

---

# SciGen: a Dataset for Reasoning-Aware Text Generation from Scientific Tables

---

Nafise Sadat Moosavi<sup>1</sup>, Andreas Rücklé<sup>1,\*</sup>, Dan Roth<sup>2</sup>, Iryna Gurevych<sup>1</sup>

<sup>1</sup>Ubiquitous Knowledge Processing Lab (UKP Lab)  
Department of Computer Science, Technical University of Darmstadt  
<https://www.ukp.tu-darmstadt.de>

<sup>2</sup>Department of Computer and Information Science, UPenn  
<https://www.seas.upenn.edu/~danroth>

## 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.<sup>2</sup> Table 1 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 against imperfect descriptions.

<u>BLEU</u>	<u>METEOR</u>	<u>MoverS</u>	<u>BertS</u>	<u>BLEURT</u>
66.73	0.614	0.98	0.99	0.95

Table 1: The automatic evaluation scores for 50 table-descriptions in which an expert has written descriptions based on table contents and without looking at the gold descriptions.

## B Impact of Table Captions

Table 2 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 few-shot setting where the captions of the tables were excluded from the input data.

## C Example outputs from all baselines

Table 3 provides the output of all the examined baselines for the table in Figure 2.

## D Examples from expert vs. automatically extracted annotations

Below we present a set of examples for comparing expert and automatically extracted description:

---

\*Contributions made prior to joining Amazon.

<sup>2</sup>They had access to the background sections of the article in case the information of the table itself is not enough to describe them.



while the throughput on the linear dataset is the lowest. This trend occurs because the maximum possible execution concurrency of a tree is affected by the balancedness of the tree. A full binary tree of  $N$  cells can be processed concurrently with at most  $N+12$  threads, because all  $N+12$  leaf nodes are mutually independent. On the other hand, an extremely unbalanced binary tree can be processed with only one or two threads at a time due to the linearity of the tree. As a result, our implementation can train input data of balanced trees with greater throughput than input data of unbalanced trees. Resource Utilization. Another interesting fact in Table 1 is that the training throughput on the linear dataset scales better than the throughput on the balanced dataset, as the batch size increases. For the balanced dataset, the recursive implementation efficiently utilizes many threads to process the data even at a small batch size of 1, and thus increasing the batch size leads to a relatively small speed boost. On the contrary, for the linear dataset, the recursive implementation fails to efficiently make use of CPU resources and thus the performance gain provided by increasing the batch size is relatively high.

**Expert:** (Vemulapalli and Agarwala, 2019) Table 3 shows the triplet prediction accuracy of median rater, FECNet-16d and AFFNet-CL-P for each triplet type in the FEC test set. [CONTINUE] the performance is best (85.1%) for two-class triplets, and is lowest (77.1%) for one-class triplets.

**Automatic:** (Vemulapalli and Agarwala, 2019) Table 3 shows the triplet prediction accuracy of median rater, FECNet-16d and AFFNet-CL-P for each triplet type in the FEC test set. As expected, the performance is best (85.1%) for two-class triplets, which are relatively the easiest ones, and is lowest (77.1%) for one-class triplets, which are relatively the most difficult ones.

## References

- Eunji Jeong, Joo Seong Jeong, Soojeong Kim, Gyeong-In Yu, and Byung-Gon Chun. 2018. Improving the expressiveness of deep learning frameworks with recursion. In *Proceedings of the Thirteenth EuroSys Conference*, pages 1–13.
- Raviteja Vemulapalli and Aseem Agarwala. 2019. A compact embedding for facial expression similarity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5683–5692.