
SciGen: Dataset Documentation

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

¹UKP Lab, Technische Universität Darmstadt

²Department of Computer and Information Science, UPenn

¹<https://www.ukp.tu-darmstadt.de>

²<https://www.seas.upenn.edu/~danroth>

1 We use the dataset documentation framework proposed by Gebru et al. (2018).

2 1 Dataset Documentation and Intended Uses

3 **For what purpose was the dataset created? Was there a specific task in mind? Was there a**
4 **specific gap that needed to be filled? Please provide a description.** This dataset is created to
5 enable research on reasoning-aware table-to-text generation based on scientific tables. Existing
6 table-to-text generation datasets mostly contain descriptions that are surface-level summaries of the
7 data (Chen and Mooney, 2008; Belz et al., 2011; Lebrecht et al., 2016; Gardent et al., 2017; Dušek
8 et al., 2018; Koncel-Kedziorski et al., 2019; Radev et al., 2020; Parikh et al., 2020). The generation
9 of descriptions in SciGen, on the other hand, requires performing arithmetic reasoning. As our results
10 show, despite the impressive progress of large pretrained language models on various NLP tasks
11 including generation tasks, their performance is severely limited when the task requires performing
12 arithmetic reasoning. The availability of SciGen opens new research directions in reasoning-aware
13 text generation and its evaluation.

14 **Who created this dataset (e.g., which team, research group) and on behalf of which entity (e.g.,**
15 **company, institution, organization)?** This dataset is created by researchers from UKP Lab from
16 the Computer Science Department at Technische Universität Darmstadt and the Department of
17 Computer and Information Science at UPenn.

18 **Who funded the creation of the dataset?** The expert annotations are done voluntarily by one of
19 the authors of the included scientific articles. The selection step—i.e., selecting table-descriptions
20 that contain arithmetic reasoning and removing table-description pairs that contain errors—, and
21 human evaluation were funded by the third parties which we acknowledged in the paper.

22 **Any other comments?** N/A

23 2 Composition

24 **What do the instances that comprise the dataset represent (e.g., documents, photos, people,**
25 **countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people and**
26 **interactions between them; nodes and edges)? Please provide a description.** Each instance in
27 SciGen presents a table-description pair. Both tables and descriptions are taken from the articles of
28 the computer science field from arXiv.org. Figure 1 shows a sample table-description pair.

29 **How many instances are there in total (of each type, if appropriate)?** SciGen is released in
30 three different settings. The few-shot setting contains 1338 instances. The medium setting contains
31 18097 instances. The large setting contains 53136 instances.

32 **Does the dataset contain all possible instances or is it a sample (not necessarily random) of**
33 **instances from a larger set? If the dataset is a sample, then what is the larger set? Is the sam-**
34 **ple representative of the larger set (e.g., geographic coverage)? If so, please describe how this**
35 **representativeness was validated/verified. If it is not representative of the larger set, please**
36 **describe why not (e.g., to cover a more diverse range of instances, because instances were with-**
37 **held or unavailable).** The instances in SciGen are mainly from the “Computation and Language”
38 and “Machine Learning” fields of “Computer Science” English articles from arXiv.org. We consider
39 the table-description pairs from “Computation and Language” as in-domain training and evaluation
40 data and those from “Machine Learning” as out-of-domain. Our automatic annotation pipeline can be
41 used to extend the data with any other English articles that are accompanied by their corresponding
42 LaTeX sources. However, the limitation is that our pipeline is developed based on the English
43 language.

44 **What data does each instance consist of? “Raw” data (e.g., unprocessed text or images) or**
45 **features? In either case, please provide a description.** Each instance—i.e., an entry in a json
46 file—represents a table-description and consist of the following fields:

- 47 • paper: the title of the corresponding article
- 48 • paper-id: the arXiv identifier of the corresponding article
- 49 • table-caption: the caption of the table
- 50 • table-column-names: the first row of the table, i.e., the header row
- 51 • table-content-values: the list of rows, excluding the header row, in the table
- 52 • text: the description of the table

53 **Is there a label or target associated with each instance? If so, please provide a description.** The
54 target of each instance is the description of the table—i.e., the “text” field described for the above
55 question.

