

SciRIFF: A Resource to Enhance Language Model Instruction-Following over Scientific Literature

Anonymous ACL submission

Abstract

We present SciRIFF (Scientific Resource for Instruction-Following and Finetuning), a dataset of 137K instruction-following instances covering 54 tasks. These tasks span five core scientific literature understanding capabilities: information extraction, summarization, question answering, claim verification, and classification. SciRIFF is unique in being the only all expert-written, high-quality instruction-following dataset designed for extracting and synthesizing information from research literature across diverse scientific fields. It features complex instructions with long input contexts, detailed task descriptions, and structured outputs. To demonstrate its utility, we finetune a series of large language models (LLMs) using a mix of general-domain and SciRIFF instructions. On nine out-of-distribution held-out tasks (referred to as SciRIFF-EVAL), LLMs finetuned on SciRIFF achieve 70.6% average improvement over our baselines trained only on general-domain instructions. SciRIFF facilitates the development and evaluation of LLMs to help researchers navigate the rapidly growing body of scientific literature.

1 Introduction

Large language models (LLMs) have the potential to advance scientific progress by helping researchers navigate and draw insights from the scientific literature. To accomplish these tasks, LLMs must be able to reliably follow a range of *instructions*—e.g. to extract information, summarize content, or answer questions—when given research articles as input. These instructions are often grounded on entire scientific articles, featuring longer inputs than other typical instruction-following resources in science domain. In addition, the model’s responses may need to be *structured* according to a specific format or schema that supports aggregation for literature review (Marshall and Wallace, 2019), or is consumable

by software components like augmented reading interfaces (Lo et al., 2023; Palani et al., 2023). For example, when analyzing clinical trials, responses should follow a PICO framework (Population, Intervention, Comparison, Outcome), or when examining methodology papers, follow a standardized format capturing study design, sample size, statistical methods, and key findings. Additionally, outputs must be machine-readable, such as JSON formats that capture relationships between entities (e.g., protein-protein interactions in biochemistry papers) or structured evidence for claims (e.g., “claim”: “Coffee consumption reduces diabetes risk”, “evidence”: [“Study A shows 23% risk reduction”, “Meta-analysis B confirms protective effect”], “confidence”: “moderate”).

While bespoke models are available for specific scientific literature understanding tasks, models that can flexibly follow instructions in domain-specific settings of science are preferable both for their ease of use (offering a unified input / output interface) and for their ability to generalize to novel applications and settings within that domain.

The general instruction-following capabilities of LLMs have advanced rapidly in recent years, largely due to the availability of general-purpose instruction datasets (Zhang et al., 2023a). In addition, some instruction-following resources are available for specific scientific and medical tasks, such as describing the properties of a molecule (Fang et al., 2024; Yu et al., 2024) or answering medical exam questions (Toma et al., 2023; Han et al., 2023) (see §5 for a review). However, few resources are available for supporting instruction-following for flexible scientific literature understanding capabilities across a range of domains.

In this work, we present SciRIFF (Scientific Resource for Instruction-Following and Finetuning), a comprehensive dataset to enable progress on instruction-following over

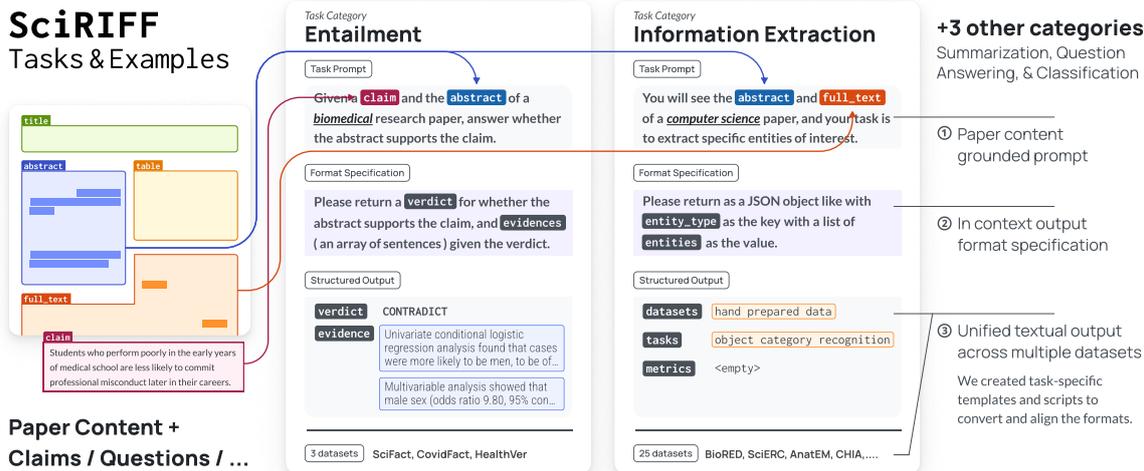


Figure 1: Example SciRIFF tasks. Given an input context from a research paper, the **text prompt** instructs an LLM to perform an operation on the input—e.g. determine whether the **abstract** entails a scientific **claim**, extract information over the **full text**, answer a question, etc. The model’s **output** must conform to a task-specific, user-specified **structure**. SciRIFF unifies 54 scientific literature understanding tasks under a common input / output format, enabling the development of LLMs that can flexibly generalize to novel scientific use cases.

scientific literature. SciRIFF includes 137K demonstrations for 54 tasks spanning five broad scientific literature understanding task categories: information extraction, summarization, question answering, claim verification, and classification.

SciRIFF covers five scientific domains, ranging from artificial intelligence to clinical medicine (Figure 2). The tasks in SciRIFF are derived from challenging scientific literature understanding datasets with real-world relevance, all of which include human-annotated inputs and responses. We opted for organic, human-annotated data rather than the synthetic or LLM-distilled instruction-following data explored in recent works (e.g., Lambert et al., 2024a), because human annotations better capture the nuanced domain expertise, complex structures and reasoning required for scientific tasks, while also providing a more reliable ground truth for evaluating model performance on real-world scientific problems.

Our resource is a unique and specialized instruction-following dataset. As illustrated in Figure 1 and with sample prompt templates provided in Appendix C, it is characterized by: (1) grounding every instance in scientific articles or texts, (2) requiring structured and complex responses, such as answers paired with attributions (i.e., tracing the source of the answer), and (3) featuring longer input contexts compared to most existing resources in the science domain (see Figure 5 and Table 5 in the Appendix). Notably, all our instruction tem-

plates are created by human experts to ensure high quality.

