 herd of models . <i>CoRR</i> ,		
826	abs/2407.21783 .	Alexander Miserlis Hoyle, Rupak Sarkar, Pranav Goel,	880
		and Philip Resnik. 2023. Natural language decom-	881
827	Darren Edge, Ha Trinh, Newman Cheng, Joshua	positions of implicit content enable better text repre-	882
828	Bradley, Alex Chao, Apurva Mody, Steven Truitt,	sentations . In <i>Proceedings of the 2023 Conference</i>	883
829	and Jonathan Larson. 2024. From local to global: A	<i>on Empirical Methods in Natural Language Process-</i>	884
830	graph RAG approach to query-focused summariza-	<i>ing, EMNLP 2023, Singapore, December 6-10, 2023</i> ,	885
831	tion . <i>CoRR</i> , abs/2404.16130 .	pages 13188–13214. Association for Computational	886
		Linguistics.	887
832	Mohamed Elhoseiny and Ahmed M. Elgammal. 2016.	Mengting Hu, Honglei Guo, Shiwan Zhao, Hang Gao,	888
833	Text to multi-level mindmaps - A novel method for	and Zhong Su. 2021. Efficient mind-map generation	889
834	hierarchical visual abstraction of natural language	via sequence-to-graph and reinforced graph refine-	890
835	text . <i>Multim. Tools Appl.</i> , 75(8):4217–4244.	ment . In <i>Proceedings of the 2021 Conference on</i>	891
836	Claire Gardent, Anastasia Shimorina, Shashi Narayan,	<i>Empirical Methods in Natural Language Processing,</i>	892
837	and Laura Perez-Beltrachini. 2017. The webnlg chal-	<i>EMNLP 2021, Virtual Event / Punta Cana, Domini-</i>	893
838	lenge: Generating text from RDF data . In <i>Proceed-</i>	<i>can Republic, 7-11 November, 2021</i> , pages 8130–	894
839	<i>ings of the 10th International Conference on Natural</i>	8141. Association for Computational Linguistics.	895

896	Haoyu Huang, Chong Chen, Conghui He, Yang Li, Jiawei Jiang, and Wentao Zhang. 2024a. Can llms be good graph judger for knowledge graph construction? <i>Preprint</i> , arXiv:2411.17388.	953
897		954
898		955
899		956
900	Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. 2024b. Mlagentbench: Evaluating language agents on machine learning experimentation . In <i>Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024</i> . OpenReview.net.	957
901		958
902		959
903		960
904		961
905		962
906	Parag Jain, Andreea Marzoca, and Francesco Piccinno. 2024. STRUCTSUM generation for faster text comprehension . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024</i> , pages 7876–7896. Association for Computational Linguistics.	963
907		964
908		965
909		966
910		967
911		968
912		969
913	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de Las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L��lio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth��e Lacroix, and William El Sayed. 2023a. Mistral 7b . <i>CoRR</i> , abs/2310.06825.	970
914		971
915		972
916		973
917		974
918		975
919		976
920		977
921	Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L��lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th��ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth��e Lacroix, and William El Sayed. 2024a. Mixtral of experts . <i>Preprint</i> , arXiv:2401.04088.	978
922		979
923		980
924		981
925		982
926		983
927		984
928		985
929		986
930		987
931		988
932	Jinhao Jiang, Kun Zhou, Zican Dong, Keming Ye, Xin Zhao, and Ji-Rong Wen. 2023b. Structgpt: A general framework for large language model to reason over structured data . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 9237–9251. Association for Computational Linguistics.	989
933		990
934		991
935		992
936		993
937		994
938		995
939		996
940	Peiwen Jiang, Xinbo Lin, Zibo Zhao, Ruhui Ma, Yvonne Chen, and Jinhua Cheng. 2024b. TKGT: redefinition and A new way of text-to-table tasks based on real world demands and knowledge graphs augmented llms . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 16112–16126. Association for Computational Linguistics.	997
941		998
942		999
943		1000
944		1001
945		1002
946		1003
947		1004
948		1005
949	Yizhu Jiao, Ming Zhong, Sha Li, Ruining Zhao, Siru Ouyang, Heng Ji, and Jiawei Han. 2023. Instruct and extract: Instruction tuning for on-demand information extraction . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023</i> , pages 10030–10051. Association for Computational Linguistics.	1006
950		1007
951		1008
952		1009
	Bowen Jin, Chulin Xie, Jiawei Zhang, Kashob Kumar Roy, Yu Zhang, Zheng Li, Ruirui Li, Xianfeng Tang, Suhang Wang, Yu Meng, and Jiawei Han. 2024. Graph chain-of-thought: Augmenting large language models by reasoning on graphs . In <i>Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024</i> , pages 163–184. Association for Computational Linguistics.	
	Zhijing Jin, Qipeng Guo, Xipeng Qiu, and Zheng Zhang. 2020. Genwiki: A dataset of 1.3 million content-sharing text and graphs for unsupervised graph-to-text generation . In <i>Proceedings of the 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020</i> , pages 2398–2409. International Committee on Computational Linguistics.	
	Shankar Kantharaj, Rixie Tiffany Ko Leong, Xiang Lin, Ahmed Masry, Megh Thakkar, Enamul Hoque, and Shafiq R. Joty. 2022. Chart-to-text: A large-scale benchmark for chart summarization . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022</i> , pages 4005–4023. 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 <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	
	Rik Koncel-Kedziorski, Dhanush Bekal, Yi Luan, Mirella Lapata, and Hannaneh Hajishirzi. 2019. Text generation from knowledge graphs with graph transformers . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)</i> , pages 2284–2293. Association for Computational Linguistics.	
	Aditya Kulkarni, Hetansh Shah, Lynette D’Mello, and Krish Shah. 2023. Flowchart generation and mind map creation using extracted summarized text . In <i>2023 International Conference on Recent Advances in Science and Engineering Technology (ICRASET)</i> , pages 1–6.	
	Anirban Laha, Parag Jain, Abhijit Mishra, and Karthik Sankaranarayanan. 2019. Scalable micro-planned generation of discourse from structured data . <i>Comput. Linguistics</i> , 45(4):737–763.	

- Rémi Lebret, David Grangier, and Michael Auli. 2016. [Neural text generation from structured data with application to the biography domain](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pages 1203–1213. The Association for Computational Linguistics. 1067–1072
- Long Li, Weiwen Xu, Jiayan Guo, Ruochen Zhao, Xingxuan Li, Yuqian Yuan, Boqiang Zhang, Yuming Jiang, Yifei Xin, Ronghao Dang, Deli Zhao, Yu Rong, Tian Feng, and Lidong Bing. 2024a. [Chain of ideas: Revolutionizing research via novel idea development with LLM agents](#). *CoRR*, abs/2410.13185. 1073–1077
- Tong Li, Zhihao Wang, Liangying Shao, Xuling Zheng, Xiaoli Wang, and Jinsong Su. 2023. [A sequence-to-sequence&set model for text-to-table generation](#). In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 5358–5370. Association for Computational Linguistics. 1078–1083
- Yinghao Li, Rampi Ramprasad, and Chao Zhang. 2024b. [A simple but effective approach to improve structured language model output for information extraction](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024, Miami, Florida, USA, November 12-16, 2024*, pages 5133–5148. Association for Computational Linguistics. 1084–1091
- Zhuoqun Li, Xuanang Chen, Haiyang Yu, Hongyu Lin, Yaojie Lu, Qiaoyu Tang, Fei Huang, Xianpei Han, Le Sun, and Yongbin Li. 2024c. [Struc-trag: Boosting knowledge intensive reasoning of llms via inference-time hybrid information structurization](#). *CoRR*, abs/2410.08815. 1092–1099
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics. 1100–1102
- Yupian Lin, Tong Ruan, Jingping Liu, and Haofen Wang. 2024. [A survey on neural data-to-text generation](#). *IEEE Trans. Knowl. Data Eng.*, 36(4):1431–1449. 1103–1106
- Jiacheng Liu, Wenya Wang, Dianzhuo Wang, Noah A. Smith, Yejin Choi, and Hannaneh Hajishirzi. 2023a. [Vera: A general-purpose plausibility estimation model for commonsense statements](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 1264–1287. Association for Computational Linguistics. 1107–1113
- Ruixue Liu, Baoyang Chen, Meng Chen, Youzheng Wu, Zhijie Qiu, and Xiaodong He. 2019a. [Mappa mundi: An interactive artistic mind map generator with artificial imagination](#). In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, pages 6545–6547. ijcai.org. 1114–1121
- Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023b. [G-eval: NLG evaluation using gpt-4 with better human alignment](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 2511–2522. Association for Computational Linguistics. 1122–1123
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019b. [Roberta: A robustly optimized BERT pretraining approach](#). *CoRR*, abs/1907.11692. 1124–1127
- Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. [Multi-task identification of entities, relations, and coreference for scientific knowledge graph construction](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 3219–3232. Association for Computational Linguistics. 1128–1135
- Linhao Luo, Yuan-Fang Li, Gholamreza Haffari, and Shirui Pan. 2024. [Reasoning on graphs: Faithful and interpretable large language model reasoning](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net. 1136–1143
- Ahmed Masry, Parsa Kavehzadeh, Do Xuan Long, Enamul Hoque, and Shafiq Joty. 2023. [Unichart: A universal vision-language pretrained model for chart comprehension and reasoning](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 14662–14684. Association for Computational Linguistics. 1144–1151
- Ahmed Masry, Mehrad Shahmohammadi, Md. Rizwan Parvez, Enamul Hoque, and Shafiq Joty. 2024. [Chartinstruct: Instruction tuning for chart comprehension and reasoning](#). In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 10387–10409. Association for Computational Linguistics. 1152–1159
- Igor Melnyk, Pierre L. Dognin, and Payel Das. 2022. [Knowledge graph generation from text](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, pages 1610–1622. Association for Computational Linguistics. 1160–1167
- Nandana Mihindukulasooriya, Sanju Tiwari, Carlos F. Enguix, and Kusum Lata. 2023. [Text2kgbench: A benchmark for ontology-driven knowledge graph generation from text](#). In *The Semantic Web - ISWC 2023 - 22nd International Semantic Web Conference, Athens, Greece, November 6-10, 2023, Proceedings, Part II*, volume 14266 of *Lecture Notes in Computer Science*, pages 247–265. Springer. 1168–1175
- Antonio Miranda-Escalada, Farrokh Mehryary, Jouni Luoma, Darryl Estrada-Zavala, Luis Gascó, Sampo 1122–1123

