Learning to Reason for Text Generation from Scientific Tables

Nafise Sadat Moosavi¹, Andreas Rücklé¹, Dan Roth², Iryna Gurevych¹

¹UKP Lab, Technische Universitat Darmstadt ²Department of Computer and Information Science, UPenn

¹https://www.ukp.tu-darmstadt.de
² https://www.seas.upenn.edu/~darroth/

Abstract

In this paper, we introduce SciGen, a new challenge dataset for the task of reasoningaware data-to-text generation consisting of tables from scientific articles and their corresponding descriptions. Describing scientific tables goes beyond the surface realization of the table content and requires reasoning over table values. The unique properties of Sci-Gen are that (1) tables mostly contain numerical values, and (2) the corresponding descriptions require arithmetic reasoning. SciGen is therefore the first dataset that assesses the arithmetic reasoning capabilities of generation models on complex input structures, i.e., tables from scientific articles. We study the effectiveness of state-of-the-art data-to-text generation models on SciGen and evaluate the results using common metrics as well as human evaluation. Our results and analyses show that (a) while humans like to reason for describing scientific tables, the ability of stateof-the-art models is severely limited on this task, (b) while adding more training data improves the results, it is not the solution for reasoning-aware text generation, and (c) one of the main bottlenecks for this task is the lack of proper automatic evaluation metrics. The data, code, and annotations for human evaluation will be available at https://github. com/UKPLab/SciGen. SciGen opens new avenues for future research in reasoning-aware text generation and evaluation.

1 Introduction

Data-to-text generation is one of the established tasks in NLP in which the input is structured data like tables and the output is a text that describes the data. In this paper, we introduce a new data-totext generation dataset, called SciGen, that contains pairs of scientific tables and their corresponding descriptions. The descriptions of a large number of tables in the scientific articles of the computer science domain require one or more types of arithmetic reasoning over table values—e.g., argMax, argMin, comparison, subtraction, etc—over table values, indicating that humans like to reason for describing scientific tables.¹ Therefore, generative models that can describe scientific tables should be able to perform arithmetic reasoning. For instance, such models can be used to automate scientific writing by generating the result section of the article using the table of experimental results, or for developing specialized chatbots that can generate answers based on the content of tables.

In this direction, we have created SciGen to enable the development and evaluation of generation models with arithmetic reasoning capabilities. For creating SciGen, we select tables and their corresponding descriptions from the computer science articles with the following properties: (1) the tables mostly contain numerical values, and (2) the corresponding descriptions are the result of arithmetic reasoning over table values.

We release the dataset in three settings based on the size of the training data. The *few-shot* setting contains table-description pairs that are annotated by experts. Since expert annotation is costly and does not scale to large data sizes, we introduce automatic pairing and pruning techniques to extract high-quality table-description pairs in an unsupervised way. We extend the expert annotated pairs in the few-shot setting using the automatically extracted pairs to create the medium and large splits of the dataset. The test data is the same in all the settings and contains expert-annotated pairs.

We study state-of-the-art data-to-text generation models—including BART (Lewis et al., 2020) and T5 (Raffel et al., 2020) pretrained language models—on SciGen, and we evaluate the results

¹For instance, we randomly select 30 tables from different articles in this domain, and 60% of these tables were described by using arithmetic reasoning.

using common evaluation metrics as well as human evaluation. Our results show that (1) none of the common metrics can properly discriminate the outputs of reasoning-aware text generation, and the lack of a proper automatic evaluation metric is the main bottleneck for the progress for this task, (2) while the generated outputs by the examined models are coherent and fluent and look like valid descriptions, they are mostly factually incorrect given the content of the table, i.e., the entities that appear in the description are relevant to the table but the lack arithmetic reasoning capabilities results in incorrect statements about the present results in the table, and (3) the addition of automatically extracted training data in the medium and large settings reduces the hallucination and improves the correctness of the generated descriptions, however, the extent that they can improve the correctness is limited, e.g., the correctness of the outputs in the large setting is lower than those in medium.

Overall, the main contributions of this paper are:

- We release SciGen, the first reasoning-aware data-to-text generation dataset based on scientific articles.
- We provide our code for extracting tabledescription pairs from latex files of scientific articles that (1) provides high-quality unsupervised training data, and (2) facilitates future annotation studies for new domains.
- We provide all the data and annotations that is used for our human evaluations to help the creation of new metrics for this task as our results suggests that existing automatic metrics are not discriminative for this task.
- We show that text generation using arithmetic reasoning is a challenging task that cannot be solved by merely adding more data or applying large pre-trained models.

2 Related Work

The task of data-to-text generation is to generate coherent, relevant, and meaningful natural language text that describes the non-linguistic input data like tables, knowledge bases, tuples, or graphs (Reiter and Dale, 2000; Gatt and Krahmer, 2018). Existing datasets for data-to-text generation cover various domains and applications including sport reports (Wiseman et al., 2017; Thomson et al., 2020; van der Lee et al., 2017), weather reports or forecast (Belz, 2008; Balakrishnan et al., 2019), restaurant descriptions (Dušek et al., 2020; Oraby et al., 2018; Reed et al., 2018), biographies (Lebret et al., 2016; Nema et al., 2018), entity descriptions (Qader et al., 2018; Wang et al., 2018), as well as open-domain datasets (Gardent et al., 2017; Parikh et al., 2020).