We also present a new benchmark dataset for evaluating instruction-following capabilities of LLMs in the science domain. We particularly hold out a subset of SciRIFF for use as an evaluation benchmark which covers a representative range of skills and tasks, which we call SciRIFF-EVAL (§3.1). To demonstrate the utility of SciRIFF in improving scientific literature instruction following, we perform supervised finetuning experiments on several LLMs ranging different sizes¹. When finetuned on a mix of SciRIFF and general open-source instruction-following data (i.e., TULU v2 (Iverson et al., 2023a)), our models show consistent improvements on SciRIFF-EVAL compared to training on general-domain instructions alone, even though these evaluation tasks test true out-of-distribution generalization with formats and templates entirely excluded from training.

In summary, our contributions are as follows:

- We introduce SciRIFF, an expert-crafted, high-quality, and comprehensive instruction-following dataset with 137K instances covering a wide range of tasks spanning five scientific domains. Many tasks in SciRIFF feature long input contexts and require structured and complex model responses.
- We supervise finetune a range of LLMs on SciRIFF and demonstrate its effectiveness on

¹Other types of post-training such as preference optimization are outside our scope.

instruction-following in scientific literature understanding.

- We release the SciRIFf dataset, our evaluation suite SciRIFf-EVAL, and the model checkpoints, and code to reproduce our results.

2 SciRIFf

SciRIFf is a comprehensive instruction-tuning resource for real-world scientific literature understanding, consisting of 137k instructions and spanning five broad task categories and five subject domains (Figure 1). Our focus is on *document-grounded* scientific literature understanding tasks, rather than tasks that evaluate scientific knowledge recall (Feng et al., 2024), or general mathematical, reasoning-related problem-solving abilities without reference to scientific literature (e.g., Sci-Instruct (Zhang et al., 2024a), MMLU (Hendrycks et al., 2021a)). In addition to coverage of a wide range of tasks, the instructions in SciRIFf often are grounded in long inputs (i.e., scientific papers), and they support *structured* outputs according to a specific schema useful for tasks in literature understanding (such as relation extraction, fact checking with rationale selection, QA with attribution, etc). The instances in SciRIFf are sourced from existing high-quality scientific datasets and converted into instructions using human expert-written and verified instruction templates. Out of 54 tasks, 50 involve templates paired with manually crafted Python scripts. These scripts serve to extract ground-truth answers, postprocess (e.g., removing duplicate named entity mentions), and normalization on the source datasets. This includes transforming raw data, such as span-level formats for encoder models, into instruction-following formats.

2.1 Dataset construction

We construct SciRIFf by sourcing from existing, high-quality scientific literature understanding datasets for instruction-following—drawing inspiration from canonical resources like Flan (Longpre et al., 2023) and Super-NaturalInstructions (Wang et al., 2022). We then ask experts to write high-quality and carefully-vetted instruction templates which will convert original dataset instances to SciRIFf instances. We chose this approach rather than the alternative recent trend of generating synthetic data using an LLM (e.g., (Köksal et al., 2023; Li et al., 2023)). We believe it is sensible to ex-

haust available human-annotated resources, which we can be fairly confident are correctly-annotated, before turning to potentially noisy synthetic data (See Appendix C). For the same reason, we would need high-quality evaluation data, which we construct by holding out nine SciRIFf tasks as an evaluation benchmark (§3.1). We hope our resource will provide valuable signals for future synthetic data generation efforts.

Dataset selection criteria In forming SciRIFf, we focus on scientific literature understanding tasks in which the model is given a portion of scientific text as input, and is instructed to produce a response derived directly from the text. The task families include summarization, reading comprehension, information extraction and other tasks, and are the most relevant setting for real-world use cases (e.g., meta-analysis of literature, clinical decision-making, augmented reading). We *exclude* datasets that require retrieval from document collections (e.g., open-domain QA), since it’s unclear how to build instruction-response pairs from them. We also exclude datasets that assess general reasoning and mathematical problem-solving skills without necessarily grounding on scientific literature, such as ScienceQA (Lu et al., 2022), SciBench (Wang et al., 2023b), and MATH (Hendrycks et al., 2021b) since such resources already exist. Additionally, we only keep datasets that are publicly available, have a permissive license, and are well-documented and actively maintained. See Appendix A.1 for the complete task list.

Instruction templates Our approach to use expert-designed templates² ensures that the instructions are tailored to each dataset’s unique characteristics, maximizing the quality and relevance of the resulting dataset. The expert annotators³ repurposed tasks for natural instruction-following using custom scripts, cleaned and refined outputs (such as deduplicating named entities), and augmented the templates with additional dataset metadata to enhance usability. We use json as the common output format for *all* structured tasks, which facilitates consistent evaluation and aligns with industry trends to request JSON outputs. Instruction templates are written in Jinja (Pallets, 2024).

²We opted for a single template per dataset. Our initial experiments with LLM-generated templates proved unsatisfactory, as they often resulted in vague or noisy instructions and failed to clearly specify the desired output format.

³The authors of the paper.

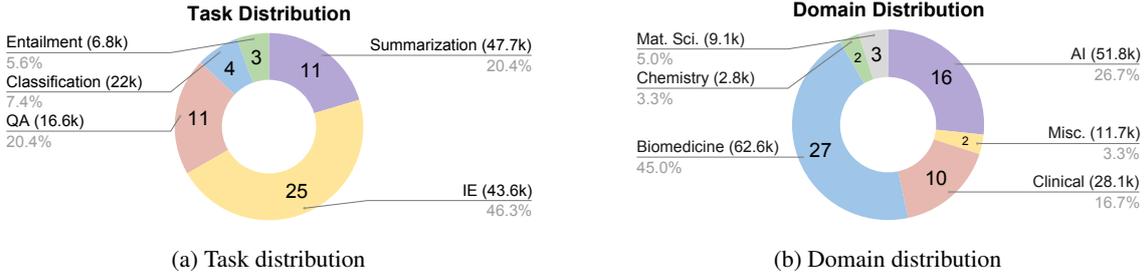


Figure 2: Distribution of task categories and domains in SciRIFF. The numbers in the pie charts indicate the number of datasets present in every task category/domain, while the numbers in brackets indicate the total number of instances per task category/domain.

Quality Verification. To ensure high-quality, each template was verified by an additional annotator for clarity and correctness. We include guidelines and best practices for prompt-writing in the release and aim to promote community contributions for actively expanding SciRIFF through our open-sourced data collection process.

2.2 Instruction Mix Statistics

Figures 2a and 2b present an overview of the SciRIFF training set distribution over task categories and domains respectively. The domain distribution reflects the current landscape of available high-quality scientific datasets (e.g., Reid et al., 2022), with a notable representation from the biomedicine and AI domain. This aligns with our dataset selection criteria, which prioritize well-documented resources with permissive licenses. Given the significant presence of information extraction tasks, a large percentage of datasets in SciRIFF (34 datasets; 63%) require structured outputs.