1124	Pyysalo, Alfonso Valencia, and Martin Krallinger.	to-end generation. In <i>Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue, Saarbrücken, Germany, August 15-17, 2017</i> , pages	1181
1125	2023. Overview of drugprot task at biocreative	<i>Saarbrücken, Germany, August 15-17, 2017</i> , pages	1182
1126	VII: data and methods for large-scale text min-	201–206. Association for Computational Linguistics.	1183
1127	ing and knowledge graph generation of heteroge-		1184
1128	nous chemical-protein relations. <i>Database J. Biol.</i>		
1129	<i>Databases Curation</i> , 2023.		
1130	Belinda Mo, Kyssen Yu, Joshua Kazdan, Proud Mpala,	Jason Obeid and Enamul Hoque. 2020. Chart-to-text: Generating natural language descriptions for charts by adapting the transformer model . In <i>Proceedings of the 13th International Conference on Natural Language Generation, INLG 2020, Dublin, Ireland, December 15-18, 2020</i> , pages 138–147. Association for	1185
1131	Lisa Yu, Chris Cundy, Charilaos Kanatsoulis, and	Computational Linguistics.	1186
1132	Sanmi Koyejo. 2025. Kggen: Extracting knowl-		1187
1133	edge graphs from plain text with language models .		1188
1134	<i>Preprint</i> , arXiv:2502.09956.		1189
1135	Sabah Mohammed and Jinan Fiaidhi. 2023. Establish-	OpenAI. 2022. Introducing chatgpt .	1192
1136	ment of a mindmap for medical e-diagnosis as a ser-		
1137	vice for graph-based learning and analytics . <i>Neural</i>	OpenAI. 2024. Gpt-4 technical report . <i>Preprint</i> ,	1193
1138	<i>Comput. Appl.</i> , 35(22):16089–16100.	arXiv:2303.08774.	1194
1139	Rizki Elisa Nalawati and Azka Dini Yuntari. 2021. Rat-	OpenAI. 2025. Introducing openai o3 and o4-mini .	1195
1140	cliff/obershelp algorithm as an automatic assessment		
1141	on e-learning . In <i>2021 4th International Conference</i>	Ankur P. Parikh, Xuezhi Wang, Sebastian Gehrmann,	1196
1142	<i>of Computer and Informatics Engineering (IC2IE)</i> ,	Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and	1197
1143	pages 244–248.	Dipanjan Das. 2020. Totto: A controlled table-to-text generation dataset . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 1173–1186. Association for Computa-	1198
1144	Linyong Nan, Dragomir R. Radev, Rui Zhang, Am-	tional Linguistics.	1199
1145	rit Rau, Abhinand Sivaprasad, Chiachun Hsieh, Xi-		1200
1146	angru Tang, Aadit Vyas, Neha Verma, Pranav Kr-	Michal Pietruszka, Michal Turski, Lukasz Borchmann,	1201
1147	ishna, Yangxiaokang Liu, Nadia Irwanto, Jessica	Tomasz Dwojak, Gabriela Nowakowska, Karolina	1202
1148	Pan, Faiaz Rahman, Ahmad Zaidi, Mutethia Mutuma,	Szyndler, Dawid Jurkiewicz, and Lukasz Garncarek.	1203
1149	Yasin Tarabar, Ankit Gupta, Tao Yu, Yi Chern Tan,	2024. Stable: Table generation framework for encoder-decoder models . In <i>Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2024 - Volume 1: Long Papers, St. Julian's, Malta, March 17-22, 2024</i> , pages 2454–2472. Association for Com-	1204
1150	Xi Victoria Lin, Caiming Xiong, Richard Socher, and	putational Linguistics.	1205
1151	Nazneen Fatema Rajani. 2021. DART: open-domain		1206
1152	structured data record to text generation . In <i>Proceed-</i>	Maja Popovic. 2015. chrF: character n-gram f-score for automatic MT evaluation . In <i>Proceedings of the Tenth Workshop on Statistical Machine Translation, WMT@EMNLP 2015, 17-18 September 2015, Lisbon, Portugal</i> , pages 392–395. The Association for	1207
1153	<i>ings of the 2021 Conference of the North American</i>	Computer Linguistics.	1218
1154	<i>Chapter of the Association for Computational Lin-</i>		1219
1155	guistics: Human Language Technologies, NAACL-	Matt Post. 2018. A call for clarity in reporting BLEU	1220
1156	HLT 2021, Online, June 6-11, 2021 , pages 432–447.	scores . In <i>Proceedings of the Third Conference on</i>	1221
1157	Association for Computational Linguistics.	<i>Machine Translation: Research Papers, WMT 2018,</i>	1222
1158	Benjamin Newman, Yoonjoo Lee, Aakanksha Naik,	<i>Belgium, Brussels, October 31 - November 1, 2018</i> ,	1223
1159	Pao Siangliulue, Raymond Fok, Juho Kim, Daniel S.	pages 186–191. Association for Computational Lin-	1224
1160	Weld, Joseph Chee Chang, and Kyle Lo. 2024. Arx-	guistics.	1225
1161	ivdigestables: Synthesizing scientific literature into		
1162	tables using language models . In <i>Proceedings of</i>	Ratish Puduppully, Li Dong, and Mirella Lapata. 2019.	1226
1163	<i>the 2024 Conference on Empirical Methods in Natu-</i>	Data-to-text generation with entity modeling . In <i>Pro-</i>	1227
1164	<i>ral Language Processing, EMNLP 2024, Miami,</i>	<i>ceedings of the 57th Conference of the Association</i>	1228
1165	<i>FL, USA, November 12-16, 2024</i> , pages 9612–9631.	<i>for Computational Linguistics, ACL 2019, Florence,</i>	1229
1166	Association for Computational Linguistics.	<i>Italy, July 28- August 2, 2019, Volume 1: Long Pa-</i>	1230
1167	David D. Nguyen, David Liebowitz, Surya Nepal,	<i>pers</i> , pages 2023–2035. Association for Computa-	1231
1168	Salil S. Kanhere, and Sharif Abuadbba. 2024. Con-	tional Linguistics.	1232
1169	textual chart generation for cyber deception . <i>CoRR</i> ,		
1170	abs/2404.04854.	Qwen. 2025. Qwq-32b: Embracing the power of rein-	1233
1171	Feng Nie, Jinpeng Wang, Jin-Ge Yao, Rong Pan, and	forcement learning .	1234
1172	Chin-Yew Lin. 2018. Operation-guided neural net-		
1173	works for high fidelity data-to-text generation . In		
1174	<i>Proceedings of the 2018 Conference on Empirical</i>		
1175	<i>Methods in Natural Language Processing, Brussels,</i>		
1176	<i>Belgium, October 31 - November 4, 2018</i> , pages		
1177	3879–3889. Association for Computational Linguis-		
1178	tics.		
1179	Jekaterina Novikova, Ondrej Dusek, and Verena Rieser.		
1180	2017. The E2E dataset: New challenges for end-		

Pritika Ramu, Aparna Garimella, and Sambaran Bandyopadhyay. 2024. Is this a bad table? A closer look at the evaluation of table generation from text. In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024</i> , pages 22206–22216. Association for Computational Linguistics.	text-to-table performance in large language models. <i>CoRR</i> , abs/2405.12174.	1293 1294
Md. Mahinur Rashid, Hasin Kawsar Jahan, Annysha Huzzat, Riyasaat Ahmed Rahul, Tamim Bin Zakir, Farhana Meem, Md. Saddam Hossain Mukta, and Swakkhar Shatabda. 2022. Text2chart: A multi-staged chart generator from natural language text. In <i>Advances in Knowledge Discovery and Data Mining - 26th Pacific-Asia Conference, PAKDD 2022, Chengdu, China, May 16-19, 2022, Proceedings, Part II</i> , volume 13281 of <i>Lecture Notes in Computer Science</i> , pages 3–16. Springer.	Aditi Singh, Abul Ehtesham, Saket Kumar, and Tala Talaie Khoei. 2025. Agentic retrieval-augmented generation: A survey on agentic RAG. <i>CoRR</i> , abs/2501.09136.	1295 1296 1297 1298
Emily Reif, Crystal Qian, James Wexler, and Minsuk Kahng. 2024. Automatic histograms: Leveraging language models for text dataset exploration. In <i>Extended Abstracts of the CHI Conference on Human Factors in Computing Systems, CHI EA 2024, Honolulu, HI, USA, May 11-16, 2024</i> , pages 53:1–53:9. ACM.	Anirudh Sundar, Christopher Richardson, William Gay, and Larry Heck. 2024a. itbbs: A dataset of interactive conversations over tabular information. <i>CoRR</i> , abs/2404.12580.	1299 1300 1301 1302
Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019</i> , pages 3980–3990. Association for Computational Linguistics.	Anirudh Sundar, Christopher Richardson, and Larry Heck. 2024b. gtbbs: Generating tables from text by conditional question answering. <i>CoRR</i> , abs/2403.14457.	1303 1304 1305 1306
Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton-Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. Code llama: Open foundation models for code. <i>CoRR</i> , abs/2308.12950.	Xiangru Tang, Yiming Zong, Jason Phang, Yilun Zhao, Wangchunshu Zhou, Arman Cohan, and Mark Gestein. 2024. Struc-bench: Are large language models good at generating complex structured tabular data? In <i>Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies: Short Papers, NAACL 2024, Mexico City, Mexico, June 16-21, 2024</i> , pages 12–34. Association for Computational Linguistics.	1307 1308 1309 1310 1311 1312 1313 1314 1315 1316
Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the conll-2003 shared task: Language-independent named entity recognition. In <i>Proceedings of the Seventh Conference on Natural Language Learning, CoNLL 2003, Held in cooperation with HLT-NAACL 2003, Edmonton, Canada, May 31 - June 1, 2003</i> , pages 142–147. ACL.	Yuan Tian, Weiwei Cui, Dazhen Deng, Xinjing Yi, Yurun Yang, Haidong Zhang, and Yingcai Wu. 2023. Chartgpt: Leveraging llms to generate charts from abstract natural language. <i>CoRR</i> , abs/2311.01920.	1317 1318 1319 1320
Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. 2020. BLEURT: learning robust metrics for text generation. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020</i> , pages 7881–7892. Association for Computational Linguistics.	Jiaqi Wang, Yuying Chang, Zhong Li, Ning An, Qi Ma, Lei Hei, Haibo Luo, Yifei Lu, and Feiliang Ren. 2024a. Techgpt-2.0: A large language model project to solve the task of knowledge graph construction. <i>CoRR</i> , abs/2401.04507.	1321 1322 1323 1324 1325
Haoxiang Shi, Jiaan Wang, Jiarong Xu, Cen Wang, and Tetsuya Sakai. 2024. Ct-eval: Benchmarking chinese	Xiao Wang, Weikang Zhou, Can Zu, Han Xia, Tianze Chen, Yuansen Zhang, Rui Zheng, Junjie Ye, Qi Zhang, Tao Gui, Jihua Kang, Jingsheng Yang, Siyuan Li, and Chunsai Du. 2023a. Instructuie: Multi-task instruction tuning for unified information extraction. <i>CoRR</i> , abs/2304.08085.	1326 1327 1328 1329 1330 1331
	Yiming Wang, Zhuosheng Zhang, and Rui Wang. 2023b. Element-aware summarization with large language models: Expert-aligned evaluation and chain-of-thought method. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 8640–8665. Association for Computational Linguistics.	1332 1333 1334 1335 1336 1337 1338 1339
	Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu, David Wadden, Kelsey MacMillan, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023c. How far can camels go? exploring the state of instruction tuning on open resources. In <i>Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023</i> .	1340 1341 1342 1343 1344 1345 1346 1347 1348 1349