The textual descriptions in the majority of existing datasets mostly contain a verbalized summary of the content in the data, and is therefore a surface-level summary of the data (Chen and Mooney, 2008; Belz et al., 2011; Lebret et al., 2016; Gardent et al., 2017; Dušek et al., 2018; Koncel-Kedziorski et al., 2019; Radev et al., 2020; Parikh et al., 2020). SciGen, on the other hand, goes beyond the surface realization of the input data and requires arithmetic reasoning for text generation. The most related dataset to SciGen is LogicNLG (Chen et al., 2020a), in which the text generation step also requires logical reasoning.

LogicNLG is based on TabFact (Chen et al., 2020b) that is a table-based fact verification datasets. For creating TabFact, annotators were asked to write refute and entailment statements based on Wikipedia tables. The statements were classified into simple and complex: simple statements are verifiable without involving logical inference, and complex statements involve multiple rows of the tables as well as logical operations such as summary, argMax, argMin, count, comparison, average, etc. LogicNLG contains the complex statements of TabFact that are labeled as entailment given their corresponding table.

Apart from their domains-i.e., Wikipedia vs. scientific texts, there are two main differences between LogicNLG and SciGen. First, annotators of TabFact were asked to generate multiple statements per table. As a result, each text only describes a part of the table-i.e., on average two rows of the table-and it often only contains one type of reasoning. The relevant rows of the table for each text are identified automatically in LogicNLG, and since identifying the relevant rows is not a trivial task, the LogicNLG examples are noisy. SciGen, on the other hand, only contains one description per table and it may contain multiple types of reasoning. SciGen is therefore more challenging than LogicNLG based on both data complexity and text complexity. For instance, LogicNLG descriptions contain 14 words on average, compared to 116 words in SciGen.

Dataset	Pairs	Cell	Num.	Text	Vocab	Domain	Source	Reasoning
WikiBIO	400K	17	3	97	400K	Biography	Crawled	No
Rotowire	11.3K	649	429	337	11.3K	Basketball	Crawled	Few
ToTTo	136K	3	1	17	136K	Open (Wikipedia)	Annotated	Few
LogicNLG	37K	91	35	14	122K	Open (Wikipedia)	Annotated	Yes
SciGen	1.3K	53	34	116	11K	Scientific	Annotated	Yes

Table 1: Comparison of SciGen to recent table-to-text generation datasets. *Pairs* shows the number of annotated pairs in each dataset. The *Cell* and *Num*. columns show the average number of total cells and cells with numerical values in the input table, respectively. *|Text|* reports the average numbers of words in descriptions. *|Vocab|* is the length of the corresponding vocabulary in each dataset.

Second, the types of logical operations that are used for creating TabFact, and therefore LogicNLG, are not limited to arithmetic operations. Based on Chen et al. (2020b)'s analyses, count is the most common logical operation in TabFact's complex statements. However, it also contains other types of reasonings like temporal reasoning—e.g., about 1200 textual descriptions in LogicNLG are generated based on before/after operations. SciGen, on the other hand, is only targeted for evaluating arithmetic reasoning.

Table 1 compares SciGen with recent table-totext generation datasets based on various properties, in particular, (a) *data complexity* that is measured by the average number of containing cells in each table, (b) *text complexity* that is measured by the average number of words and the size of the vocabulary in the target text, and (c) the *reasoning* requirement to generate target texts.

3 Dataset and Problem Definition

3.1 **Problem Definition**

SciGen is a dataset for generating descriptions from scientific tables by reasoning over their content. An input in SciGen is a table T extracted from a scientific article with its corresponding caption C, which is a word sequence containing one or few sentences about the table. $T = \{R_1, \ldots, R_n\}$ is represented as a list of lists, in which each R_i represents a row of the table. The task is to generate textual description D for the table that describes the most important findings of T by reasoning over its content.

For instance, Figure 1 shows a table from Sci-Gen alongside its annotated description.³ As we see, generating the given description requires comparing the values of "MQAN" and "ESIM (ELMO)" rows with their corresponding "+coverage" rows, as well as comparing the values of "SNLI", "Glockner", and "SICK" columns.

We can generate other descriptions from this table. For instance, "scores of all systems are higher in SNLI compared to SICK", or, "except for MQAN+coverage, the rest of examined models have higher scores on SNLI compared to Glockner", etc. However, since the table is about evaluating the impact of coverage features, as mentioned in the caption, only a specific subset of values in the table are used for generating the description. This indicates that, apart from reasoning, relevant *content selection* is also an important step in generating scientific descriptions.

3.2 Annotation Process

For creating SciGen, we have selected scientific articles from arXiv.org that are accompanied by their corresponding latex sources. The selected articles are mainly from "Computation and Language" and "Machine Learning" fields of "Computer Science". We extract tables from the selected articles using the AxCell tool (Kardas et al., 2020). Expert annotators-i.e., PhD students or postdocs in the corresponding field, and in many cases one of the co-authors of the selected papers-have annotated the description of each table as follows: (1) the description is selected from the text of the article, i.e., text spans of the article that describe the content of the table, and (2) only text spans that can be generated given the table and its caption are annotated as descriptions. For instance, the "The reason is that the dataset creation process, and therefore the task formulation is similar in SNLI and MultiNLI, but are different from SICK." text span in Figure 1 provides additional insights about the results of the table. However, it is not annotated because it cannot be generated given the table and caption.

Two expert annotators then examine the resulting annotations from the first step and remove incorrect annotations—e.g., pairs in which the table

³The table and description are from Moosavi et al. (2019).

	in-domain	out-of-domain			
	MultiNLI	SNLI	Glockner	SICK	
MQAN	72.30	60.91	41.82	53.95	
+ coverage	73.84	65.38	78.69	54.55	
ESIM (ELMO)	80.04	68.70	60.21	51.37	
+ coverage	80.38	70.05	67.47	52.65	

Table 2: Impact of using coverage for improving generalization across different datasets of the same task (NLI). All models are trained on MultiNLI.