We construct three instruction mixes from this dataset collection, with maximum context lengths (input + output tokens) of 4,096, 8,192 and 16,382 per instance (longer instances are truncated where possible and discarded otherwise; see Appendix A.3). Due to model and hardware limitations, we conduct experiments in this work using the SciRIFF-4096 mixture, and make the longer mixtures available to enable future research. In what follows, we refer to SciRIFF-4096 simply as SciRIFF.

3 Experiment setup

We conduct supervised finetuning experiments to evaluate the effectiveness of SciRIFF in improving LLM performance on scientific instruction-following tasks across various model families and sizes. Our experiments explore different

data configurations and their impact on scientific instruction-following as measured through SciRIFF-EVAL described in §3.1.

3.1 Evaluation

SciRIFF-EVAL We selected a set of 9 tasks from SciRIFF for evaluation, designed to cover a diverse range of task categories and scientific domains. SciRIFF tests true out-of-distribution generalization with formats and templates entirely excluded from training. The inputs, outputs, and evaluation metrics for each task are detailed in Table 1. Additional details are included in Appendix E.

3.2 Scientific Instruction Finetuning

Using SciRIFF, our goal is to adapt pre-trained LLMs to the scientific literature domain. We conduct supervised fine-tuning experiments using a range of models and data configurations to assess the effectiveness of SciRIFF. In §4.2, we present an additional analysis examining the potential of using SciRIFF for further finetuning of instruction-tuned models, exploring a compute-efficient strategy for adapting generic models to scientific literature domains.

Data sources We finetune using two primary datasets: (1) SciRIFF, and (2) TüLU V2 Mix (Iverson et al., 2023b), an open-source high-quality general-domain instruction-following dataset that includes demonstrations from various sources, both human-written (e.g., Flan (Wei et al., 2022)) and distilled from proprietary LLMs (e.g., ShareGPT⁴, Open Assistant⁵). The original TüLU V2 Mix contains 326,154 examples, including 7.5K scientific literature understanding demonstrations which overlap (i.e. contaminated) with our evaluation set

⁴<https://sharegpt.com/>

⁵<https://github.com/LAION-AI/Open-Assistant>

Name	Type	Input	Output	Metrics
BioASQ List QA	QA	Question, paper excerpts	Answer entities	Exact match F1
BioRED	IE(NER)	Biomedical abstract	6 entity types	Exact match F1
DiSCoMaT	IE(Table)	LaTeX table excerpt	Table entries	BLEU score
MultiCite (MC)	Classification	Citation context	Citation intents	Exact match F1
MuP	Summarization	ML paper full text	Peer review summary	LLM judge similarity
Qasper	QA	NLP paper question	Answer / Attribution	LLM judge similarity / Token F1
SciERC	IE(Rel)	CS abstract	6 entity types	Exact match F1
SciFact	Entailment	Claim, abstract	Verdict / Evidence	Label F1 / Token F1

Table 1: Evaluation tasks included in SciRIFF-EVAL. “/” separators indicate two separate subtasks. We use GPT-4o as our LLM judge and evaluate similarity on a 1-5 scale; see Appendix E for details.

SciRIFF-EVAL. We remove those 7.5K examples for clean experiments and to avoid contamination with SciRIFF-EVAL. For all experiments, we consistently use this filtered version and refer to this as TüLU V2 MIX to maintain controlled finetuning and unbiased evaluations.

Base models We use several base LLMs as starting points for our finetuning experiments: Llama 3.1-8B (Touvron et al., 2023), Llama 3.2-3B (Dubey et al., 2024), and Qwen 2.5-1.5B (Yang et al., 2024).⁶ While our primary focus is on improving base models, we also experiment with models that have undergone proprietary instruction tuning and preference optimization (“-instruct” versions) (Ouyang et al., 2022). Although direct comparisons with these proprietary posttrained models are complicated by unknown training details, we show that SciRIFF can provide additional value even in these cases. We note, however, that our main results and analyses focus on the controlled experiments with base models where we can fully account for all training conditions.

Finetuning data configurations For each model, we explore three data configurations: (1) TüLU V2 MIX only, to establish a **baseline** for general instruction-following; (2) SciRIFF only, to assess the impact of scientific instruction data in isolation; and (3) SciRIFF+TüLU, combining the general and scientific instruction data.

4 Results

This section discusses our key results and findings.

4.1 Main Results

We report our main experimental results in Table 2. For fair comparison, all models are finetuned on the

⁶Due to compute resource constraints, we were not able to train larger models for all the model families. However, as shown in §4.1 improvements are consistent across sizes/families.

same data mixes. We show that including SciRIFF instances results in the best average performance in each model family. We report GPT-4o and GPT-4o-mini as strong baselines. Our key observations and findings are below:

SciRIFF enhances scientific literature understanding Table 2 shows that finetuning on SciRIFF consistently enhances the overall performance on SciRIFF-EVAL. Compared to the corresponding base models finetuned on Tülu, SciRIFF-trained models achieve, on average, 70.6% performance gain. For example, the Qwen 2.5 1.5B model improves from 29.1 to 57.2 in average score with SciRIFF alone, and further to 59.1 with SciRIFF+Tülu. Furthermore, without exception, SciRIFF also adds values when used to finetune “-instruct” versions of models (44.6% on average). Across all model groups, the “-instruct” variants trained exclusively on SciRIFF achieve the highest average scores within their respective groups. Finally, while GPT-4o and GPT-4o-mini serve as strong proprietary baselines with average scores of 60.4 and 58.0 respectively, out of the twelve models trained with the inclusion of SciRIFF instances, nine outperform GPT-4o on SciRIFF-EVAL, with Qwen 2.5 1.5B showing the most significant improvement. These results indicate that our specialized SciRIFF can substantially enhance scientific literature understanding and extraction capabilities beyond what general or proprietary instruction data provides.