1350	Yuhao Wang, Ruiyang Ren, Junyi Li, Xin Zhao,	Ruobing Xie, Zhiyuan Liu, Fen Lin, and Leyu Lin. 2018.	1407
1351	Jing Liu, and Ji-Rong Wen. 2024b. REAR: A	Does william shakespeare REALLY write hamlet?	1408
1352	relevance-aware retrieval-augmented framework for	knowledge representation learning with confidence.	1409
1353	open-domain question answering. In <i>Proceedings</i>	In <i>Proceedings of the Thirty-Second AAAI Confer-</i>	1410
1354	<i>of the 2024 Conference on Empirical Methods in</i>	<i>ence on Artificial Intelligence, (AAAI-18), the 30th in-</i>	1411
1355	<i>Natural Language Processing, EMNLP 2024, Miami,</i>	<i>novative Applications of Artificial Intelligence (IAAI-</i>	1412
1356	<i>FL, USA, November 12-16, 2024,</i> pages 5613–5626.	<i>18), and the 8th AAAI Symposium on Educational</i>	1413
1357	Association for Computational Linguistics.	<i>Advances in Artificial Intelligence (EAAI-18), New</i>	1414
1358	Zilong Wang, Hao Zhang, Chun-Liang Li, Julian Martin	<i>Orleans, Louisiana, USA, February 2-7, 2018,</i> pages	1415
1359	Eisenschlos, Vincent Perot, Zifeng Wang, Lesly Mi-	4954–4961. AAAI Press.	1416
1360	culichich, Yasuhisa Fujii, Jingbo Shang, Chen-Yu Lee,	Derong Xu, Wei Chen, Wenjun Peng, Chao Zhang, Tong	1417
1361	and Tomas Pfister. 2024c. Chain-of-table: Evolving	Xu, Xiangyu Zhao, Xian Wu, Yefeng Zheng, Yang	1418
1362	tables in the reasoning chain for table understanding.	Wang, and Enhong Chen. 2024. Large language mod-	1419
1363	In <i>The Twelfth International Conference on Learning</i>	els for generative information extraction: a survey.	1420
1364	<i>Representations, ICLR 2024, Vienna, Austria, May</i>	<i>Frontiers Comput. Sci.,</i> 18(6):186357.	1421
1365	<i>7-11, 2024.</i> OpenReview.net.	Wenyi Xu, Yuren Mao, Xiaolu Zhang, Chao Zhang,	1422
1366	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	Xuemei Dong, Mengfei Zhang, and Yunjun Gao.	1423
1367	Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,	2025. Dagent: A relational database-driven data anal-	1424
1368	and Denny Zhou. 2022. Chain-of-thought prompting	ysis report generation agent. <i>CoRR</i> , abs/2503.13269.	1425
1369	elicits reasoning in large language models. In <i>Ad-</i>	Rui Yang, Boming Yang, Aosong Feng, Sixun Ouyang,	1426
1370	<i>advances in Neural Information Processing Systems 35:</i>	Moritz Blum, Tianwei She, Yuang Jiang, Freddy	1427
1371	<i>Annual Conference on Neural Information Process-</i>	Lécué, Jinghui Lu, and Irene Li. 2024. Graphusion:	1428
1372	<i>ing Systems 2022, NeurIPS 2022, New Orleans, LA,</i>	A RAG framework for knowledge graph construction	1429
1373	<i>USA, November 28 - December 9, 2022.</i>	with a global perspective. <i>CoRR</i> , abs/2410.17600.	1430
1374	Yang Wei, Honglei Guo, Jin-Mao Wei, and Zhong Su.	Yuan Yao, Deming Ye, Peng Li, Xu Han, Yankai Lin,	1431
1375	2019. Revealing semantic structures of texts: Multi-	Zhenghao Liu, Zhiyuan Liu, Lixin Huang, Jie Zhou,	1432
1376	grained framework for automatic mind-map genera-	and Maosong Sun. 2019. Docred: A large-scale	1433
1377	tion. In <i>Proceedings of the Twenty-Eighth Interna-</i>	document-level relation extraction dataset. In <i>Pro-</i>	1434
1378	<i>tional Joint Conference on Artificial Intelligence, IJ-</i>	<i>ceedings of the 57th Conference of the Association</i>	1435
1379	<i>CAI 2019, Macao, China, August 10-16, 2019,</i> pages	<i>for Computational Linguistics, ACL 2019, Florence,</i>	1436
1380	5247–5254. <i>ijcai.org.</i>	<i>Italy, July 28- August 2, 2019, Volume 1: Long Pa-</i>	1437
1381	Sam Wiseman, Stuart M. Shieber, and Alexander M.	<i>pers,</i> pages 764–777. Association for Computational	1438
1382	Rush. 2017. Challenges in data-to-document gen-	<i>Linguistics.</i>	1439
1383	eration. In <i>Proceedings of the 2017 Conference on</i>	Simge Yildirim, Yunus Santur, and Murat Aydoğan.	1440
1384	<i>Empirical Methods in Natural Language Processing,</i>	2024. Nlp in fintech: Developing a lightweight text-	1441
1385	<i>EMNLP 2017, Copenhagen, Denmark, September</i>	to-chart application for financial analysis. <i>2024 8th</i>	1442
1386	<i>9-11, 2017,</i> pages 2253–2263. Association for Com-	<i>International Artificial Intelligence and Data Pro-</i>	1443
1387	putational Linguistics.	<i>cessing Symposium (IDAP),</i> pages 1–6.	1444
1388	Xueqing Wu, Jiacheng Zhang, and Hang Li. 2022. Text-	Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021.	1445
1389	to-table: A new way of information extraction. In	Bartscore: Evaluating generated text as text genera-	1446
1390	<i>Proceedings of the 60th Annual Meeting of the As-</i>	tion. In <i>Advances in Neural Information Processing</i>	1447
1391	<i>sociation for Computational Linguistics (Volume 1:</i>	<i>Systems 34: Annual Conference on Neural Informa-</i>	1448
1392	<i>Long Papers), ACL 2022, Dublin, Ireland, May 22-27,</i>	<i>tion Processing Systems 2021, NeurIPS 2021, De-</i>	1449
1393	2022, pages 2518–2533. Association for Computa-	<i>cember 6-14, 2021, virtual,</i> pages 27263–27277.	1450
1394	tional Linguistics.	Fatemeh Pesaran Zadeh, Juyeon Kim, Jin-Hwa Kim,	1451
1395	xAI. 2025. Grok 3 beta — the age of reasoning agents.	and Gunhee Kim. 2024. Text2chart31: Instruction	1452
1396	Renqiu Xia, Bo Zhang, Hancheng Ye, Xiangchao Yan,	tuning for chart generation with automatic feedback.	1453
1397	Qi Liu, Hongbin Zhou, Zijun Chen, Min Dou, Bo-	In <i>Proceedings of the 2024 Conference on Empirical</i>	1454
1398	tian Shi, Junchi Yan, and Yu Qiao. 2024. Chartx	<i>Methods in Natural Language Processing, EMNLP</i>	1455
1399	& chartvlm: A versatile benchmark and founda-	<i>2024, Miami, FL, USA, November 12-16, 2024,</i> pages	1456
1400	tion model for complicated chart reasoning. <i>CoRR</i> ,	11459–11480. Association for Computational Lin-	1457
1401	abs/2402.12185.	guistics.	1458
1402	Shishi Xiao, Suizi Huang, Yue Lin, Yilin Ye, and Wei	Liang Zhang, Anwen Hu, Haiyang Xu, Ming Yan,	1459
1403	Zeng. 2024. Let the chart spark: Embedding seman-	Yichen Xu, Qin Jin, Ji Zhang, and Fei Huang.	1460
1404	tic context into chart with text-to-image generative	2024a. Tinychart: Efficient chart understanding with	1461
1405	model. <i>IEEE Trans. Vis. Comput. Graph.,</i> 30(1):284–	program-of-thoughts learning and visual token merg-	1462
1406	294.	ing. In <i>Proceedings of the 2024 Conference on Em-</i>	1463
		<i>pirical Methods in Natural Language Processing,</i>	1464

EMNLP 2024, Miami, FL, USA, November 12-16, 2024, pages 1882–1898. Association for Computational Linguistics.