Table 2 shows the performance for both systems for in-domain (the MultiNLI development set) as well as out-of-domain evaluations on SNLI, Glockner, and SICK datasets.

The results show that coverage information considerably improves the generalization of both examined models across various NLI datasets. The resulting cross-dataset improvements on the SNLI and Glockner datasets are larger than those on the SICK dataset. The reason is that the dataset creation process and therefore, the task formulation is similar in SNLI and MultiNLI, but they are different from SICK. In particular, in the neutral pairs

Figure 1: An example of an expert table-description annotation. The annotation process contains marking the text spans from the article that (1) describes the table, and (2) can be generated using the table and its caption.

Domain	Article	Pairs	Cell	Num.	Cap. Sent	Cap. Word	Desc. Sent	Desc. Word
Computation and Language	299	792	54.6	35.7	2.5	35.1	5.7	113.6
Machine Learning	191	410	49.9	32.1	2.4	32.5	6.0	117.2
Others	57	136	60.1	36.7	1.9	25.5	6.1	126.9

Table 2: The statistic of the expert-annotated data in SciGen. All domains consist of the "Computer Science" articles from arXiv.org. The *Article* and *Pairs* columns show the number of annotated articles and tables in each domain, respectively. *Cell* shows the average number of cells in the annotated tables, and *Num.* shows the average number of cells containing numerical values. The *Cap. sent* and *Cap. word* columns show the average number of containing sentences and words in captions, respectively. Similarly, *Desc. sent* and *Desc. word* report the average number of sentences and words that are included in the annotated descriptions.²

was not extracted correctly by the tool or the annotated description was incorrect—and those in which the description does not involve reasoning, e.g., descriptions like "The dataset statistics are summarized in Table 2".⁴

3.3 Dataset Statistics

Table 2 presents the statistics of the expertannotated table-description pairs in SciGen. The table shows the statistics for different domains of the dataset, i.e., "Computation and Language", "Machine Learning", and "Others". The articles of the "Others" domain belong to various fields like "Computational Geometry", "Distributed, Parallel, and Cluster Computing", and "Networking and Internet Architecture".

3.4 Extending the Dataset with Automatic Annotations

Annotating descriptions of tables in scientific articles requires expert knowledge about the content of the article and is very time-consuming. Therefore, there is a limit to the amount of training data that can be created using expert annotations. In this section, we propose an automatic table-description extraction process using the corresponding latex files of scientific articles. Using latex sources, we can

easily locate paragraphs of the article that include a reference to a table. This way, we can collect an arbitrarily large number of table-description pairs from any scientific domain.

We then post-process the automatically collected data by discarding pairs in which (a) the description does not reason over the table content-for this step, based on the retained and discarded pairs in the data-cleaning step of the expert annotations, we construct a set of heuristics rules for determining reasoning-aware descriptions based on captions and descriptions—, (b) the table does not contain numerical values, (c) the description describes multiple tables or figures, (d) the table is not extracted correctly by the Axcell tool, and (e) the description is too short-i.e., less than 15 words-or too long, i.e., longer than 400 words. Additionally, we shorten the paragraphs that consist of multiple subsections to only contain the one that relates to the target table.

To assess the quality of the resulting data, we automatically extract table-description pairs from the articles in our expert-annotated data and compare the automatically extracted descriptions with those that are manually annotated. Based on this comparison:

• 20% of these pairs in the automatically extracted data do not exist in the expert anno-

⁴This step results in removing 970 table-descriptions.

	Pairs	Text	BLEU	METEOR	MScore	BertS	BLEURT
automatic annotations	950	182	31.38	0.64	0.37	0.90	-0.34
+post-processing	380	123	48.36	0.70	0.44	0.92	-0.13

Table 3: Comparing the similarity of the automatically extracted table-description pairs, before and after postprocessing, to the expert annotations based on BLEU, METEOR, MoverScore, BertScore, and BLEURT metrics. *Pairs* shows the number of common tables in the expert and automatically extracted annotations. *|Text|* shows the average number of words in descriptions.

tations. Based on our analysis, most of these pairs were discarded from the expert annotations as a result of errors in the pdf annotation extraction $tool.^{5}$

- 27% of the pairs in the expert annotations do not exist in the automatic data because they were discarded by our heuristic rules in postprocessing.
- Assuming C is the set of common tables in expert and automatic annotations. In the majority of the tables in C—i.e., 71%—automatically extracted descriptions contain the annotated descriptions by experts. The average length of descriptions in the expert annotations of C's tables is 95 words while it is 113 words for automatic annotations. This indicates that automatically extracted descriptions often contain additional information, e.g., the extracted paragraph may contain explanations about the reasons that a system performs better/worse in some settings, which cannot be generated only using the table content and its caption.
- In 29% of tables in *C*, the automatic description is shorter than the expert annotation. Such cases include expert annotations that spread over multiple paragraphs among which only one or few have a reference to the table.⁶ In such cases, the automatic description is a partial description of the table.

Table 3 reports the similarity of automatic annotations, before and after post-processing, to expert annotations according to the evaluation metrics of § 4.3. As we see, the post-processing step considerably improves the similarity of automatically extracted annotations with the expert annotations.

Our automatic table-description extraction and post-processing scripts (1) will make it possible to collect high-quality unsupervised table-description pairs from any new domain, for which we have access to latex files, and (2) will facilitate expert annotations by suggesting related paragraphs to each table as well as identifying potential reasoningaware descriptions. This way, annotators will not need to read the whole article or section for finding the descriptions.