Task-specific impacts and room for improvement Beyond the overall improvements, SciRIFF training achieves particularly large gains on the three IE tasks (BioRED, DiSCoMaT, and SciERC). Relative to their Tülu-only counterparts, SciRIFF-finetuned base models improve IE performance by, on average, 200.4%, while “-instruct” models see an average 139.8% improvement. In

Model	Data	BioASQ	BioR	DiscMT	EI	MC	MuP	Qasper	SciERC	SciFact	Avg.
GPT-4o	-	48.3	63.6	71.3	25.9	62.0	88.3	54.0 / 55.0	40.3	85.5 / 70.4	60.4
GPT-4o-mini	-	49.6	53.7	75.6	27.7	54.8	88.8	61.7 / 46.7	33.1	82.7 / 63.6	58.0
Qwen 2.5 1.5B--Instruct	-	38.9	19.7	35.5	10.5	36.9	80.8	38.8 / 39.4	20.8	55.0 / 31.5	37.1
	SciRIFF	48.1	79.7	80.6	20.9	70.9	67.3	42.8 / 54.3	52.0	80.9 / 68.9	60.6
	SciRIFF+Tülu	49.3	80.1	79.5	21.3	70.8	61.3	45.8 / 48.6	51.0	78.5 / 70.1	59.7
Qwen 2.5 1.5B	Tülu	35.7	23.4	31.8	7.6	6.6	73.0	25.0 / 23.2	12.0	52.4 / 29.5	29.1
	SciRIFF	43.6	81.8	45.6	18.9	71.2	67.8	47.0 / 51.4	52.7	78.8 / 70.5	57.2
	SciRIFF+Tülu	46.5	79.0	78.3	19.4	70.2	63.8	40.4 / 49.7	51.7	80.9 / 70.6	59.1
Llama 3.2 3B--Instruct	-	42.9	35.9	61.0	11.2	47.3	86.0	43.9 / 35.8	20.8	59.5 / 40.0	44.0
	SciRIFF	42.7	84.0	83.4	25.5	71.4	64.8	50.0 / 57.1	58.2	86.8 / 70.5	63.1
	SciRIFF+Tülu	43.0	83.3	82.9	21.7	72.2	69.0	51.9 / 58.2	53.3	85.6 / 70.3	62.8
Llama 3.2 3B	Tülu	35.5	30.1	46.7	3.1	44.0	75.6	47.4 / 34.4	20.3	55.4 / 36.6	39.0
	SciRIFF	43.6	84.2	83.2	25.2	71.7	64.3	46.0 / 57.2	57.2	81.6 / 65.8	61.8
	SciRIFF+Tülu	46.0	84.3	83.3	24.6	72.7	65.5	47.7 / 56.3	57.0	82.7 / 71.2	62.8
Llama 3.1 8B--Instruct	-	43.7	48.8	62.2	17.8	48.8	88.3	54.0 / 43.0	30.6	66.7 / 51.7	50.5
	SciRIFF	45.9	86.0	83.7	25.0	71.4	70.5	53.3 / 54.1	56.8	85.8 / 72.5	64.1
	SciRIFF+Tülu	48.8	84.7	83.6	26.6	71.3	66.0	50.9 / 55.2	54.4	85.5 / 70.2	63.4
Llama 3.1 8B	Tülu	44.4	42.8	51.8	1.1	39.4	80.8	42.8 / 28.6	24.3	50.0 / 33.6	40.0
	SciRIFF	46.2	84.2	83.9	23.5	71.0	68.5	49.8 / 52.2	56.2	83.3 / 71.9	62.8
	SciRIFF+Tülu	41.6	85.2	78.7	28.2	71.6	70.5	47.9 / 61.0	58.1	87.4 / 71.2	63.8

Table 2: Performance on SciRIFF-EVAL tasks across model families and training configurations. For base models, we compare SciRIFF against TüLU on scientific understanding tasks. GPT-4o and GPT-4o-mini serve as strong baselines. Best performance per model group is **bolded**. Columns with a “/” indicate two evaluation metrics as described in §3.1.

contrast, performance on the summarization task (MuP) generally shows decreases after SciRIFF finetuning. This suggests that while SciRIFF is particularly effective for enhancing IE capabilities, it may not provide additional benefits for summarization tasks that are likely well-covered in general instruction-following. Additionally, we hypothesize that MuP’s evaluation approach (LLM-as-Judge against a gold reference) might penalize stylistic or structural shifts in answers that deviate from the reference. The fact that GPT-4o and our strong finetuned models achieve only an average score of around 60 highlights the difficulty of SciRIFF-EVAL. Model performance remains relatively low on tasks like EI; This is due to a combination of task difficulty and evaluation challenges, which we discuss in §6.

Balancing scientific and general data As shown in Table 2, combining SciRIFF and TüLU V2 Mix training data (SciRIFF+TüLU) yields the best performance on SciRIFF-EVAL for *base* models. This suggests that incorporating general instruction-following data may provide some broader capability transfer, which base models particularly benefit

from, though the impact remains limited (within 2.2%). On the other hand, for “--instruct” models training exclusively on SciRIFF data proves more effective, but similarly trivially (within 1% on average).

Checkpoint	SciRIFF	TüLU-V2	Total
Llama 2 base	35,357	318,686	354,043
TüLU V2	35,357	35,357	70,714

Table 3: SciRIFF and TüLU V2 Mix instances used for finetuning described in §4.2, with $n_{sci} = 1000$.

4.2 Continual Finetuning Analysis

In early phase of our experiments, we study strategies for efficient adaptation. Specifically, we examined whether starting from an existing instruction-tuned checkpoint (on general instruction-following data) could provide compute advantages over training from scratch, without hurting SciRIFF-EVAL performance. For this controlled experiment, we selected two starting points: (1) Llama 2 base and (2) the same model already finetuned on science-decontaminated TüLU V2 Mix (referred to as TüLU V2).

Model	Data	7B	Sci. 70B
Llama 2	TÜLU	36.7	47.5
	SciRIFF	48.0	51.1
	SciRIFF+TÜLU	46.0	50.8
TÜLU V2	SciRIFF	47.0	48.8
	SciRIFF+TÜLU	47.0	50.7

Table 4: Comparison of SciRIFF-EVAL (Sci.) performance for models finetuned from Llama 2 base and TÜLU V2 (science-decontaminated), at both 7B and 70B scales. Best performance in each group is **bolded**.

We explored different training approaches for each starting point. For Llama 2 base, we train on all available TÜLU V2 Mix demonstrations, combined with 1000 instances per SciRIFF task, given the *empirical findings* in Section 4.3. For the TÜLU V2 starting point, we perform continual finetuning using 1000 instances per SciRIFF task, together with a *matching number* (1000) of instances sampled from TÜLU V2 Mix.