Ruichen Zhang, Shunpu Tang, Yinqiu Liu, Dusit Niyato, Zehui Xiong, Sumei Sun, Shiwen Mao, and Zhu Han. 2025. [Toward agentic AI: generative information retrieval inspired intelligent communications and networking](#). *CoRR*, abs/2502.16866.

Songheng Zhang, Lei Wang, Toby Jia-Jun Li, Qiaomu Shen, Yixin Cao, and Yong Wang. 2024b. [Chartify-text: Automated chart generation from data-involved texts via LLM](#). *CoRR*, abs/2410.14331.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with BERT](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Weinan Zhang, Junwei Liao, Ning Li, and Kounianhua Du. 2024c. [Agentic information retrieval](#). *CoRR*, abs/2410.09713.

Yunjia Zhang, Jordan Henkel, Avriella Floratou, Joyce Cahoon, Shaleen Deep, and Jignesh M. Patel. 2024d. [Reactable: Enhancing react for table question answering](#). *Proc. VLDB Endow.*, 17(8):1981–1994.

Zhuowei Zhang, Mengting Hu, Yinhao Bai, and Zhen Zhang. 2024e. [Coreference graph guidance for mind-map generation](#). In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2024, February 20-27, 2024, Vancouver, Canada*, pages 19623–19631. AAAI Press.

Feng Zhao, Hongzhi Zou, and Cheng Yan. 2023. [Structure-aware knowledge graph-to-text generation with planning selection and similarity distinction](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 8693–8703. Association for Computational Linguistics.

Jiawen Zhu, Jinye Ran, Roy Ka-Wei Lee, Zhi Li, and Kenny T. W. Choo. 2021. [Autochart: A dataset for chart-to-text generation task](#). In *Proceedings of the International Conference on Recent Advances in Natural Language Processing (RANLP 2021), Held Online, 1-3 September, 2021*, pages 1636–1644. INCOMA Ltd.

Ziyuan Zhuang, Zhiyang Zhang, Sitao Cheng, Fangkai Yang, Jia Liu, Shujian Huang, Qingwei Lin, Saravan Rajmohan, Dongmei Zhang, and Qi Zhang. 2024. [Efficientrag: Efficient retriever for multi-hop question answering](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 3392–3411. Association for Computational Linguistics.

Appendices

A Task Definition for Graph Generation

This provides an extended definition of graph generation, where graphs can be further categorized into knowledge graphs and mind maps.

Knowledge Graph A knowledge graph is a directed graph $G_{kg} = (\mathcal{V}, \mathcal{E}, \mathcal{A}, \mathcal{L})$, where: \mathcal{V} represents entities, \mathcal{E} represents direct relationships, $\mathcal{A}(v)$ and $\mathcal{A}(e)$ are the attributes of entity types and relation types, $\mathcal{L}(v)$ and $\mathcal{L}(e)$ provide contextual meaning for nodes and edges.

Mind Map A mind map is a hierarchical directed graph $G_{mm} = (\mathcal{V}, \mathcal{E}, \mathcal{A}, \mathcal{L})$, where: \mathcal{V} represents concepts with a unique central node, \mathcal{E} represents directed governing relations, $\mathcal{A}(v)$ distinguish root and subordinate nodes, $\mathcal{L}(v)$ and $\mathcal{L}(e)$ provide contextual meaning for concepts and relationships. There is a hierarchy constraint that each non-root node has a unique path to the root (Wei et al., 2019; Hu et al., 2021).

B Detail of Current Datasets

B.1 Dataset Categorization

We use T2S to differentiate native text-to-structure datasets from those repurposed from structure-to-text tasks, as this fundamental design distinction impacts model performance evaluation. Regarding annotation practices, we consider any human verification or correction, even in crawled datasets, as constituting annotation. Schema constraints distinguish between limited schemas (domain-specific, fixed attributes like sports terminology) and unlimited schemas (open-domain, free-form relations, such as Wikipedia tables). Text inputs are classified by length into sentence, paragraph, and document scales. Reasoning complexity indicates the cognitive effort required to transform textual information into structured formats: (1) Low difficulty involves direct extraction of explicit entities/relations, populating a table with stated attributes, or creating a simple bar chart from numeric mentions; (2) Medium difficulty involves context-bound operation, such as resolving coreferences, linking “it” to a named entity in a knowledge graph, or aggregating values for chart axes; (3) High difficulty involves global reasoning, for example, temporal reasoning (e.g., constructing timelines from events), or inferring implicit hierarchies (e.g., deducing parent-child relationships in taxonomies for graph struc-

tures). The structural complexity can be systematically defined across multiple axes: (1) Tables: low complexity means less than two rows or columns; medium complexity means larger grids with uniform cell structures; high complexity means complex layouts (e.g., merged cells, nested tables, hierarchical headers); (2) Graphs: the complexity are graded by average node degree and total node count in ascending order; (3) Charts: the complexity are categorized by the diversity of chart types (e.g., bar, line, scatter) in the dataset. The Gold Test Set (GTS) serves as a critical quality indicator, distinguishing between fully human-annotated test sets and those with alternative provenance.

B.2 Common Practices for Data Construction

LLM Paraphrasing During the data generation process, leveraging LLMs for paraphrasing ensures alignment with instructional tone, enhances sentence diversity, and improves the overall data quality (Jiao et al., 2023; Deng et al., 2024).

Quality Control Various methods have been employed to ensure the quality of custom-built datasets. Some studies randomly select instances for annotation to analyze results (Gui et al., 2023; Tang et al., 2024; Deng et al., 2024). Others design iterative prompts to ensure data quality (Li et al., 2024b). Moreover, a combined approach leveraging LLMs and human collaboration has been applied for data cleaning (Shi et al., 2024).

Categorization by Difficulty Level It is a common practice that the proposed datasets are pre-divided based on subtask difficulty levels (Jiao et al., 2023; Deng et al., 2024). It allows for a more granular assessment of the model’s strengths and weaknesses across different complexity levels.

C Details of T2S Benchmark

C.1 Details of Evaluation Criteria

As mentioned in Section 6.1, we propose an evaluation framework consisting of content faithfulness (precision and recall) and structure coherence (alignment and clarity). To integrate these metrics into the context of agentic AI systems, we assess their application across four key aspects of text-to-structure generation, which is a critical process for enhancing agents’ autonomy and effectiveness. Figure 5 shows the prompting template for the evaluation of structure outputs. For the complete version of the evaluation criteria regarding tables, graphs,

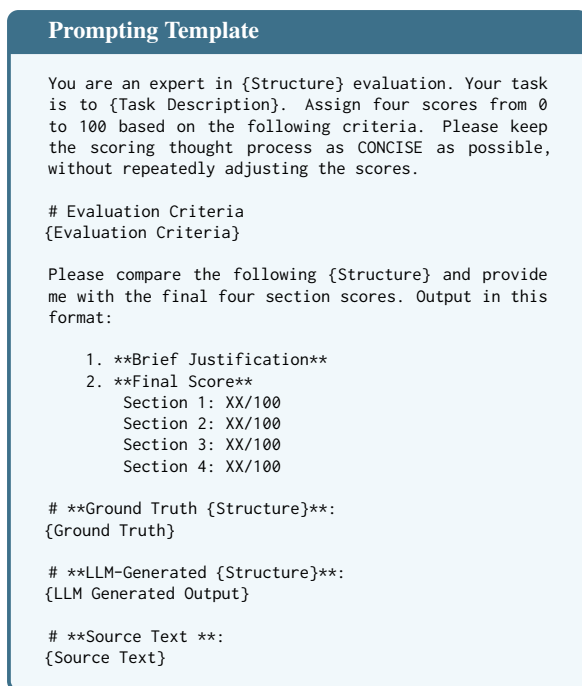


Figure 5: Prompting template for structure outputs evaluation.

and charts in the prompting template, please refer to Figure 6, 7, 8, respectively. Below, we outline how each metric applies to this process:

Precision Precision evaluates the structured outputs by identifying conflicts with the ground truth and source text, such as factual errors or mismatches. In applications like automated report generation or document extraction, precision ensures that the AI-generated tables and charts do not introduce mistakes. The detailed calculation rules for precision are outlined in the prompt – it is defined as the percentage of incorrect information relative to the total generated information, multiplied by 100 to obtain the precision score. When comparing whether two pieces of information are equivalent, numerical values, textual content, and other data types are evaluated separately.

Recall Recall evaluates the completeness of the structured outputs, assessing how thoroughly the AI captures details from the source text. It is crucial for ensuring that no critical information is omitted during the text-to-structure conversion. The detailed calculation rules for recall are detailed in the prompt – it is determined by the percentage of information from the ground truth that appears in the output, multiplied by 100 to obtain the recall score.

Alignment Alignment evaluates how well the structured outputs conform to expected formats,

such as consistent schemas or data types, ensuring seamless processing by the AI or downstream applications. The detailed calculation rules for alignment are outlined in the prompt – it is generally assessed by categorizing the degree of alignment between the output and the instruction across expected formats, data types, and other relevant dimensions into four scoring tiers.

Clarity Clarity evaluates the logical organization and presentation of the structured outputs, ensuring they are unambiguous and free of redundancy. This will enhance the efficiency of data interpretation and avoid confusion to maintain usability. The detailed calculation rules for clarity are outlined in the prompt – the outputs will be evaluated and scored across four tiers based on degree of ambiguity, comprehensibility, and redundancy.

Taking applying to agentic workflows as an example, this scoring framework optimizes the performance of text-to-structure generation: *precision* and *recall* validate extracted data fidelity, ensuring agents develop an accurate and comprehensive understanding of data, thereby facilitating the generation of reliable and holistic decisions; while *alignment* and *clarity* govern integration with downstream modules, enhancing interoperability within agentic AI systems and optimizing the efficiency of data interpretation and action.