3.5 Dataset Splits

We release the dataset in three different settings: (1) *few-shot*, (2) *medium*, and (3) *large*.

The data splits in *few-shot* only contain tabledescription pairs from expert annotations. The training and development sets in this setting only contain pairs from the "Computation and Language" (C&L) articles. We split the test set into two "C&L" and "Other" domains, in which the "Other" domain mainly contains examples from the "Machine Learning" (ML) articles.

The training and development sets in the *medium* setting contain those in *few-shot* plus automatically extracted pairs from additional "C&L" articles. Similarly, the training and development sets in the *large* setting contain those in *medium* in addition to automatically extracted pairs from additional "ML" articles. The test data is the same in all three settings. The "Other" test set can be considered as an out-of-domain evaluation set for *few-shot* and *medium*. Table 4 reports the statistics of the three settings.

Setting	Domain	Train	Dev	Test
Few-shot	C&L	200	100	492
	Others	0	0	546
Medium	C&L	200+13407	100+3352	492
	Other	0	0	546
Large	C&L	200+13407	100+3352	492
	Other	26362	8677	546

Table 4: Number of table-description pairs in the training, development, and test sets of each setting. *C&L* consists of "Computation and Language" articles and *Other* mostly contains "Machine Learning" articles.

 $^{^5}Expert$ annotators have annotated the descriptions on the PDF files in the annotation process of §3.2.

⁶E.g., the annotation in the example of Figure 1.

	ellipsis (infl.)	ellipsis (VP)
baseline	53.0	28.4
concat	76.2	76.6
s-hier-to-2.tied	66.4	65.6
CADec	72.2	80.0

Table 8: Accuracy on ellipsis test set.

 65.6
 <C> 66.4 <C> 65.6 <R> <C> CADec <C> 72.2

 80.0
 <C> [BOLD] 80.0 <CAP> Table 8: Accuracy on ellipsis test set.

<R> <C> [EMPTY] <C> [BOLD] ellipsis (infl.) <C> [BOLD] ellipsis (VP) <R> <C> baseline <C> 53.0 <C> 28.4 <R> <C> concat <C> [BOLD] 76.2 <C> 76.6 <R> <C> s-hier-to-2.tied

Figure 2: A sample table—from Voita et al. (2019)—with its corresponding input representation. "<R>", "<C>", and "<CAP>" are special tokens that specify the start of the row, cell, and caption of the tables, respectively.

4 Experimental Setup

In this section, we describe the baselines, input representation format, and automatic evaluation metrics that we use in our experiments.

4.1 Baselines

Motivated by the results of Ribeiro et al. (2020) that show the BART (Lewis et al., 2020) and T5 (Raffel et al., 2020) pretrained language models consistently outperform recent specialized data-to-text models on various benchmarks, we study the effectiveness of these two models on our dataset.

BART is a Transformer-based encoder-decoder that is pretrained as a text-to-text denoising autoencoder. For BART pre-training, input texts—taken from books and Wikipedia—are corrupted with a noising function so that the model learns to reconstruct original inputs. Lewis et al. (2020) show that while BART works well for comprehension tasks, it is in particular effective for text generation. In our experiments, we use the facebook/bart-large model from HuggingFace's Transformers (Wolf et al., 2019) with 400M parameters.

T5 models various NLP tasks into a unified textto-text format, and has an on-par performance with state-of-the-art models on various NLP datasets including text classification, summarization, question answering, etc. In our experiments, we use T5-base and T5-large models from HuggingFace's Transformers with 220M, and 770M parameters, respectively.⁷

4.2 Input Representation

For using text-to-text generation baselines, we should convert input tables into a text sequence. In order to preserve the structure of the table, we use three special tokens to specify the beginning of rows, cells, and the caption of the table, namely "<R>", "<C>", and "<CAP>" tokens, respectively. Figure 2 shows an input table with its corresponding input representation.

4.3 Automatic Evaluation Metrics

We use the following evaluation metrics:

BLEU (Papineni et al., 2002) is one of the most common evaluation metrics for text generation. It computes the geometric average of the precision over output text's n-grams. We use SacreBLEU (Post, 2018)⁸ that produces comparable and reproducible BLEU scores.

METEOR (Denkowski and Lavie, 2014) aligns the output text to the reference text and calculates sentence-level similarity scores for the alignments.

BertScore (Zhang et al., 2020) uses BERT embeddings and matches words in the output and reference sentences by cosine similarity. It then computes precision, recall, and F_1 measures based on the resulting matching.

MoverScore (Zhao et al., 2019) computes the distance between the contextual representation of the output and reference texts. It captures the amount of shared content between two texts as well as how much the output texts deviate from the reference. It uses BERT embeddings for computing contextualized representations.

BLEURT (Sellam et al., 2020) is a learned evaluation metric based on BERT. It is first pre-trained on synthetic examples and then fine-tuned on human judgments for the task of machine translation.

The above metrics mostly measure the surface similarity of generated descriptions to gold ones, and they cannot evaluate the factual correctness of

⁷All models are available at https://github.com/ huggingface/transformers/tree/v2.10.0

⁸https://github.com/mjpost/sacreBLEU

- D1 The results of Table 2 shows that the addition of coverage features improves the performance of MQAN by 1.54, 4.47, 36.87, and 0,6 points on MultiNLI, SNLI, Glockner, and SICK, respectively. Similarly, it improves the performance of ESIM (ELMO) by 0.34, 1.35, 7,26, and 1,28 on the mentioned datasets, respectively. We observe the highest improvements of both systems on the Glockner dataset.
- D₂ The results of using coverage for MQAN and ESIM (ELMO) systems on various datasets are reported in Table 2. The results show that the addition of coverage features significantly decrease the performance for both MQAN and ESIM (ELMO) baselines. We observe the highest drop in the Glockner dataset.