Table 4 reports average SciRIFF-EVAL performance for our two starting checkpoints using three data configurations. Despite using a smaller total amount of training data, starting from TÜLU V2 performs comparably to Llama 2 base while requiring only 20% of the compute. When trained on SciRIFF+TÜLU data, models from both starting points achieve similar performance: TÜLU V2 is slightly better on science at 7B and nearly identical at 70B. Given that finetuning TÜLU V2 requires only 20% of the data (Table 3), this suggests a potential for compute-efficient adaptation of strong instruction-following models to scientific domains. This aligns with findings from prior works, e.g. Dong et al. (2024); Shi et al. (2023). While our main experiments (§3.2) now use other model architectures⁷, this analysis demonstrates how practitioners can study the learning behavior of a specific model family to identify compute-efficient adaptation strategies. When combined with the analysis in the following Section 4.3, these insights can help determine optimal training configurations for SciRIFF when working with a fixed model architecture.

4.3 Instruction Data Scale

We define n_{sci} as the number of instances per SciRIFF task. Figure 3 shows that performance on SciRIFF-EVAL increases sharply as n_{sci} rises from 100 to 500 and levels off subsequently. We found

⁷Due to compute limitation, we do not perform similar continual finetuning analysis on all other models.

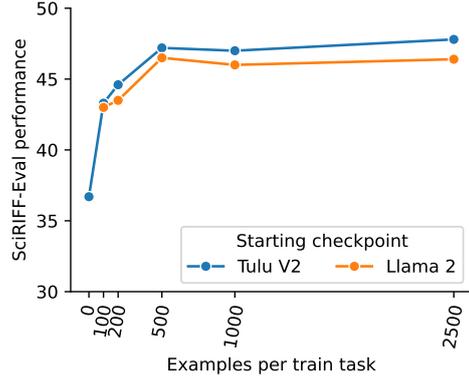


Figure 3: Performance on SciRIFF-EVAL as a function of n_{sci} , the number of science instances per task. Performance gains largely saturate by $n_{sci} = 1000$. Experiments are done on Llama 2 and our TÜLU V2 7B models. We extrapolate the finding to settings in §3.2.

that 1,000 instances per science task are sufficient for peak performance for Llama 2 models.

Therefore, we set $n_{sci} = 1000$ across our experiments in the continual finetuning analysis (§4.2) based on these findings.

5 Related Work

Strategies for creation of instruction-following resources. Instruction tuning has emerged as a crucial technique for enhancing generalizability and controllability of LLMs (Wei et al., 2022; Sanh et al., 2022; Mishra et al., 2022; Ivison et al., 2023b). Instruction-following resource creation strategies have been explored, such as repurposing existing datasets using human-written instruction templates (Wei et al., 2022; Chung et al., 2024; Sanh et al., 2022), crowdsourcing instructions [Databricks (2023); Zhou et al. (2023); Mishra et al. (2021), ShareGPT⁸] and generating synthetic data. As LLM capabilities grow, synthetic instruction generation approaches, often including humans in the loop as correctors, have shown promising results. Broadly, these approaches use LLMs to either generate new dataset/task instances alongside instructions (Wang et al., 2023c; Xu et al., 2024; Nayak et al., 2024; Lou et al., 2024), or to “back-translate” existing datasets into instructions (Yin et al., 2023; Köksal et al., 2023; Li et al., 2023). In this work, we create instructions using human-written templates (§2.1) for quality assurance.

Instruction-following resources for scientific literature. While numerous open-domain

⁸<https://sharegpt.com/>

instruction-following collections exist, resources for enhancing and evaluating LLMs’ instruction-following capabilities on scientific literature are limited. Such resources are critical for developing models that can assist researchers and accelerate scientific discovery (Taylor et al., 2022; Xie et al., 2023). Recent work has taken steps in this direction with the development of instruction-following datasets for specific domains such as mathematics (Yue et al., 2024a,b; Shao et al., 2024; Luo et al., 2023; Tang et al., 2024; Toshniwal et al., 2024), medicine (Parmar et al., 2022; Wu et al., 2024; Rohanian et al., 2023), chemistry (Yu et al., 2024; Zhang et al., 2024b), molecular biology (Fang et al., 2024; Tran et al., 2023), materials science (Song et al., 2023), and college-level foundational science (Zhang et al., 2024a). Besides domain limitations, these resources primarily focus on improving LLMs’ abilities to solve college-level science problems or reasoning tasks (see also, MMLU (Hendrycks et al., 2021a), SciEval (Sun et al., 2023), TheoremQA (Chen et al., 2023), SciBench (Wang et al., 2023b), and GPQA (Rein et al., 2023)). In contrast, SciRIFf both covers a broader set of scientific domains and focuses on document-grounded scientific literature understanding tasks that can power real-world scientific use cases. Another distinguishing factor of our work is our inclusion of tasks that require structured outputs, following a uniform JSON output format, besides text-to-text tasks. Some instruction-tuning resources have explored structured output formats (Zhang et al., 2023b; Wang et al., 2023a; Jiao et al., 2023; Gao et al., 2023), but not with a focus on scientific literature. Finally, most datasets in SciRIFf require long input contexts, leading to longer instruction contexts than prior works (see Appendix Table 5 for a comparison).

Other scientific literature benchmarks. Prior works have developed benchmarks to improve and assess scientific literature understanding. Notable efforts in the biomedical domain include BLUE (Peng et al., 2019), BLURB (Gu et al., 2021), InBoXBART (Parmar et al., 2022), and BigBio (Fries et al., 2022); SciRIFf covers a broader set of domains than these resources. Other efforts such as SciRepEval (Singh et al., 2023), Galactica (Taylor et al., 2022), and AcademicGPT (Wei et al., 2023) cover domains beyond biomedicine, but are not suitably formatted for training or

evaluating instruction-following models. SciASSESS (Cai et al., 2024) evaluates LLMs’ proficiency in scientific literature analysis, focusing on tasks like memorization and reasoning. Li et al. (2024) introduce a hybrid strategy that combines continual pretraining and supervised finetuning to specialize LLMs for scientific literature understanding, along with SciLitIns - a *synthetically* generated instruction dataset. In contrast, SciRIFf provides fully *expert-written* instructions with structured outputs, serving both as a benchmark and training resource for advancing LLMs in scientific literature tasks and downstream applications.

6 Conclusion and Future Work

In this work, we introduced SciRIFf, a resource to facilitate progress on LLM instruction-following over scientific literature. We demonstrated that training on SciRIFf leads to significant improvement of model performance on held-out scientific tasks (on average 70.6% over baselines), with especially large improvements on tasks requiring structured extraction or attribution.