56 **Is any information missing from individual instances? If so, please provide a description,**
57 **explaining why this information is missing (e.g., because it was unavailable). This does not**
58 **include intentionally removed information, but might include, e.g., redacted text.** No, the
59 dataset does not contain missing information.

60 **Are relationships between individual instances made explicit (e.g., users’ movie ratings, social**
61 **network links)? If so, please describe how these relationships are made explicit.** Yes. The
62 table-description pairs the belong to the same paper can be identified by both of “paper-id” and
63 “paper” data fields. There is no explicit relation between table-description pairs that belong to different
64 papers.

65 **Are there recommended data splits (e.g., training, development/validation,testing)? If so,**
66 **please provide a description of these splits, explaining the rationale behind them.** Yes. We
67 have provided a fixed training, development, and test splits regarding two different configurations:
68 (1) based on the amount of available training data, and (2) the domain of training vs. test data. The
69 provided splits enable evaluating examined models in two different scenarios: (1) when the amount
70 of training data increases, and (2) when evaluated on in-domain vs. out-of-domain instances.

71 **Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide**
72 **a description.** A subset of the training data in the medium and large settings that are extracted
73 automatically may contain additional information—e.g., describing the intuitions and reasons for
74 the results of the tables—that cannot be generated by only using the table content and its caption. In
75 addition, the table-description pairs are extracted automatically from LaTeX files using the AxCell
76 tool (Kardas et al., 2020). While extracting tables and texts from LaTeX files is straightforward and
77 accurate, there may be some cases in which the extracted content is noisy. The test set, on the other
78 hand, is annotated and validated by experts and we are not aware of any errors or noise in the test set.

79 **Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,**
80 **websites, tweets, other datasets)?** SciGen is self-contained.

81 **Does the dataset contain data that might be considered confidential (e.g., data that is protected**
82 **by legal privilege or by doctor-patient confidentiality, data that includes the content of individ-**
83 **uals' non-public communications)? If so, please provide a description.** SciGen does not contain
84 any confidential data.

85 **Does the dataset contain data that, if viewed directly, might be offensive, insulting, threatening,**
86 **or might otherwise cause anxiety? If so, please describe why.** SciGen contains table-description
87 pairs with tables that mostly contain numerical values and their corresponding descriptions. There
88 are few tables, which are related to articles about toxicity or emotion detection, in which some of the
89 header cells in the table may contain swear words.

90 **Does the dataset relate to people? If not, you may skip the remaining questions in this section.**
91 No, SciGen does not relate to people.

92 **3 Collection Process**

93 **How was the data associated with each instance acquired? Was the data directly observable**
94 **(e.g., raw text, movie ratings), reported by subjects (e.g., survey responses), or indirectly in-**
95 **ferred/derived from other data (e.g., part-of-speech tags, model-based guesses for age or lan-**
96 **guage)? If data was reported by subjects or indirectly inferred/derived from other data, was**
97 **the data validated/verified? If so, please describe how.** The table-description pairs are all directly
98 selected (observed) from the content of scientific articles.

99 **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or sen-**
100 **sor, manual human curation, software program, software API)? How were these mechanisms**
101 **or procedures validated?** The data is collected in two different ways: (1) manual annotations by
102 experts, and (2) automatic annotation extraction. Manual annotations are done in two different steps
103 (a) annotating all the table-descriptions pairs in the articles by one of the authors of the articles, and
104 (b) selecting annotated table-description pairs that require arithmetic reasoning and discarding noisy
105 pairs. Both steps of the manual annotation are done by experts.

106 For the automatic annotation, we have used the AxCell tool (Kardas et al., 2020) and a set of heuristic
107 rules for extracting table-description pairs from LaTeX sources of scientific articles. As discussed in
108 Section 3.3, we have validated our pipeline by comparing them with expert annotations. To do so,
109 we extract the table-description pairs from those articles for which we have the expert annotations.
110 We then compare the automatically extracted pairs with those that are manually annotated and have
111 reported the results.