Table 5: Sample descriptions for the table in Figure 1, where D_1 is correct and D_2 includes incorrect conclusions.

the generated descriptions given their corresponding tables. For instance, consider the sample descriptions in Table 5, where D_1 accurately describe the results while D_2 includes completely incorrect conclusions from the table. The BLEU, METEOR, BertScore, MoverScore, and BLEURT scores for D_1 are 11.65, 0.35, 0.86, 0.27, and -0.57, respectively. These scores for D_2 are 12.18, 0.30, 0.87, 0.30, and -0.54, respectively. As we see, results for both systems based on all these evaluation metrics are in the same range, and in cases higher for the incorrect description. As a result, we also perform human evaluation apart from reporting these standard metrics.

5 Results

We report the results based on both automatic evaluation metrics as well as human evaluation.

5.1 Automatic Evaluation

We report the results of our baselines—i.e., BARTlarge, T5-base, and T5-large—on different splits of the SciGen dataset in Table 6 using the evaluation metrics of §4.3. Based on the results: (1) except for BertScore, the value range for the rest of the metrics is very low, (2) BertScore values are very high for all the experiments, however, as we will see in § 5.2, generated descriptions are far from being accurate, (3) there is not a clear agreement between the rankings of different metrics for the examined models and settings, and (4) according to automatic metrics BART performs better than the two other models, however, as we will see in § 5.2, T5-large performs on-par with or in cases better than BART according to human evaluation.

As an example of model outputs, Table 7 shows the outputs of the BART-large model on one of the "C&L" test examples, i.e., the table in Figure 2.

5.2 Human Evaluation

For human evaluation, we select 58 tabledescription pairs from the SciGen "C&L" test set

Setting	Model	BLEU	METEOR	MoverS	BertS	BLEURT			
	Test(C&L)								
	BART	4.73	0.22	0.14	0.84	-0.66			
Few	T5-base	2.59	0.13	0.02	0.79	-1.09			
	T5-large	3.16	0.16	0.06	0.81	-0.95			
	BART	5.30	0.23	0.13	0.84	-0.72			
Medium	T5-base	3.32	0.15	0.05	0.82	-0.89			
	T5-large	3.65	0.17	0.10	0.83	-0.77			
	BART	5.04	0.22	0.14	0.84	-0.71			
Large	T5-base	3.38	0.15	0.06	0.82	-0.85			
	T5-large	3.84	0.18	0.10	0.83	-0.79			
			Test(Other)						
	BART	4.26	0.22	0.13	0.84	-0.66			
Few	T5-base	2.12	0.11	0.00	0.78	-1.13			
	T5-large	2.58	0.14	0.03	0.81	-1.01			
	BART	4.64	0.22	0.12	0.84	-0.71			
Medium	T5-base	3.17	0.14	0.04	0.81	-0.91			
	T5-large	3.00	0.16	0.07	0.82	-0.81			
	BART	4.66	0.20	0.12	0.84	-0.72			
Large	T5-base	3.66	0.14	0.05	0.81	-0.86			
-	T5-large	3.61	0.17	0.07	0.82	-0.81			

Table 6: Performance of the BART-large, T5-base, and T5-large models on the *few-shot*, *Medium*, and *Large* settings based on BLEU, METEOR, Mover-Score, BertScore, and BLEURT metrics. We report the scores on "C&L" and "Other" test domains separately. The "Other" examples are out-of-domain test data for the *few-shot* and *Medium* settings. Highest scores for each domain are boldfaced.

and their corresponding system-generated descriptions from the BART and T5-large models for the three settings.⁹ We break down each description, both gold and system-generated ones—i.e., $58 \times 2 \times 3$ descriptions—to a list of individual statements. For instance, the corresponding statements with the gold description in Table 7 are (a) "For ellipsis, both models improve substantially over the baseline (by 19-51 percentage points)", (b) "concat is stronger for inflection tasks", and (c) "CADec is stronger for VPellipsis".

We assign one of the following labels to each of the extracted statements from system-generated descriptions: (1) *entailed*: a generated statement that is entailed by the corresponding gold description,

⁹Based on our initial analysis, the quality of the descriptions of these two models were higher than those of T5-base.

Setting	Description
Gold	For ellipsis, both models improve substantially over the baseline (by 19-51 percentage points), with concat stronger for inflection tasks and CADec stronger for VPellipsis.
Few	table 8 shows the performance of our model on the ellipsis test set. our concat model outperforms all the other models that do not use concat embeddings.
Med.	the results show that both concat and cadec significantly outperform the baseline and s-hier-to-2.tied on the ellipsis test set. the concat model is better than the baseline on both inflected and verb-inflated forms, while cadec is better on the vp version. on the inflated version, cadec outperforms the baseline by a large margin and is within 3 points of the best model. the model with the best accuracy on the verb portion of the test set (80.0%) is very close to the best performance on the noun portion (75.6%). the results indicate that cadec has better generalizability than the other models.
Large	the concatenation model concatenates the output of the concatenated encoder and decoder, and the s-hier-to-2.tied model performs slightly better than concat on both ellipsis (infl.) and vice versa. cadec outperforms the other two models by a large margin, achieving 80.0% accuracy on the vp ellipsoids, which is the new state of the art.