As observed in §4.1, neither our best-performing fine-tuned models nor the GPT-4o on SciRIFf-EVAL is sufficiently strong (around 60%). This is partly due to the difficulty of the tasks, but also due the challenges associated with evaluating structured LLM responses in cases where the predicted surface form does not match the reference, but the underlying meaning is the same (Wadhwa et al., 2023). Utilizing LLMs to perform more flexible and fine-grained evaluations (Kim et al., 2024) represents a promising direction. Future work could focus on reliably generating multiple templates for such complex tasks in a more controlled and principled manner to help models learn from a richer set of demonstrations and improve their generalization to unseen tasks. Incorporating reliable synthetic data generation techniques and preference data (Lambert et al., 2024b) for scientific literature understanding tasks is also a promising avenue. In conclusion, we are optimistic that the SciRIFf data and evaluations, as well as the model checkpoints, will serve as valuable resources to build systems which can boost the productivity of scientific researchers.

587 **Limitations**

588 While we demonstrated the effectiveness of SciR-
589 IFF and the value of SciRIFF-EVAL, we note the
590 following limitations about our work: Although our
591 resource is constructed through human-annotated
592 templates, the source of our data is existing la-
593 beled datasets from scientific literature understand-
594 ing and synthesizing tasks. Although we included
595 a wide range of datasets, this still could limit the
596 open-ended tasks that could involve literature un-
597 derstanding. Additionally, more sophisticated it-
598 erative or conversational interactions mimicking
599 interactions with a research assistant are not cap-
600 tured with SciRIFF. Furthermore, the largest model
601 our compute could afford to finetune on is 8B in
602 size. We suspect that training larger open-source
603 models (e.g., Llama 3.1 405B) can close the gap.

604 **Ethical Statement**

605 The ethical risks associated with this work are min-
606 imal. As we source the data from existing datasets
607 and we work in the science domain, we do not sus-
608 pect major risks are involved in the curation of our
609 dataset. However, potential biases might still exist
610 in some datasets. For example, one of the source
611 datasets is paper summarization which is sourced
612 from OpenReview.net peer reviews by the original
613 authors. And peer reviews might inherently occa-
614 sionally include biases or unhelpful languages. As
615 with all LLMs, our trained models are still prone
616 to issues such as hallucinations, so users should
617 exercise caution when interpreting model outputs,
618 particularly in downstream applications in science,
619 and verify any generated content for accuracy and
620 relevance.

621 **References**

622 Hengxing Cai, Xiaochen Cai, Junhan Chang, Sihang
623 Li, Lin Yao, Changxin Wang, Zhifeng Gao, Hong-
624 shuai Wang, Yongge Li, Mujie Lin, Shuwen Yang,
625 Jiankun Wang, Mingjun Xu, Jin Huang, Fang Xi, Ji-
626 axi Zhuang, Yuqi Yin, Yaqi Li, Changhong Chen,
627 Zheng Cheng, Zifeng Zhao, Linfeng Zhang, and
628 Guolin Ke. 2024. Sciassess: Benchmarking llm profi-
629 ciency in scientific literature analysis. *arXiv preprint*
630 *arXiv: 2403.01976*.

631 Wenhu Chen, Ming Yin, Max Ku, Pan Lu, Yixin Wan,
632 Xueguang Ma, Jianyu Xu, Xinyi Wang, and Tony
633 Xia. 2023. TheoremQA: A theorem-driven question
634 answering dataset. In *Proceedings of the 2023 Con-*
635 *ference on Empirical Methods in Natural Language*

Processing, pages 7889–7901, Singapore. Associa-
tion for Computational Linguistics.

Hyung Won Chung, Le Hou, Shayne Longpre, Barret
Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi
Wang, Mostafa Dehghani, Siddhartha Brahma, Al-
bert Webson, Shixiang Shane Gu, Zhuyun Dai,
Mirac Suzgun, Xinyun Chen, Aakanksha Chowdh-
ery, Alex Castro-Ros, Marie Pellat, Kevin Robinson,
Dasha Valter, Sharan Narang, Gaurav Mishra, Adams
Yu, Vincent Zhao, Yanping Huang, Andrew Dai,
Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Ja-
cob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le,
and Jason Wei. 2024. Scaling instruction-finetuned
language models. *Journal of Machine Learning Re-*
search, 25(70):1–53.

Databricks. 2023. [Free dolly: Introducing the world’s
first truly open instruction-tuned llm](#). Blog post.

Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng
Li, Mingfeng Xue, Dayiheng Liu, Wei Wang, Zheng
Yuan, Chang Zhou, and Jingren Zhou. 2024. How
abilities in large language models are affected by
supervised fine-tuning data composition.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
Akhil Mathur, Alan Schelten, Amy Yang, Angela
Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang,
et al. 2024. The llama 3 herd of models. *arXiv*
preprint arXiv: 2407.21783.

Yin Fang, Xiaozhuan Liang, Ningyu Zhang, Kangwei
Liu, Rui Huang, Zhuo Chen, Xiaohui Fan, and Hua-
jun Chen. 2024. Mol-instructions: A large-scale
biomolecular instruction dataset for large language
models. In *International Conference on Learning*
Representations (ICLR).

Kehua Feng, Keyan Ding, Weijie Wang, Xiang Zhuang,
Zeyuan Wang, Ming Qin, Yu Zhao, Jianhua Yao,
Jiang Zhang, and Huajun Chen. 2024. [Sciknoweval:](#)
[Evaluating multi-level scientific knowledge of large](#)
[language models](#). *Preprint*, arXiv:2406.09098.

Jason Fries, Leon Weber, Natasha Seelam, Gabriel Al-
tay, Debajyoti Datta, Samuele Garda, Sunny Kang,
Rosaline Su, Wojciech Kusa, Samuel Cahyawijaya,
Fabio Barth, Simon Ott, Matthias Samwald, Stephen
Bach, Stella Biderman, Mario Sanger, Bo Wang,
Alison Callahan, Daniel Leon Perrian, Theo Gi-
gant, Patrick Haller, Jenny Chim, Jose Posada, John
Giorgi, Karthik Rangasai Sivaraman, Marc Pamies,
Marianna Nezhurina, Robert Martin, Michael Cul-
lan, Moritz Freidank, Nathan Dahlberg, Shubhan-
shu Mishra, Shamik Bose, Nicholas Broad, Yanis
Labrak, Shlok Deshmukh, Sid Kiblawi, Ayush Singh,
Minh Chien Vu, Trishala Neeraj, Jonas Golde, Albert
Villanova del Moral, and Benjamin Beilharz. 2022.
Bigbio: A framework for data-centric biomedical
natural language processing. In *Advances in Neural*
Information Processing Systems, volume 35, pages
25792–25806. Curran Associates, Inc.