112 **If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deterministic,**
113 **probabilistic with specific sampling probabilities)?** The candidate arXiv articles are selected
114 from the period 2010-2020 and mainly from the “Computation and Language” and “Machine
115 Learning” fields of “Computer Science” English articles from arXiv.org.

116 **Who was involved in the data collection process (e.g., students, crowdworkers, contractors) and**
117 **how were they compensated (e.g., how much were crowdworkers paid)?** The expert annotations of
118 table-description pairs are done voluntarily by one of the authors of the included scientific articles.
119 We paid 10 per hour for the table-description pair selection step, and human evaluation.

120 **Over what timeframe was the data collected? Does this timeframe match the creation time-**
121 **frame of the data associated with the instances (e.g., recent crawl of old news articles)? If not,**
122 **please describe the time-frame in which the data associated with the instances was created.**
123 The manual annotations were collected from January 2020 to June 2020. The automatic annotations
124 are extracted in July 2020. SciGen is self-contained and its content will not change over time.

125 **Were any ethical review processes conducted (e.g., by an institutional review board)? If so,**
126 **please provide a description of these review processes, including the outcomes, as well as a link**
127 **or other access point to any supporting documentation.** No.

128 **Does the dataset relate to people? If not, you may skip the remainder of the questions in this**
129 **section.** No.

130 **4 Preprocessing/cleaning/labeling**

131 **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,**
132 **tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing**
133 **of missing values)? If so, please provide a description. If not, you may skip the remainder of**
134 **the questions in this section.** Yes. As described in Section 3.2, for the manual annotations, 970
135 table-description pairs were discarded from the manual annotations that were collected in the first
136 step. The purpose of this selection step is to only include pairs that contain arithmetic reasoning
137 and to discard pairs that contain errors. Also, as described in Section 3.3, our annotation extraction
138 pipeline includes a post-processing step that discards a subset of initially extracted pairs based on
139 a set of rules. For instance, we can see from the analysis of Table 3 that the post-processing step
140 reduces the extracted table-description pairs from 950 to 380.

141 **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support**
142 **unanticipated future uses)? If so, please provide a link or other access point to the “raw” data.**
143 No. The discarded pairs are those that do not require arithmetic reasoning or contain an error.
144 Therefore, they are not related to the SciGen dataset, i.e., a reasoning-aware data-to-text generation
145 dataset.

146 **Is the software used to preprocess/clean/label the instances available? If so, please provide**
147 **a link or other access point.** The `postprocessing.py` for the post-processing step of automatic
148 annotations is available in the extraction-pipeline directory of the supplementary materials. The data
149 cleaning for the manual annotations was done manually and by experts.

150 **5 Uses**

151 **Has the dataset been used for any tasks already? If so, please provide a description** We have
152 used ScGen for the experiments of Section 5 for the task of reasoning-aware data-to-text generation.

153 **Is there a repository that links to any or all papers or systems that use the dataset? If so, please**
154 **provide a link or other access point.** <https://paperswithcode.com/dataset/scigen>

155 **What (other) tasks could the dataset be used for?** If we also include the whole text of the articles
156 to SciGen, the annotations can be used for the task of relevant content selection. The input will be
157 a table and the text of its corresponding article, and the task is to select all the text spans from the
158 article that describe the table.

159 **Is there anything about the composition of the dataset or the way it was collected and prepro-**
160 **cessed/cleaned/labeled that might impact future uses? For example, is there anything that a**
161 **future user might need to know to avoid uses that could result in unfair treatment of individ-**
162 **uals or groups (e.g., stereotyping, quality of service issues) or other undesirable harms (e.g.,**
163 **financial harms, legal risks) If so, please provide a description. Is there anything a future user**
164 **could do to mitigate these undesirable harms?** No, there is minimal risk for undesirable harm
165 from arXiv papers.