C.2 Selection of Datasets

Regarding the selection of datasets for our Text-to-Structure benchmark, we first ensure that the chosen datasets are open-source and that their test sets are manually annotated to guarantee accuracy. For the text-to-table generation task, we select three datasets, each with distinct characteristics:

- **InstructIE** (Jiao et al., 2023): Each (text, table) pair is accompanied by a complex instruction, with a unique instruction for different instances.
- **Struc-Bench Table** (Tang et al., 2024): A cleaned version of RotoWire dataset (Wiseman et al., 2017), serving as a traditional benchmark since the beginning of text-to-table tasks.
- **LiveSum** (Deng et al., 2024): Focuses on evaluating a model’s ability to integrate information from text into tables, a capability currently lacking in many LLMs.

Evaluation Criteria for Tables

1. **Precision (100 points)**
 - Evaluates whether generated cells (with its corresponding row/column headers) conflict with ground truth. If no corresponding information exists in the ground truth, compare with the original text. Incorrect information is deducted proportionally from the total score.
 - Error Rate = (Number of conflicting cells) / (Total cells)
 - Final Score = $100 * (1 - \text{Error Rate})$
 - - **Correctness Guidelines**
 - Numerical values: Strict equality required.
 - Textual content: Semantic equivalence required.
 - Empty cells: Equivalent to “not specified”, “N/A”, or “-”.
2. **Recall (100 points)**
 - For each cell in the ground truth (with its corresponding row/column headers), evaluate the percentage of information captured in the LLM-generated output. Missing or incorrect information is deducted proportionally from the total score.
 - Rate = (Number of missing or incorrect cells) / (Total cells)
 - Final Score = $100 * (1 - \text{Rate})$
 - For **Correctness Guidelines**, please follow the same criteria as Precision.
3. **Alignment (100 points)**
 - Evaluates alignment with expected formats, schema relations, and data types (transposed rows/columns allowed without penalty):
 - Perfect Match (90-100): Rows/columns fully match reference structure.
 - Minor Deviations (80-89): Few row/column discrepancies; requires minimal adjustments.
 - Noticeable Differences (60-79): Multiple added/removed rows/columns; needs manual reorganization.
 - Severe Deviation (0-59): Missing headers/data type errors; structure unusable.
4. **Clarity (100 points)**
 - Assesses overall intelligibility of the table and clarity and redundancy in headers/cell values:
 - Perfect Clarity (90-100): No ambiguity or redundancy.
 - Minor Issues (80-89): Few redundant/ambiguous terms; core meaning intact.
 - Significant Issues (60-79): Multiple ambiguous terms requiring inference; extractable core data.
 - Critical Flaws (0-59): Conflicting headers/unreadable data; table unreliable.

Figure 6: Evaluation Criteria for Tables

Evaluation Criteria for Knowledge Graphs

1. **Precision (100 points)**
 - Evaluates whether generated triples conflict with ground truth. If no corresponding information exists in the ground truth, compare with the original text. Incorrect information is deducted proportionally from the total score.
 - Error Rate = (Number of conflicting triples) / (Total triples)
 - Final Score = $100 * (1 - \text{Error Rate})$
 - - **Correctness Guidelines**
 - Numerical values: Strict equality required.
 - Textual content: Semantic equivalence required.
2. **Recall (100 points)**
 - For each triple in the ground truth, evaluate the percentage of information captured in the LLM-generated triples. Missing or incorrect information is deducted proportionally from the total score.
 - Rate = (Number of missing or incorrect triples) / (Total triples)
 - Final Score = $100 * (1 - \text{Rate})$
 - For **Correctness Guidelines**, please follow the same criteria as Precision.
3. **Alignment (100 points)**
 - Evaluates alignment with expected formats, schema relations, and data types (transposed rows/columns allowed without penalty):
 - Perfect Match (90-100): Exact schema/format match.
 - Minor Deviations (80-89): extra/missing triples; parsable.
 - Noticeable Differences (60-79): Multiple missing triples/schema gaps; needs manual fixes.
 - Severe Deviation (0-59): Invalid predicates/scrambled data; KG unusable.
4. **Clarity (100 points)**
 - Assesses clarity, standardization, and absence of ambiguity/redundancy in entities and relations:
 - Perfect Clarity (90-100): Standardized terms; zero redundancy.
 - Minor Issues (80-89): Rare ambiguous labels; core clear.
 - Significant Issues (60-79): Frequent ambiguous terms; extractable core.
 - Critical Flaws (0-59): Uninterpretable relations; KG unreliable

Figure 7: Evaluation Criteria for Knowledge Graphs

For the text-to-graph generation task, the available datasets are limited. We select two open-domain datasets: **DART** (Nan et al., 2021) and **Instruc-**

IE (Gui et al., 2023), where the latter includes half of its data in Chinese, enabling evaluation of multilingual performance. In the text-to-chart

Evaluation Criteria for Charts

1. Precision (100 points)

- Evaluates whether generated chart elements (data points, labels, categories) conflict with the ground truth. If no ground truth chart exists, compare with the original text.
 - Error Rate = (Number of conflicting data points/labels/categories) / (Total data points/labels/categories)
 - Final Score = $100 * (1 - \text{Error Rate})$
- - **Correctness Guidelines**
 - Numerical values (e.g., axis values, percentages): Strict equality required (up to two decimal places allowed).
 - Categorical labels (e.g., axis labels, legend entries): Semantic equivalence required (e.g., "Q1 2023" vs. "First Quarter" is acceptable; "Male" vs. "Female" is not).
 - Missing elements: Treat as errors if present in the ground truth but omitted in the generated chart.

2. Recall (100 points)

- Measures how well the generated chart captures all data points, trends, and categories from the ground truth.
 - Rate = (Number of missing or misrepresented data points/categories/labels) / (Total data points/categories/labels)
 - Final Score = $100 * (1 - \text{Rate})$
- For **Correctness Guidelines**, please follow the same criteria as Precision, prioritizing critical elements (e.g., peaks, outliers, primary categories).

3. Alignment (100 points)

- Alignment with expected chart structure:
 - Perfect (90-100): Chart type, headers/labels, and groupings fully match reference.
 - Minor flaws (80-89): Formatting/positioning deviations (e.g., axis units, legend placement) with minimal impact.
 - Major errors (60-79): Incorrect chart type, missing key labels, or misaligned data groups requiring manual fixes.
 - Invalid (0-59): Broken structure (e.g., inverted axes, data-trend conflicts).

4. Clarity (100 points)

- Clarity and truthfulness:
 - Perfect (90-100): Visual encodings (color/size/position) logically match data; no ambiguity.
 - Minor issues (80-89): Non-critical redundancy/design flaws (e.g., label overlaps) with preserved meaning.
 - Confusing (60-79): Ambiguous scales/axes requiring user guesses.
 - Misleading (0-59): Distorted representations (e.g., truncated axes, false color mappings) altering insights.

Figure 8: Evaluation Criteria for Charts

generation task, due to the scarcity of dedicated datasets, we choose **Text2Chart31** (Zadeh et al., 2024), the most comprehensive dataset in terms of chart type coverage, which effectively tests the generalization capability of models. To standardize the datasets, we fix the training set sample size to 1,000 instances during SFT and sample 250 instances from each dataset’s test set for evaluation.

C.3 Selection of Models

The selected LLMs in our study can be categorized into two types: non-thinking models and thinking models. Specifically, the non-thinking models include: Mistral-7B (Jiang et al., 2023a), Llama-3.1-8B (Dubey et al., 2024), R1-Distill-Llama-8B (DeepSeek-AI et al., 2025), Phi-3-medium (Abdin et al., 2024), Mixtral-8x22B (Jiang et al., 2024a), Llama-3.1-405B (Dubey et al., 2024), Llama-3.3-70B (Dubey et al., 2024), GPT-3.5-turbo (OpenAI, 2022) and GPT-4o (OpenAI, 2024); while the thinking models comprise QwQ-32B (Qwen, 2025), DeepSeek-R1 (DeepSeek-AI et al., 2025), Grok-3-mini (xAI, 2025), and O4-mini (OpenAI, 2025). For SFT experiments, we evaluate models Mistral-7B, Llama-3.1-8B, R1-Distill-Llama-8B. The remaining models are as-

essed under a zero-shot setting. Additionally, for non-thinking models, we further examine their performance when augmented with CoT prompting. For the SFT models, our experiments are conducted on a single H20 GPU card. For other LLMs, we access their capabilities through API calls. Specifically, we conduct model inference with temperature set to 0 for non-thinking models, while adopting the officially recommended hyperparameters for thinking models as specified in their respective documentation, with all results averaged over five independent runs. When employing DeepSeek-R1 and GPT-4o as LLM-as-judge, we fix the temperature parameter at 0 to ensure that every evaluation remains deterministic.

D Result Analysis

D.1 Text-to-Table Generation

As shown in Table 7, on the InstructIE dataset, the thinking models demonstrate strong performance, with Grok-3-mini achieving the highest precision of 81.82 and R1 achieving the highest recall of 78.43; however, CoT does not produce consistent improvements across all models. Regarding coherence, while the thinking models maintain excellent