693	Chang Gao, Wenxuan Zhang, Guizhen Chen, and Wai Lam. 2023. Jsontuning: Towards generalizable, robust, and controllable instruction tuning. <i>arXiv preprint arXiv:2310.02953</i> .		
694			
695			
696			
697	Yu Gu, Robert Tinn, Hao Cheng, Michael Lucas, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. 2021. Domain-specific language model pretraining for biomedical natural language processing. <i>ACM Transactions on Computing for Healthcare</i> , 3(1):1–23.		
698			
699			
700			
701			
702			
703	Tianyu Han, Lisa C. Adams, Jens-Michalis Papaioannou, Paul Grundmann, Tom Oberhauser, Alexander Löser, Daniel Truhn, and Keno Kyrill Bressemer. 2023. Medalpaca - an open-source collection of medical conversational ai models and training data. <i>ArXiv</i> , abs/2304.08247.		
704			
705			
706			
707			
708			
709	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021a. Measuring massive multitask language understanding. In <i>International Conference on Learning Representations</i> .		
710			
711			
712			
713			
714	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021b. Measuring mathematical problem solving with the math dataset. <i>NeurIPS</i> .		
715			
716			
717			
718	Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023a. Camels in a changing climate: Enhancing lm adaptation with tulu 2. <i>arXiv preprint arXiv: 2311.10702</i> .		
719			
720			
721			
722			
723			
724	Hamish Ivison, Yizhong Wang, Valentina Pyatkin, Nathan Lambert, Matthew Peters, Pradeep Dasigi, Joel Jang, David Wadden, Noah A. Smith, Iz Beltagy, and Hannaneh Hajishirzi. 2023b. Camels in a changing climate: Enhancing lm adaptation with tulu 2. <i>arXiv preprint arXiv: 2311.10702</i> .		
725			
726			
727			
728			
729			
730	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</i> , pages 10030–10051, Singapore. Association for Computational Linguistics.		
731			
732			
733			
734			
735			
736			
737	Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoon Yun, Seongjin Shin, Sungdong Kim, James Thorne, and Minjoon Seo. 2024. Prometheus: Inducing fine-grained evaluation capability in language models. In <i>The Twelfth International Conference on Learning Representations</i> .		
738			
739			
740			
741			
742			
743			
744	Abdullatif Köksal, Timo Schick, Anna Korhonen, and Hinrich Schütze. 2023. Longform: Effective instruction tuning with reverse instructions. <i>arXiv preprint arXiv: 2304.08460</i> .		
745			
746			
747			
	Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. 2024a. Tulu 3: Pushing frontiers in open language model post-training. <i>Preprint</i> , arXiv:2411.15124.		748 749 750 751 752 753 754 755 756 757
	Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. 2024b. RewardBench: Evaluating reward models for language modeling. <i>arXiv preprint arXiv: 2403.13787</i> .		758 759 760 761 762 763
	Sihang Li, Jian Huang, Jiayi Zhuang, Yaorui Shi, Xiaochen Cai, Mingjun Xu, Xiang Wang, Linfeng Zhang, Guolin Ke, and Hengxing Cai. 2024. Scilittlem: How to adapt llms for scientific literature understanding. <i>arXiv preprint arXiv: 2408.15545</i> .		764 765 766 767 768
	Xian Li, Ping Yu, Chunting Zhou, Timo Schick, Omer Levy, Luke Zettlemoyer, Jason Weston, and Mike Lewis. 2023. Self-alignment with instruction back-translation. <i>arXiv preprint arXiv: 2308.06259</i> .		769 770 771 772
	Kyle Lo, Joseph Chee Chang, Andrew Head, Jonathan Bragg, Amy X. Zhang, Cassidy Trier, Chloe Anastasiades, Tal August, Russell Authur, Danielle Bragg, Erin Bransom, Isabel Cachola, Stefan Candra, Yoganand Chandrasekhar, Yen-Sung Chen, Evie (Yu-Yen) Cheng, Yvonne Chou, Doug Downey, Rob Evans, Raymond Fok, F.Q. Hu, Regan Huff, Dongyeop Kang, Tae Soo Kim, Rodney Michael Kinney, Aniket Kittur, Hyeonsu B Kang, Egor Kl-evak, Bailey Kuehl, Michael Langan, Matt Latzke, Jaron Lochner, Kelsey MacMillan, Eric Stuart Marsh, Tyler C. Murray, Aakanksha Naik, Ngoc-Uyen Nguyen, Srishti Palani, Soya Park, Caroline Paulic, Napol Rachatasumrit, Smita R Rao, Paul Sayre, Ze-jiang Shen, Pao Siangliulue, Luca Soldaini, Huy Tran, Madeleine van Zuylen, Lucy Lu Wang, Christopher Wilhelm, Caroline M Wu, Jiangjiang Yang, Angele Zamarron, Marti A. Hearst, and Daniel S. Weld. 2023. The semantic reader project: Augmenting scholarly documents through ai-powered interactive reading interfaces. <i>ArXiv</i> , abs/2303.14334.		773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793
	Shayne Longpre, Le Hou, Tu Vu, Albert Webson, Hyung Won Chung, Yi Tay, Denny Zhou, Quoc V Le, Barret Zoph, Jason Wei, and Adam Roberts. 2023. The flan collection: Designing data and methods for effective instruction tuning. In <i>Proceedings of the 40th International Conference on Machine Learning</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 22631–22648. PMLR.		794 795 796 797 798 799 800 801
	Renze Lou, Kai Zhang, Jian Xie, Yuxuan Sun, Jan-ice Ahn, Hanzi Xu, Yu su, and Wenpeng Yin. 2024. MUFFIN: Curating multi-faceted instructions for improving instruction following. In <i>The Twelfth International Conference on Learning Representations</i> .		802 803 804 805 806