166 **Are there tasks for which the dataset should not be used? If so, please provide a description.**
167 The dataset should not be used for generating fake scientific articles.

168 **6 Distribution**

169 **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution,**
170 **organization) on behalf of which the dataset was created? If so, please provide a description.**
171 Yes, the dataset is publicly available at <https://github.com/UKPLab/SciGen>. We will update

172 the repository with human evaluation, the extraction pipeline, and baseline codes, which are already
173 available in the supplementary materials.

174 **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the**
175 **dataset have a digital object identifier (DOI)?** The dataset is distributed on GitHub. We will
176 make the DOI available after finalizing the GitHub repository.

177 **When will the dataset be distributed?** The dataset itself is already distributed. The repository
178 will be updated with the human evaluation, the extraction pipeline, and baseline codes after the
179 NeurIPS 2021 camera-ready deadline.

180 **Will the dataset be distributed under a copyright or other intellectual property (IP) license,**
181 **and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU, and**
182 **provide a link or other access point to, or otherwise reproduce, any relevant licensing terms or**
183 **ToU, as well as any fees associated with these restrictions.** SciGen is licensed under a Creative
184 Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

185 **Have any third parties imposed IP-based or other restrictions on the data associated with the**
186 **instances? If so, please describe these restrictions, and provide a link or other access point to,**
187 **or otherwise reproduce, any relevant licensing terms, as well as any fees associated with these**
188 **restrictions.** No.

189 **Do any export controls or other regulatory restrictions apply to the dataset or to individual**
190 **instances? If so, please describe these restrictions, and provide a link or other access point to,**
191 **or otherwise reproduce, any supporting documentation.** No.

192 7 Maintenance

193 **Who is supporting/hosting/maintaining the dataset?** The dataset will be supported, hosted, and
194 maintained by the UKP Lab at Technische Universität Darmstadt.

195 **How can the owner/curator/manager of the dataset be contacted (e.g., email ad-**
196 **dress)?** All comments and questions about SciGen can be sent to Nafise Sadat Moosavi:
197 moosavi@ukp.informatik.tu-darmstadt.de Other contacts can be found at: [https://www.](https://www.informatik.tu-darmstadt.de/ukp/ukp_home/staff_ukp/index.en.jsp)
198 [informatik.tu-darmstadt.de/ukp/ukp_home/staff_ukp/index.en.jsp](https://www.informatik.tu-darmstadt.de/ukp/ukp_home/staff_ukp/index.en.jsp).

199 **Is there an erratum? If so, please provide a link or other access point.** All changes to the
200 dataset will be announced on <https://github.com/UKPLab/SciGen>.

201 **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete in-**
202 **stances)? If so, please describe how often, by whom, and how updates will be communicated**
203 **to users (e.g., mailing list, GitHub)?** Yes. Any updates to the dataset will be communicated via
204 the GitHub repository, and the UKP Lab website. Potential updates are creating new training settings
205 that contain larger automatically extracted training data.

206 **If the dataset relates to people, are there applicable limits on the retention of the data associ-**
207 **ated with the instances (e.g., were individuals in question told that their data would be retained**
208 **for a fixed period of time and then deleted)?** N/A

209 **Will older versions of the dataset continue to be supported/hosted/maintained? If so, please**
210 **describe how. If not, please describe how its obsolescence will be communicated to users** All
211 versions of SciGen will be continue to be supported and maintained on [https://github.com/](https://github.com/UKPLab/SciGen)
212 [UKPLab/SciGen](https://github.com/UKPLab/SciGen). We will post the updates on the GitHub repository.