Model	InstructIE (Jiao et al., 2023)						Struc-Bench (Tang et al., 2024)						LiveSum (Deng et al., 2024)					
	Prec.	Rec.	FAITH.	Align.	Clr.	COH.	Prec.	Rec.	FAITH.	Align.	Clr.	COH.	Prec.	Rec.	FAITH.	Align.	Clr.	COH.
◆ SFT																		
Mistral-7B	<u>67.84</u>	55.93	58.08	<u>67.02</u>	<u>74.60</u>	68.85	<u>91.14</u>	86.73	87.65	<u>95.47</u>	<u>87.62</u>	90.05	53.11	53.87	53.31	<u>99.81</u>	<u>97.78</u>	98.54
Llama-3.1-8B	<u>64.38</u>	<u>60.34</u>	59.31	<u>64.61</u>	<u>72.28</u>	66.44	<u>87.02</u>	<u>83.94</u>	85.01	<u>90.38</u>	<u>82.67</u>	85.95	<u>61.69</u>	<u>61.46</u>	61.34	<u>90.56</u>	<u>90.65</u>	90.30
R1-Distill-Llama-8B	<u>70.39</u>	<u>61.44</u>	62.86	<u>64.69</u>	<u>75.12</u>	68.08	<u>63.97</u>	<u>54.84</u>	54.56	<u>62.55</u>	<u>57.34</u>	58.97	<u>52.93</u>	<u>43.21</u>	43.27	<u>86.28</u>	<u>86.44</u>	86.22
◆ Zero-Shot																		
Phi-3-medium	43.19	28.79	29.15	40.42	42.75	39.89	<u>74.68</u>	<u>75.36</u>	74.04	<u>91.79</u>	<u>78.68</u>	83.99	41.84	42.27	41.61	<u>93.64</u>	<u>90.62</u>	91.73
Phi-3-medium (CoT)	42.85	29.82	30.36	37.83	39.56	37.23	<u>74.91</u>	<u>73.83</u>	73.42	<u>88.81</u>	<u>75.92</u>	81.06	46.57	46.59	46.36	<u>94.40</u>	<u>92.54</u>	93.19
Mistral-7B	58.53	53.29	53.00	57.05	<u>65.53</u>	59.73	55.17	<u>61.17</u>	56.19	<u>72.52</u>	<u>64.05</u>	67.05	31.25	30.56	30.14	<u>90.48</u>	85.78	87.26
Mistral-7B (CoT)	59.60	53.82	54.02	57.30	<u>66.83</u>	60.37	54.35	<u>62.38</u>	56.25	<u>76.51</u>	<u>66.44</u>	70.34	42.38	42.14	41.27	89.75	83.43	85.97
Mixtral-8x22B	<u>66.92</u>	<u>54.34</u>	56.35	<u>66.39</u>	<u>73.42</u>	68.25	84.95	<u>84.71</u>	84.16	89.00	<u>82.74</u>	85.20	49.74	48.97	49.12	<u>87.54</u>	88.91	87.91
Mixtral-8x22B (CoT)	<u>66.82</u>	56.05	58.27	<u>64.97</u>	<u>74.73</u>	68.29	82.81	<u>80.26</u>	80.55	<u>74.98</u>	<u>75.91</u>	74.94	55.68	54.43	54.87	<u>86.75</u>	88.51	87.37
Llama-3.1-405B	<u>71.31</u>	<u>62.27</u>	63.88	<u>65.36</u>	<u>77.31</u>	69.63	<u>90.20</u>	<u>89.32</u>	89.43	<u>91.57</u>	<u>90.49</u>	90.43	<u>62.46</u>	<u>62.38</u>	62.29	<u>99.55</u>	<u>98.94</u>	99.11
Llama-3.1-405B (CoT)	<u>72.43</u>	<u>62.39</u>	64.54	<u>69.43</u>	<u>78.69</u>	72.80	<u>90.69</u>	<u>88.69</u>	89.38	<u>91.72</u>	<u>90.51</u>	90.74	<u>66.27</u>	<u>66.33</u>	66.25	<u>98.87</u>	<u>98.51</u>	98.64
Llama-3.3-70B	<u>71.10</u>	<u>60.09</u>	62.88	<u>64.63</u>	<u>76.50</u>	69.19	82.00	<u>80.98</u>	80.58	<u>84.68</u>	<u>77.58</u>	80.00	58.55	58.63	58.48	<u>99.40</u>	<u>99.13</u>	99.19
Llama-3.3-70B (CoT)	<u>71.71</u>	<u>63.35</u>	65.34	<u>64.79</u>	<u>77.99</u>	69.46	87.53	<u>86.47</u>	86.47	<u>87.16</u>	<u>85.71</u>	85.84	<u>62.66</u>	<u>63.00</u>	62.79	<u>98.98</u>	<u>98.47</u>	98.68
GPT-3.5-turbo	<u>65.29</u>	53.80	55.09	<u>67.05</u>	<u>74.99</u>	69.01	79.60	<u>79.61</u>	78.44	<u>92.09</u>	<u>81.98</u>	85.37	46.08	45.39	45.38	<u>77.44</u>	<u>82.92</u>	79.18
GPT-3.5-turbo (CoT)	<u>68.95</u>	55.61	58.44	<u>65.69</u>	<u>75.21</u>	68.80	<u>72.86</u>	<u>75.03</u>	71.88	<u>89.55</u>	<u>76.30</u>	81.54	45.14	45.57	45.09	<u>85.88</u>	<u>87.71</u>	86.37
GPT-4o	<u>75.17</u>	<u>71.67</u>	70.72	<u>72.97</u>	<u>82.35</u>	76.36	88.00	<u>87.44</u>	87.35	<u>96.34</u>	<u>89.22</u>	92.22	<u>63.46</u>	<u>63.09</u>	63.21	<u>91.46</u>	<u>92.96</u>	91.99
GPT-4o (CoT)	<u>75.91</u>	<u>67.93</u>	68.11	<u>71.97</u>	<u>80.15</u>	74.40	87.37	<u>85.46</u>	86.05	<u>94.29</u>	<u>87.24</u>	90.03	<u>61.96</u>	<u>62.88</u>	62.16	<u>98.76</u>	<u>98.38</u>	98.50
QwQ-32B	<u>75.53</u>	<u>66.66</u>	68.23	<u>63.08</u>	<u>77.90</u>	68.79	81.15	<u>81.29</u>	81.21	<u>90.97</u>	<u>84.03</u>	87.12	58.72	58.47	58.06	<u>95.94</u>	<u>96.95</u>	96.20
DeepSeek-R1	<u>79.05</u>	<u>78.72</u>	76.86	<u>70.04</u>	<u>82.93</u>	74.75	88.19	<u>87.95</u>	87.77	<u>94.79</u>	<u>89.32</u>	91.05	<u>64.37</u>	<u>64.79</u>	64.89	<u>96.74</u>	<u>97.88</u>	96.85
Grok-3-mini	<u>81.34</u>	<u>72.42</u>	75.57	<u>68.93</u>	<u>82.91</u>	75.40	86.17	<u>85.45</u>	85.60	<u>94.61</u>	<u>89.59</u>	92.28	<u>65.81</u>	<u>65.83</u>	66.35	<u>96.45</u>	<u>94.27</u>	95.26
O4-mini	<u>80.86</u>	<u>73.26</u>	75.71	<u>68.73</u>	<u>83.57</u>	74.85	84.13	<u>84.83</u>	83.45	<u>94.52</u>	<u>86.30</u>	90.02	<u>84.20</u>	<u>84.43</u>	84.48	<u>98.64</u>	<u>98.65</u>	99.16

Table 7: Performance of different models on three text-to-table datasets, with the highest scores in each metric underlined.

Model	DART (Nan et al., 2021)						InstructIE (Gui et al., 2023)					
	Prec.	Rec.	FAITH.	Align.	Clr.	COH.	Prec.	Rec.	FAITH.	Align.	Clr.	COH.
◆ SFT												
Mistral-7B	<u>80.99</u>	<u>70.81</u>	75.07	<u>83.15</u>	<u>83.54</u>	83.17	59.75	45.35	44.58	<u>72.21</u>	<u>69.31</u>	69.67
Llama-3.1-8B	<u>63.63</u>	35.92	38.94	55.56	56.75	55.13	55.37	35.73	34.93	<u>66.42</u>	<u>65.70</u>	64.99
R1-Distill-Llama-8B	<u>68.98</u>	54.83	59.03	<u>68.41</u>	<u>71.40</u>	69.43	<u>71.78</u>	33.35	37.27	<u>70.89</u>	<u>69.91</u>	69.73
◆ Zero-Shot												
Phi-3-medium	<u>73.80</u>	<u>60.80</u>	64.60	<u>72.36</u>	<u>73.26</u>	72.59	<u>74.64</u>	39.84	43.40	<u>74.08</u>	<u>71.68</u>	71.44
Phi-3-medium (CoT)	<u>73.69</u>	<u>61.10</u>	64.88	<u>71.77</u>	<u>73.52</u>	72.32	<u>74.54</u>	39.10	41.28	<u>74.39</u>	<u>70.40</u>	70.44
Mistral-7B	<u>74.22</u>	<u>69.81</u>	70.27	<u>73.34</u>	<u>73.44</u>	73.06	<u>65.88</u>	41.21	42.74	<u>69.34</u>	<u>66.85</u>	66.54
Mistral-7B (CoT)	<u>76.98</u>	59.92	64.68	<u>70.70</u>	<u>71.96</u>	70.96	<u>66.78</u>	30.86	33.35	<u>66.36</u>	<u>67.40</u>	66.01
Mixtral-8x22B	<u>70.58</u>	<u>65.30</u>	65.70	<u>71.72</u>	<u>71.07</u>	71.07	<u>68.69</u>	38.65	40.07	<u>71.32</u>	<u>68.57</u>	68.89
Mixtral-8x22B (CoT)	<u>78.14</u>	<u>69.19</u>	71.27	<u>74.93</u>	<u>74.43</u>	74.53	<u>75.68</u>	43.30	47.47	<u>72.92</u>	<u>71.31</u>	71.29
Llama-3.1-405B	<u>75.77</u>	<u>65.24</u>	68.74	<u>75.59</u>	<u>76.11</u>	75.61	<u>73.79</u>	43.30	46.72	<u>78.16</u>	<u>75.52</u>	75.82
Llama-3.1-405B (CoT)	<u>83.56</u>	<u>66.49</u>	72.31	<u>77.99</u>	<u>78.57</u>	78.09	<u>79.32</u>	43.46	47.94	<u>79.90</u>	<u>78.10</u>	78.48
Llama-3.3-70B	<u>79.22</u>	<u>66.95</u>	71.09	<u>78.34</u>	<u>77.21</u>	77.54	<u>77.23</u>	44.64	47.11	<u>78.70</u>	<u>76.36</u>	76.30
Llama-3.3-70B (CoT)	<u>85.74</u>	<u>71.69</u>	77.11	<u>79.02</u>	<u>79.24</u>	78.99	<u>76.47</u>	41.80	46.56	<u>80.23</u>	<u>78.08</u>	78.51
GPT-3.5-turbo	<u>80.95</u>	<u>64.02</u>	69.75	<u>76.56</u>	<u>77.27</u>	76.72	<u>74.63</u>	38.16	41.89	<u>75.66</u>	<u>73.32</u>	73.34
GPT-3.5-turbo (CoT)	<u>79.76</u>	<u>65.20</u>	70.17	<u>76.20</u>	<u>76.76</u>	75.99	<u>74.78</u>	37.01	41.78	<u>75.08</u>	<u>73.40</u>	73.51
GPT-4o	<u>80.51</u>	<u>71.42</u>	74.15	<u>77.51</u>	<u>77.33</u>	77.18	<u>76.02</u>	43.13	46.77	<u>78.81</u>	<u>76.85</u>	77.00
GPT-4o (CoT)	<u>87.02</u>	<u>76.47</u>	79.89	<u>79.44</u>	<u>79.45</u>	79.33	<u>81.21</u>	49.30	53.28	<u>79.38</u>	<u>77.72</u>	77.65
QwQ-32B	<u>76.96</u>	<u>65.79</u>	69.42	<u>78.64</u>	<u>78.83</u>	78.02	<u>86.42</u>	50.02	56.30	<u>80.67</u>	<u>82.82</u>	82.15
DeepSeek-R1	<u>84.67</u>	<u>74.82</u>	78.25	<u>80.04</u>	<u>79.42</u>	79.89	<u>80.93</u>	47.77	52.15	<u>80.27</u>	<u>80.20</u>	79.23
Grok-3-mini	<u>82.41</u>	<u>68.94</u>	73.32	<u>79.11</u>	<u>80.38</u>	80.34	<u>87.66</u>	<u>50.06</u>	56.03	<u>81.18</u>	<u>83.90</u>	82.68
O4-mini	<u>84.03</u>	<u>71.06</u>	74.82	<u>79.28</u>	<u>80.52</u>	80.19	<u>89.89</u>	49.91	56.46	<u>81.51</u>	<u>83.33</u>	82.32