Table 7: Sample outputs from the BART model on few-shot, medium, and large settings on the corresponding table from Figure 2. Factually correct statements are marked with green, factually incorrect statements and hallucinations are marked with red and blue, respectively.

Model	Setting	Rec.	Prec.	Corr.	Hall.
	Few	0.0	0.0	0.1	0.4
BART	Medium	0.1	0.1	0.4	0.2
	Large	0.1	0.1	0.2	0.3
	Few	0.0	0.0	0.1	0.3
T5-Large	Medium	0.1	0.1	0.4	0.2
	Large	0.1	0.1	0.3	0.1

Table 8: The results of human evaluation for the BART and T5-large models according to the recall, precision, correctness, and hallucination metrics. Best values i.e., highest coverage, precision, and correctness values and the lowest hallucination—are bold-faced.

i.e., is equivalent to one of the extracted statements from the gold description, (2) *extra*: a statement that is not entailed by the gold description but is correct based on the table content, (3) *incorrect*: a statement that is relevant to the table but is factually incorrect—e.g., "the s-hier-to-2.tied model performs slightly better than concat on both ellipsis (infl.) and vice versa." in Table 7 contains relevant entities that are mentioned in the table, but the statement is incorrect—, and (4) *hallucinated*: a statement that is irrelevant to the table.

Based on the above labels, we compute four metrics as follows:

Recall: the ratio of the statements in the gold description that are covered by the system-generated description, i.e., $\frac{|\text{entailed statements}|}{|\text{gold statements}|}$ per description.

Precision: the ratio of the statements in the system-generated description that exist in the gold description, i.e., $\frac{|\text{entailed statements}|}{|\text{generated statements}|}$ per description.

Correctness: the ratio of the statements in the system-generated description that are factually correct, i.e., $\frac{|\text{entailed statements}|+|\text{extra statements}|}{|\text{generated statements}|}$.

hallucination: the ratio of irrelevant statements with regard to the table that is computed as <u>hallucinated statements</u>.

Table 8 presents the results of the human evaluation. Based on the results, we observe that:

- The addition of automatically extracted pairs in the medium and large settings improves the recall, precision, and correctness of the generated descriptions and decreases their hallucination.
- Compared to the medium setting, the generated descriptions in the large setting contain a larger number of *factually-incorrect* facts, i.e., lower correctness. This could be due to the fact that the additional table-description pairs in large are from a different domain, i.e., "ML".
- There is still a large gap for solving the Sci-Gen dataset, i.e., in the best case, only 40% of the generated statements are correct and only 10% of the gold statements are covered by the generated description.

6 Conclusions

In this paper, we introduce SciGen that is a challenging dataset for reasoning-aware data-to-text generation. The input data in SciGen are tables from scientific articles and generating their corresponding descriptions requires arithmetic reasoning over table values. Annotating scientific articles requires expert knowledge and is therefore very costly and does not scale to large data sizes. Therefore, we provide a pipeline to extract high-quality unsupervised table-description pairs from the corresponding latex files of scientific articles. We release SciGen in three different settings-i.e., few-shot, medium, and large-based on the size of the available training data. The few-shot setting and the test set contain expert-annotated table-description pairs while the training data in medium and large settings contain automatically extracted table-description pairs. We study two state-of-the-art data-to-text generation models on SciGen and evaluate the results using both common automatic evaluation metrics as well as human evaluation. Our results show that (1) common metrics are not reliable for evaluating reasoning-aware text generation, i.e., they do not correlate with human evaluation and they also do not agree with each other, and (2) adding automatically extracted annotations improves the correctness of the generated descriptions and reduces the hallucination, however, there is still a large gap to solve the dataset. We ensure reproducibility at all levels by releasing our data, full extraction and post-processing pipeline, and by publishing the raw responses from our human evaluation.

Acknowledgements

The authors would like to thank Xia Zeng, Dennis Kuhn, Ragini Monen, and Gisela Vallejo for their great efforts in the data collection process. We gratefully acknowledge the help of numerous members of this research community in helping us with the annotations. This work was possible thanks to all these contributing researchers. We would also like to thank Michael Bugert, Jan-Christoph Kile, Ji-Ung Lee, Yevgeniy Puzikov, Kevin Stowe, and Ajie Utama for their valuable feedbacks on the paper. This work has been funded by the German Federal Ministry of Education and Research and the Hessian Ministry of Higher Education, Research, Science and the Arts within their joint support of the National Research Center for Applied Cybersecurity ATHENE.

References

Anusha Balakrishnan, Jinfeng Rao, Kartikeya Upasani, Michael White, and Rajen Subba. 2019. Constrained decoding for neural NLG from compositional representations in task-oriented dialogue. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 831– 844, Florence, Italy. Association for Computational Linguistics.