213 **If others want to extend/augment/build on/contribute to the dataset, is there a mechanism**
214 **for them to do so? If so, please provide a description. Will these contributions be vali-**
215 **dated/verified? If so, please describe how. If not, why not? Is there a process for commu-**
216 **nicating/distributing these contributions to other users? If so, please provide a description.**

Filenames	dataset/train/few-shot/train.json dataset/train/medium/train.json dataset/train/large/train.json dataset/development/few-shot/dev.json dataset/development/medium/dev.json dataset/development/large/dev.json dataset/test/test-CL.json dataset/test/test-Other.json
Format	json
URL	https://github.com/UKPLab/SciGen
Domain	Scientific articles
Keywords	data-to-text generation, arithmetic reasoning, scientific articles
License	CC BY-NC-SA 4.0

Table 1: Metadata of SciGen

217 Yes. Others can contact the authors of this paper describing their proposed extension or contribution.
218 We would discuss their proposed contribution to confirm its validity, and if confirmed, we will release
219 a new version of the dataset on GitHub and will announce it accordingly.

220 8 Metadata

221 Table 1 presents the metadata of the SciGen dataset.

222 9 Responsibility

223 The authors bear all responsibility in case of violation of rights, etc. We confirm that the dataset
224 is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International
225 License.

226 References

- 227 Anja Belz, Mike White, Dominic Espinosa, Eric Kow, Deirdre Hogan, and Amanda Stent. 2011. The
228 first surface realisation shared task: Overview and evaluation results. In *Proceedings of the 13th*
229 *European Workshop on Natural Language Generation*, pages 217–226, Nancy, France. Association
230 for Computational Linguistics.
- 231 David L Chen and Raymond J Mooney. 2008. Learning to sportscast: a test of grounded language
232 acquisition. In *Proceedings of the 25th international conference on Machine learning*, pages
233 128–135.
- 234 Ondřej Dušek, Jekaterina Novikova, and Verena Rieser. 2018. Findings of the E2E NLG challenge.
235 In *Proceedings of the 11th International Conference on Natural Language Generation*, pages
236 322–328, Tilburg University, The Netherlands. Association for Computational Linguistics.
- 237 Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. The
238 WebNLG challenge: Generating text from RDF data. In *Proceedings of the 10th International*
239 *Conference on Natural Language Generation*, pages 124–133, Santiago de Compostela, Spain.
240 Association for Computational Linguistics.
- 241 Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wal-
242 lach, Hal Daumé III, and Kate Crawford. 2018. Datasheets for datasets. *arXiv preprint*
243 *arXiv:1803.09010*.
- 244 Marcin Kardas, Piotr Czapla, Pontus Stenetorp, Sebastian Ruder, Sebastian Riedel, Ross Taylor,
245 and Robert Stojnic. 2020. Axccl: Automatic extraction of results from machine learning papers.

- 246 In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*
247 (*EMNLP*). Association for Computational Linguistics.
- 248 Rik Koncel-Kedziorski, Dhanush Bekal, Yi Luan, Mirella Lapata, and Hannaneh Hajishirzi. 2019.
249 Text Generation from Knowledge Graphs with Graph Transformers. In *Proceedings of the*
250 *2019 Conference of the North American Chapter of the Association for Computational Linguistics:*
251 *Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2284–2293, Minneapolis,
252 Minnesota. Association for Computational Linguistics.
- 253 Rémi Lebret, David Grangier, and Michael Auli. 2016. Neural text generation from structured data
254 with application to the biography domain. In *Proceedings of the 2016 Conference on Empirical*
255 *Methods in Natural Language Processing*, pages 1203–1213, Austin, Texas. Association for
256 Computational Linguistics.
- 257 Ankur P Parikh, Xuezhi Wang, Sebastian Gehrmann, Manaal Faruqui, Bhuvan Dhingra, Diyi Yang,
258 and Dipanjan Das. 2020. Totto: A controlled table-to-text generation dataset. In *Proceedings of the*
259 *2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association
260 for Computational Linguistics.
- 261 Dragomir Radev, Rui Zhang, Amrit Rau, Abhinand Sivaprasad, Chiachun Hsieh, Nazneen Fatema
262 Rajani, Xiangru Tang, Aadit Vyas, Neha Verma, Pranav Krishna, et al. 2020. Dart: Open-domain
263 structured data record to text generation. *arXiv preprint arXiv:2007.02871*.