Table 8: Performance of different models on two text-to-KG datasets, with the highest scores in each metric underlined.

performance, GPT-4o achieves the highest score. Among SFT models, the R1-distilled Llama-8B exhibits the best overall performance, although it still falls short of the thinking models’ capabilities. On the Struc-Bench dataset, most LLMs show strong baseline performance since the tasks primarily involve information replication rather than complex

reasoning. In terms of specific metrics, Mistral-7B (SFT) achieves the highest precision score, while Llama-3.1-405B achieves the highest recall score. For coherence evaluation, GPT-4o exhibits superior capabilities. Notably, the Mistral-7B (SFT) model produces great performance gains, with improvements of 56% in faithfulness and 34% in coher-

ence metrics compared to zero-shot approaches, empirically demonstrating that SFT can significantly enhance model performance on information-replication text-to-table tasks. On the Livesum dataset, both precision and recall scores are generally low due to the complex information integration and reasoning operations required by the task. O4-mini achieves exceptional performance with a faithfulness score of 83.99. Regarding coherence, all models attain relatively high scores due to the dataset’s standardized table format, with Llama-3.3-70B demonstrating the best performance. The results indicate that thinking models maintain a competitive advantage on this dataset, while among SFT models, Llama-3.1-8B achieves comparable performance to larger LLMs.

D.2 Text-to-Graph Generation

As shown in Table 8, our experimental results demonstrate that GPT-4o (CoT) achieves the highest faithfulness scores on the DART dataset, while Mistral-7B (SFT) shows superior coherence performance that significantly outperforms all other LLMs. The application of CoT prompting generally improves precision and recall metrics by 6-10% across most LLMs, though it does not yield significant improvement in coherence metrics. On the InstructIE dataset, inherent dataset characteristics lead to lower recall scores overall, but QwQ-32B maintains relatively better performance in this metric, with O4-mini achieving the highest precision score of 89.45. The coherence metrics show remarkable consistency with the patterns observed on the DART dataset, suggesting stable model behavior across different evaluation benchmarks.

D.3 Text-to-Chart Generation

As shown in Table 9, in addition to the aforementioned metrics, we also evaluate the success rate of the generated Python code executing correctly to produce the desired images, denoted as SUCC. It is observed that the majority of LLMs achieve a success rate of over 95%. Among them, Grok-3-mini demonstrates the highest success rate, reaching 99.2%. In terms of precision and recall, most LLMs achieve scores above 80, with Llama-3.3-70B (CoT) performing the highest, scoring 88.18 and 90.16, respectively. For alignment and clarity, Grok-3-mini exhibits the best performance, achieving scores of 88.34 and 89.65, respectively. Among the SFT models, Llama-3.1-8B demonstrated the highest overall performance. Mistral-7B achieves

Model	Text2Chart31 (Zadeh et al., 2024)						
	SUCC.	Prec.	Rec.	FAITH.	Align.	Clr.	COH.
◆ <i>SFT</i>							
Mistral-7B	72.8%	64.04	65.66	64.54	63.38	63.70	63.41
Llama-3.1-8B	87.2%	76.25	79.16	77.06	77.07	77.70	77.25
R1-Distill-Llama-8B	47.6%	41.36	42.46	41.66	39.97	41.11	40.45
◆ <i>Zero-Shot</i>							
Phi-3-medium	68.8%	59.62	62.20	60.24	60.26	60.89	60.45
Phi-3-medium (CoT)	67.6%	56.93	59.05	57.75	58.17	58.84	58.27
Mistral-7B	64.0%	51.72	53.84	52.25	52.03	52.70	52.03
Mistral-7B (CoT)	59.2%	44.78	46.38	45.13	46.90	47.00	46.51
Mixtral-8x22B	72.4%	60.89	63.05	61.62	62.09	62.51	62.15
Mixtral-8x22B (CoT)	67.6%	58.23	60.49	58.92	58.61	59.37	58.89
Llama-3.1-405B	98.0%	85.64	88.00	86.49	86.47	87.26	86.62
Llama-3.1-405B (CoT)	98.4%	86.39	88.71	87.23	86.40	87.15	86.40
Llama-3.3-70B	98.4%	86.77	89.37	87.44	86.54	88.00	87.12
Llama-3.3-70B (CoT)	98.8%	88.18	90.16	88.78	87.19	88.80	87.83
GPT-3.5-turbo	94.4%	84.53	86.23	85.17	84.20	85.63	84.79
GPT-3.5-turbo (CoT)	93.2%	82.48	84.27	83.06	83.63	84.20	83.62
GPT-4o	95.6%	85.16	86.35	85.48	85.91	86.83	86.01
GPT-4o (CoT)	96.0%	85.48	86.82	85.99	86.24	87.26	86.63
QwQ-32B	98.4%	86.28	88.46	86.71	86.61	88.04	86.94
DeepSeek-R1	95.6%	85.36	87.16	85.79	85.71	87.64	86.56
Grok-3-mini	99.2%	87.44	89.19	88.05	88.34	89.65	88.74
O4-mini	95.6%	84.92	87.30	85.56	85.84	87.16	86.17

Table 9: Performance of different models on text-to-chart dataset, with the highest scores in each metric underlined.

significant improvements compared to its zero-shot counterpart, with a 13.8% increase in success rate, 23.5% and 21.9% enhancements in faithfulness and coherence scores, respectively, showing the effectiveness of SFT in model optimization.

E Case Study

This section illustrates five representative failure cases of traditional evaluation metrics through systematic analysis. We empirically demonstrate scenarios where traditional metrics (BERTScore, ROUGE, etc.) contradict human judgment while our framework correctly identifies superior outputs. Each case dissects specific error types and metric failures, providing granular evidence for our approach’s validity across diverse tasks. To reflect the scoring process of *Content Faithfulness* and *Structure Coherence*, we append the intermediate reasoning steps from the LLM as a brief justification following the output for each case.

Case 1 (Figure 9) It is notable that the table generated by GPT-3.5 lacks two critical rows of information, and its single extracted row contains multiple errors compared to the ground truth. Conversely, the table produced by R1 contains comprehensive information, fully capturing all details present in the ground truth while additionally extracting more complete information. Despite R1’s table being demonstrably superior in this context,

both BERTScore and ROUGE metrics indicate GPT-3.5’s output is preferable.

Case 2 (Figure 10) The table generated by GPT-3.5 introduced erroneous information in several cells (e.g., misrepresenting reactants and products). In contrast, R1’s extracted table maintained complete accuracy in all key information. Paradoxically, BERTScore and ROUGE metrics significantly favor the former.

Case 3 (Figure 11) GPT-3.5’s table exhibits numerical inaccuracies across multiple cells, with an error rate 60% higher than R1’s table. Nevertheless, BERTScore metrics consistently rate GPT-3.5’s output as superior.

Case 4 (Figure 12) GPT-3.5’s edge information contains instances where predicates and objects are entirely confused, accompanied by structural disarray. Conversely, R1 extracts a structurally correct graph. Contradictorily, both BERTScore and ROUGE metrics indicate GPT-3.5’s graph extraction performed better.

Case 5 (Figure 13) GPT-3.5’s generated chart contains numerical errors and fails to reflect logical pairing in its layout. R1’s chart, however, maintains numerical accuracy and replicates the ground truth’s logical structure. Despite this qualitative advantage, METEOR and CodeBLEU metrics favor GPT-3.5’s output.

Case 1: Instance #068 from InstructIE

Comparison of Evaluation Metrics

Output	BERTScore	ROUGE	Content Faithfulness	Structure Coherence
GPT-3.5-turbo	0.8602	0.7052	40.00	64.62
DeepSeek-R1	0.8364	0.5974	75.00	80.00

Ground Truth

Number of Bedrooms	Number of Bathrooms	Special Rooms/Spaces	Floor Material	Outdoor Spaces
3	2.5	Breakfast nook, playroom	Hardwood	Fenced backyard, deck, patio
3	N/A	Den, indoor rock garden	New flooring	Rear balconies
5	3.5	Family room, playroom, laundry	Wood	Large back patio, storage shed, fenced yard

GPT-3.5-turbo Output

Number of Bedrooms	Number of Bathrooms	Special Rooms/Spaces	Floor Material	Outdoor Spaces
3	2.5	Playroom, Master Suite, Den	Hardwood floors, New flooring	Deck, Patio, Fenced backyard, Front patio area

Brief Justification

- **Precision**: 3/5 cells conflicting (Special Rooms, Floor Material, Outdoor Spaces) → 40.
- **Recall**: 3/5 ground truth cells missing/incorrect (Breakfast nook, Hardwood, exact Outdoor Spaces) → 40.
- **Alignment**: Missing two rows (ground truth has 3 rows, LLM has 1) → Noticeable Differences (70).
- **Clarity**: Merged data from multiple sources into one row creates ambiguity → Significant Issues (60).