- Anja Belz. 2008. Automatic generation of weather forecast texts using comprehensive probabilistic generation-space models. *Natural Language Engineering*, 14(4):431.
- Anja Belz, Mike White, Dominic Espinosa, Eric Kow, Deirdre Hogan, and Amanda Stent. 2011. The first surface realisation shared task: Overview and evaluation results. In *Proceedings of the 13th European Workshop on Natural Language Generation*, pages 217–226, Nancy, France. Association for Computational Linguistics.
- David L Chen and Raymond J Mooney. 2008. Learning to sportscast: a test of grounded language acquisition. In *Proceedings of the 25th international conference on Machine learning*, pages 128–135.
- Wenhu Chen, Jianshu Chen, Yu Su, Zhiyu Chen, and William Yang Wang. 2020a. Logical natural language generation from open-domain tables. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7929– 7942, Online. Association for Computational Linguistics.
- Wenhu Chen, Hongmin Wang, Jianshu Chen, Yunkai Zhang, Hong Wang, Shiyang Li, Xiyou Zhou, and William Yang Wang. 2020b. Tabfact : A large-scale dataset for table-based fact verification. In *International Conference on Learning Representations* (*ICLR*), Addis Ababa, Ethiopia.
- Michael Denkowski and Alon Lavie. 2014. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 376–380, Baltimore, Maryland, USA. Association for Computational Linguistics.
- Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2018. Findings of the E2E NLG challenge. In Proceedings of the 11th International Conference on Natural Language Generation, pages 322–328, Tilburg University, The Netherlands. Association for Computational Linguistics.
- Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2020. Evaluating the State-of-the-Art of End-to-End Natural Language Generation: The E2E NLG Challenge. *Computer Speech & Language*, 59:123–156.
- Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. The WebNLG challenge: Generating text from RDF data. In *Proceedings of the 10th International Conference on Natural Language Generation*, pages 124–133, Santiago de Compostela, Spain. Association for Computational Linguistics.
- Albert Gatt and Emiel Krahmer. 2018. Survey of the state of the art in natural language generation: Core tasks, applications and evaluation. *Journal of Artificial Intelligence Research*, 61:65–170.

- Marcin Kardas, Piotr Czapla, Pontus Stenetorp, Sebastian Ruder, Sebastian Riedel, Ross Taylor, and Robert Stojnic. 2020. Axcell: Automatic extraction of results from machine learning papers. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). Association for Computational Linguistics.
- Rik Koncel-Kedziorski, Dhanush Bekal, Yi Luan, Mirella Lapata, and Hannaneh Hajishirzi. 2019. Text Generation from Knowledge Graphs with Graph Transformers. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), pages 2284–2293, Minneapolis, Minnesota. Association for Computational Linguistics.
- 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*, pages 1203–1213, Austin, Texas. Association for Computational Linguistics.
- Chris van der Lee, Emiel Krahmer, and Sander Wubben. 2017. PASS: A Dutch data-to-text system for soccer, targeted towards specific audiences. In *Proceedings of the 10th International Conference on Natural Language Generation*, pages 95–104, Santiago de Compostela, Spain. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pretraining for natural language generation, translation, and comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Nafise Sadat Moosavi, Prasetya Ajie Utama, Andreas Rücklé, and Iryna Gurevych. 2019. Improving generalization by incorporating coverage in natural language inference. *arXiv preprint arXiv:1909.08940*.
- Preksha Nema, Shreyas Shetty, Parag Jain, Anirban Laha, Karthik Sankaranarayanan, and Mitesh M. Khapra. 2018. Generating descriptions from structured data using a bifocal attention mechanism and gated orthogonalization. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 1539–1550, New Orleans, Louisiana. Association for Computational Linguistics.
- Shereen Oraby, Lena Reed, Shubhangi Tandon, Sharath T.S., Stephanie Lukin, and Marilyn Walker. 2018. Controlling personality-based stylistic variation with neural natural language generators. In Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue, pages 180–190, Melbourne,

Australia. Association for Computational Linguistics.

- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Ankur P Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuwan Dhingra, Diyi Yang, and Dipanjan Das. 2020. Totto: A controlled table-totext generation dataset. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). Association for Computational Linguistics.
- Matt Post. 2018. A call for clarity in reporting BLEU scores. In Proceedings of the Third Conference on Machine Translation: Research Papers, pages 186– 191, Belgium, Brussels. Association for Computational Linguistics.
- Raheel Qader, Khoder Jneid, François Portet, and Cyril Labbé. 2018. Generation of company descriptions using concept-to-text and text-to-text deep models: dataset collection and systems evaluation. In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 254–263. Association for Computational Linguistics.
- Dragomir Radev, Rui Zhang, Amrit Rau, Abhinand Sivaprasad, Chiachun Hsieh, Nazneen Fatema Rajani, Xiangru Tang, Aadit Vyas, Neha Verma, Pranav Krishna, et al. 2020. Dart: Open-domain structured data record to text generation. *arXiv preprint arXiv:2007.02871*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the limits of transfer learning with a unified text-totext transformer. *Journal of Machine Learning Research*, 21(140):1–67.
- Lena Reed, Shereen Oraby, and Marilyn Walker. 2018. Can neural generators for dialogue learn sentence planning and discourse structuring? In *Proceedings of the 11th International Conference on Natural Language Generation*, pages 284–295, Tilburg University, The Netherlands. Association for Computational Linguistics.
- Ehud Reiter and Robert Dale. 2000. *Building natural language generation systems*. Cambridge university press.
- Leonardo FR Ribeiro, Martin Schmitt, Hinrich Schütze, and Iryna Gurevych. 2020. Investigating pretrained language models for graph-to-text generation. *arXiv* preprint arXiv:2007.08426.

- Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. BLEURT: Learning robust metrics for text generation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.
- Craig Thomson, Ehud Reiter, and Somayajulu Sripada. 2020. Sportsett: Basketball - a robust and maintainable dataset for natural language generation. Intel-LanG : Intelligent Information Processing and Natural Language Generation ; Conference date: 07-09-2020 Through 07-09-2020.
- Elena Voita, Rico Sennrich, and Ivan Titov. 2019. When a good translation is wrong in context: Context-aware machine translation improves on deixis, ellipsis, and lexical cohesion. In *Proceedings* of the 57th Annual Meeting of the Association for Computational Linguistics, pages 1198–1212, Florence, Italy. Association for Computational Linguistics.
- Qingyun Wang, Xiaoman Pan, Lifu Huang, Boliang Zhang, Zhiying Jiang, Heng Ji, and Kevin Knight. 2018. Describing a knowledge base. In Proceedings of the 11th International Conference on Natural Language Generation, pages 10–21, Tilburg University, The Netherlands. Association for Computational Linguistics.
- Sam Wiseman, Stuart Shieber, and Alexander Rush. 2017. Challenges in data-to-document generation. In Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, pages 2253–2263, Copenhagen, Denmark. Association for Computational Linguistics.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2019. Huggingface's transformers: State-of-the-art natural language processing. *ArXiv*, abs/1910.03771.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. BERTScore: Evaluating text generation with bert. In International Conference on Learning Representations (ICLR).
- Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, and Steffen Eger. 2019. MoverScore: Text generation evaluating with contextualized embeddings and earth mover distance. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 563–578, Hong Kong, China. Association for Computational Linguistics.

A Upper-Bound Estimation

To estimate an upper bound for the performance of generation models on our dataset based on automatic metrics, we randomly select 50 tables from expert-annotations. We provide these tables to a new expert annotator and ask them to describe the tables in their own words without looking at the gold descriptions or the result sections of the corresponding articles.¹⁰ Table 9 shows the scores of these 50 annotated tables compared to their corresponding gold annotations. The high-value range for all the metrics shows that automatic metrics can acknowledge the correctness of accurate descriptions, i.e., those written by an expert. However, as we see in § 5.1, they cannot discriminate imperfect descriptions.

BLEU	METEOR	MoverS	BertS	BLEURT
66.73	0.614	0.98	0.99	0.95

Table 9: The automatic evaluation scores for 50 tabledescriptions in which an expert has written descriptions based on table contents and without looking at the gold descriptions.

B Impact of Table Captions

Table 10 shows the impact of captions on automatic evaluation metrics. The *caption* row shows the results when the caption is considered as the description, i.e., evaluating captions compared to gold descriptions. The *BART* and *T5-large* rows show the result of these two models in the fewshot setting where the captions of the tables were excluded from the input data.

Model	BLEU	METEOR	MoverS	BertS	BLEURT
caption	1.68	0.11	0.12	0.84	-0.98
BART	2.39	0.17	0.08	0.83	-0.72
T5-large	1.50	0.12	0.06	0.81	-0.98

Table 10: Evaluating the impact of table captions. The reported results are equivalent to the reported few-shot results in Table 6 but with the difference that the table captions are excluded from the input. The caption rows show the values of evaluation metrics when the caption is considered as the description.

¹⁰They had access to the background sections of the article in case the information of the table itself is not enough to describe them.

Setting	Model	Description
Gold		For ellipsis, both models improve substantially over the baseline (by 19-51 percentage points), with concat stronger for inflection tasks and CADec stronger for VPellipsis.
Few	BART	table 8 shows the performance of our model on the ellipsis test set. our concat model outperforms all the other models that do not use concat embeddings.
	T5-base	$<$ extra_id_0 > c >[bold] ellipsis (infl .) c >[s - hier - to - 2 . tied] c >(s - hier - to - 2 . tied) c >s - hier - to - 2 . tied c >(s - hier - to - 2 . tied) c >80 % accuracy on ellipsis test set .
	T5-large	$<$ extra_id_0 >vp) and concatenation (s - hier - to - 2. tied , infl.), respectively , are shown in table 8. as expected , concatenation improves the accuracy on the infl. ellipsis test set .
Medium	BART	the results show that both concat and cadec significantly outperform the baseline and s-hier-to-2.tied on the ellipsis test set. the concat model is better than the baseline on both inflected and verb-inflated forms, while cadec is better on the vp version. on the inflated version, cadec outperforms the baseline by a large margin and is within 3 points of the best model. the model with the best accuracy on the verb portion of the test set (80.0%) is very close to the best performance on the noun portion (75.6%). the results indicate that cadec has better generalizability than the other models.
	T5-base	$<$ extra_id_0 > the s - hier - to - 2. tied model outperforms the baseline on both infl. and vp ellipsis datasets. cadec outperforms s - hier - to - 2. tied on both infl. and vp ellipsis datasets. cadec achieves the best performance on both infl. and vp ellipsis datasets.
	T5-large	can see that both concatenation and s - hier - to - 2. tied approaches are significantly better than the baseline . cadec outperforms s - hier - to - 2. tied on both inflectional and verbal ellipsis and is comparable to the baseline on verb ellipsis .
Large	BART	the concatenation model concatenates the output of the concatenated encoder and decoder, and the s-hier-to-2.tied model performs slightly better than concat on both ellipsis (infl.) and vice versa. cadec outperforms the other two models by a large margin, achieving 80.0% accuracy on the vp ellipsoids, which is the new state of the art.
	T5-base	cadec outperforms the baseline and s - hier - to - 2. tied by a large margin. cadec also outperforms s - hier - to - 2. tied by a large margin.
	T5-large	concatenated embeddings outperform the s - hier - to - 2 . tied baseline by 6 $\%$ on infrequent words , and by 5 $\%$ on vp ellipsis . the concatenated embeddings outperform the s - hier - to - 2 . tied baseline by 10 $\%$ on infrequent words , and by 5 $\%$ on vp ellipsis .

Table 11: Sample outputs from the examined models in Table 6 on the corresponding table from Figure 2.

C Example outputs from all baselines

Table 11 provides the output of all the examined baselines for the table in Figure 2.