DeepSeek-R1 Output

#	Bedrooms	Bathrooms	Special Rooms/Spaces	Floor Material	Outdoor Spaces
1	3	2.5	Pantry, kitchen island, breakfast nook, master suite (walk-in closets), playroom, laundry area, garage/workshop	Hardwood	Fenced backyard, deck, patio, front patio
2	3	N/A	Two-story living room, conversation pit, indoor rock garden, den, updated kitchen	New flooring	Rear balconies, golf course views, lake views
3	5	3.5	Kitchen island, pantry, laundry room, primary suite (walk-in closet), finished basement (family room, playroom)	Wood floors	Back patio, storage shed with loft, fenced yard

Brief Justification

- **Precision**: LLM-added details (e.g., extra rooms/outdoor features) not in GT create 6/15 conflicts (60).
- **Recall**: All GT cells are captured despite extra LLM additions (100).
- **Alignment**: Extra row numbers and header wording deviations (80).
- **Clarity**: Clear headers but redundant cell details (80).

Figure 9: Misalignment of Traditional Metrics and Our Framework: Case 1

Case 2: Instance #127 from InstructIE

Comparison of Evaluation Metrics

Output	BERTScore	ROUGE	Content Faithfulness	Structure Coherence
GPT-3.5-turbo	0.9229	0.9012	61.46	67.41
DeepSeek-R1	0.7231	0.4086	100.00	76.77

Ground Truth

Reactants	Solvent	Product	Yield
CyPBn-Cy, NiCl ₂ (DME)	THF	(CyPBn-Cy)NiCl ₂	68%
(CyPBn-Cy)NiCl ₂ , (o-tol)MgCl	THF	(CyPBn-Cy)Ni(o-tol)Cl	95%

GPT-3.5-turbo Output

Reactants	Products	Solvent	Yield	Analysis Method
CyPBn-Cy, NiCl ₂ (DME), THF	(CyPBn-Cy)NiCl ₂ , dichloromethane	THF	68%	Elemental, NMR, X-ray
(CyPBn-Cy)NiCl ₂ , (o-tol)MgCl, THF	(CyPBn-Cy)Ni(o-tol)Cl	THF	95%	—

Brief Justification

- **Precision**: 4/10 cells incorrect (extra "dichloromethane" in Product, Analysis Method column invalid, THF mislabeled as Reactant).
- **Recall**: 3/8 GT cells missing/incorrect (Reactants misrepresented in both rows, Product error in row 1).
- **Alignment**: Added Analysis Method column, Products vs Product header mismatch (Noticeable Differences).
- **Clarity**: Ambiguous Reactants (THF included redundantly), redundant Analysis Method (Significant Issues).

DeepSeek-R1 Output

Reaction Equation	Reactants	Solvent	Workup/Isolation	Yield	Analytical Data
eq1	CyPBn-Cy, NiCl ₂ (DME)	THF	Dichloromethane solvate	68%	Elemental analysis, NMR, X-ray crystallography
eq2	(CyPBn-Cy)NiCl ₂ , (o-tol)MgCl	THF	Isolated as pure solid	95%	Elemental analysis

Brief Justification

- **Precision**: All corresponding cells (Reactants, Solvent, Product, Yield) match GT and source text. Extra columns (Workup/Isolation, Analytical Data) align with source and do not conflict.
- **Recall**: All GT cells are captured in the generated table.
- **Alignment**: Added columns (Reaction Equation, Workup/Isolation, Analytical Data) deviate from GT schema.
- **Clarity**: Clear headers and data, but redundant columns introduce minor deviations from GT focus.

Figure 10: Misalignment of Traditional Metrics and Our Framework: Case 2

Case 3: Instance #42842344 from LiveSum

Comparison of Evaluation Metrics

Output	BERTScore	ROUGE	Content Faithfulness	Structure Coherence
GPT-3.5-turbo	0.9303	0.7750	50.00	100.00
DeepSeek-R1	0.9294	0.7750	69.00	100.00

Ground Truth

Team	Goals	Shots	Fouls	Yellow Cards	Red Cards	Corner Kicks	Free Kicks	Offsides
Home Team	5	15	8	2	0	4	8	1
Away Team	1	10	8	0	0	6	8	3

GPT-3.5-turbo Output

	Goals	Shots	Fouls	Yellow Cards	Red Cards	Corner Kicks	Free Kicks	Offsides
Home Team	5	12	6	1	0	5	4	1
Away Team	1	8	7	0	0	4	6	3

Brief Justification

- **Precision**: 8/16 cells have errors (numerical mismatches), leading to 50% error rate.
- **Recall**: Same 8 missing/correct GT cells, 50% recall.
- **Alignment**: Columns/headers match; row order irrelevant per criteria (no penalty).
- **Clarity**: Headers clear, data unambiguous despite errors.

DeepSeek-R1 Output

Stat	Home Team	Away Team
Goals	5	1
Shots	12	8
Fouls	8	12
Yellow Cards	2	0
Red Cards	0	0
Corner Kicks	5	5
Free Kicks	12	8
Offsides	1	3

Brief Justification

- **Precision/Recall**: 5 cell errors (Shots, Fouls, Corner Kicks, Free Kicks).
- **Alignment**: Transposed correctly without penalties.
- **Clarity**: Clear headers and layout despite data inaccuracies.

Figure 11: Misalignment of Traditional Metrics and Our Framework: Case 3

Case 4: Instance #048 from DART

Comparison of Evaluation Metrics

Output	BERTScore	ROUGE	Content Faithfulness	Structure Coherence
GPT-3.5-turbo	0.7165	0.6465	40.48	74.67
DeepSeek-R1	0.7052	0.6166	100.00	80.00

Ground Truth

Subject	Predicate	Object
The Vaults	eatType	pub
The Vaults	food	Japanese
The Vaults	priceRange	£20–25
The Vaults	customer rating	high
The Vaults	area	city centre
The Vaults	familyFriendly	yes
The Vaults	near	Rainbow Vegetarian Café

GPT-3.5-turbo Output

Subject	Predicate	Object
The Vaults	has price range	£20–25
The Vaults	is child friendly	
The Vaults	is a Japanese pub	
The Vaults	is near	Rainbow Vegetarian Café
The Vaults	is located in	city centre
The Vaults	is highly rated	

Brief Justification

- **Precision**: 2/6 triples conflict (missing "yes" and "high" values).
- **Recall**: 5/7 ground truth triples missing/incorrect (combined or incomplete).
- **Structural**: Missing objects in 2 triples and merged triples, parsable but flawed.
- **Semantic**: Non-standard predicates (e.g., "is child friendly") but core meaning preserved.

DeepSeek-R1 Output

Subject	Relation	Object
The Vaults	is	child friendly
The Vaults	is a	highly rated Japanese pub
The Vaults	has price range	£20–25
The Vaults	is located in	city centre
The Vaults	is near	Rainbow Vegetarian Café

Brief Justification

All ground truth information is captured without conflicts, but relations use non-standard phrasing (e.g., "is a" vs "eatType"). Structural/semantic scores reflect schema deviations and non-standard predicates.

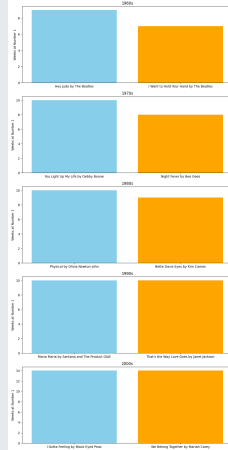
Figure 12: Misalignment of Traditional Metrics and Our Framework: Case 4

Case 5: Instance #0208113110 from Text2Chart31

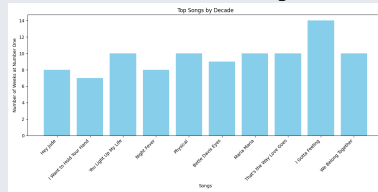
Comparison of Evaluation Metrics

Output	METEOR	CodeBLEU	Content Faithfulness	Structure Coherence
GPT-3.5-turbo	0.4480	0.3701	74.67	69.64
DeepSeek-R1	0.3148	0.2263	100.00	72.41

Ground Truth



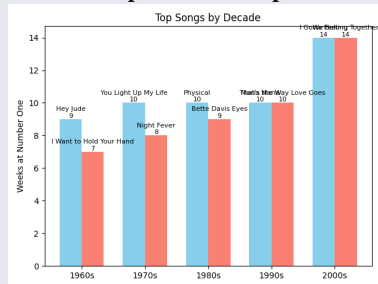
GPT-3.5-turbo Output



Brief Justification

- **Precision**: The LLM-generated chart has discrepancies in the number of weeks for some songs compared to the ground truth. For example, "Hey Jude" is shown with 8 weeks instead of 9, and "I Gotta Feeling" with 14 weeks instead of 10.
- **Recall**: The generated chart includes all songs and their respective weeks, but the values are not always accurate, affecting recall.
- **Alignment**: The chart type is consistent, but the LLM-generated chart lacks the decade grouping present in the ground truth, which is a significant structural difference.
- **Clarity**: The generated chart is clear but lacks the decade separation, which could lead to confusion about the time periods of the songs.

DeepSeek-R1 Output



Brief Justification

- **Precision/Recall**: All correct.
- **Alignment**: The chart type is consistent (bar chart), but the structure differs. The ground truth uses separate subplots for each decade, while the generated chart combines them into one. This structural difference impacts the score.
- **Clarity**: The generated chart is clear and logical in its presentation, but the overlapping labels and combined decades create some confusion, affecting the semantic score.

Figure 13: Misalignment of Traditional Metrics and Our Framework: Case 5