807	Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In <i>The 36th Conference on Neural Information Processing Systems (NeurIPS)</i> .	861
808		862
809		863
810		864
811		865
812		866
813	Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct. <i>arXiv preprint arXiv: 2308.09583</i> .	867
814		868
815		869
816		
817		
818		
819	Iain James Marshall and Byron C. Wallace. 2019. Toward systematic review automation: a practical guide to using machine learning tools in research synthesis. <i>Systematic Reviews</i> , 8.	
820		
821		
822		
823	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2021. Cross-Task Generalization via Natural Language Crowdsourcing Instructions . <i>arXiv</i> .	
824		
825		
826		
827	Swaroop Mishra, Daniel Khashabi, Chitta Baral, and Hannaneh Hajishirzi. 2022. Cross-task generalization via natural language crowdsourcing instructions. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3470–3487, Dublin, Ireland. Association for Computational Linguistics.	
828		
829		
830		
831		
832		
833		
834	Nihal V. Nayak, Yiyang Nan, Avi Trost, and Stephen H. Bach. 2024. Learning to generate instruction tuning datasets for zero-shot task adaptation. <i>arXiv preprint arXiv: 2402.18334</i> .	
835		
836		
837		
838	Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback . <i>Preprint</i> , arXiv:2203.02155.	
839		
840		
841		
842		
843		
844		
845		
846	Srishti Palani, Aakanksha Naik, Doug Downey, Amy X. Zhang, Jonathan Bragg, and Joseph Chee Chang. 2023. Relatedly: Scaffolding literature reviews with existing related work sections. <i>Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems</i> .	
847		
848		
849		
850		
851		
852	Pallets. 2024. Jinja. https://github.com/pallets/jinja/ . GitHub repository.	
853		
854	Mihir Parmar, Swaroop Mishra, Mirali Purohit, Man Luo, Murad Mohammad, and Chitta Baral. 2022. In-BoXBART: Get instructions into biomedical multi-task learning. In <i>Findings of the Association for Computational Linguistics: NAACL 2022</i> , pages 112–128, Seattle, United States. Association for Computational Linguistics.	
855		
856		
857		
858		
859		
860		
	Yifan Peng, Shankai Yan, and Zhiyong Lu. 2019. Transfer learning in biomedical natural language processing: An evaluation of bert and elmo on ten benchmarking datasets. In <i>Proceedings of the 2019 Workshop on Biomedical Natural Language Processing (BioNLP 2019)</i> .	861
		862
		863
		864
		865
		866
	Machel Reid, Victor Zhong, Suchin Gururangan, and Luke Zettlemoyer. 2022. M2d2: A massively multi-domain language modeling dataset.	867
		868
		869
	David Rein, Betty Li Hou, Asa Cooper Stickland, Jackson Petty, Richard Yuanzhe Pang, Julien Dirani, Julian Michael, and Samuel R. Bowman. 2023. Gpqa: A graduate-level google-proof q&a benchmark. <i>arXiv preprint arXiv: 2311.12022</i> .	870
		871
		872
		873
		874
	Omid Rohanian, Mohammadmahdi Nouriborji, and David A. Clifton. 2023. Exploring the effectiveness of instruction tuning in biomedical language processing. <i>arXiv preprint arXiv: 2401.00579</i> .	875
		876
		877
		878
	Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Evfry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. 2022. Multi-task Prompted Training Enables Zero-Shot Task Generalization. In <i>International Conference on Learning Representations (ICLR)</i> .	879
		880
		881
		882
		883
		884
		885
		886
		887
		888
		889
		890
		891
		892
		893
		894
	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Mingchuan Zhang, Y.K. Li, Y. Wu, and Daya Guo. 2024. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. <i>CoRR</i> , abs/2402.03300.	895
		896
		897
		898
		899
	Chufan Shi, Yixuan Su, Cheng Yang, Yujiu Yang, and Deng Cai. 2023. Specialist or generalist? instruction tuning for specific NLP tasks . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 15336–15348, Singapore. Association for Computational Linguistics.	900
		901
		902
		903
		904
		905
	Amanpreet Singh, Mike D’Arcy, Arman Cohan, Doug Downey, and Sergey Feldman. 2023. SciRepEval: A multi-format benchmark for scientific document representations. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 5548–5566, Singapore. Association for Computational Linguistics.	906
		907
		908
		909
		910
		911
		912
	Yu Song, Santiago Miret, Huan Zhang, and Bang Liu. 2023. HoneyBee: Progressive instruction finetuning of large language models for materials science. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 5724–5739, Singapore. Association for Computational Linguistics.	913
		914
		915
		916
		917
		918

919	Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhe-Wei Shen, Baocai Chen, Lu Chen, and Kai Yu. 2023. Sci-eval: A multi-level large language model evaluation benchmark for scientific research. <i>AAAI Conference on Artificial Intelligence</i> .	models with self-generated instructions. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13484–13508, Toronto, Canada. Association for Computational Linguistics.	975 976 977 978 979
924	Zhengyang Tang, Xingxing Zhang, Benyou Wang, and Furu Wei. 2024. Mathscale: Scaling instruction tuning for mathematical reasoning. <i>arXiv preprint arXiv: 2403.02884</i> .	Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, Eshaan Pathak, Giannis Karamanolakis, Haizhi Lai, Ishan Purohit, Ishani Mondal, Jacob Anderson, Kirby Kuznia, Krima Doshi, Kuntal Kumar Pal, Maitreya Patel, Mehrad Moradshahi, Mihir Parmar, Mirali Purohit, Neeraj Varshney, Phani Rohitha Kaza, Pulkit Verma, Ravsehaj Singh Puri, Rushang Karia, Savan Doshi, Shailaja Keyur Sampat, Siddhartha Mishra, Sujan Reddy A, Sumanta Patro, Tanay Dixit, and Xudong Shen. 2022. Super-NaturalInstructions: Generalization via declarative instructions on 1600+ NLP tasks. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 5085–5109, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997
928	Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. 2022. Galactica: A large language model for science. <i>arXiv preprint arXiv: 2211.09085</i> .	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M. Dai, and Quoc V. Le. 2022. Finetuned language models are zero-shot learners. <i>International Conference on Learning Representations</i> .	998 999 1000 1001 1002
933	Augustin Toma, Patrick R. Lawler, Jimmy Ba, Rahul G. Krishnan, Barry Rubin, and Bo Wang. 2023. Clinical camel: An open-source expert-level medical language model with dialogue-based knowledge encoding. <i>ArXiv</i> , abs/2305.12031.	Shufa Wei, Xiaolong Xu, Xianbiao Qi, Xi Yin, Jun Xia, Jingyi Ren, Peijun Tang, Yuxiang Zhong, Yihao Chen, Xiaoqin Ren, Yuxin Liang, Liankai Huang, Kai Xie, Weikang Gui, Wei Tan, Shuanglong Sun, Yongquan Hu, Qinxian Liu, Nanjin Li, Chihao Dai, Lihua Wang, Xiaohui Liu, Lei Zhang, and Yutao Xie. 2023. Academicgpt: Empowering academic research . <i>Preprint</i> , arXiv:2311.12315.	1003 1004 1005 1006 1007 1008 1009 1010
938	Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Gitman. 2024. Openmathinstruct-1: A 1.8 million math instruction tuning dataset. <i>arXiv preprint arXiv: Arxiv-2402.10176</i> .	Chaoyi Wu, Weixiong Lin, Xiaoman Zhang, Ya Zhang, Weidi Xie, and Yanfeng Wang. 2024. PMC-LLaMA: toward building open-source language models for medicine . <i>Journal of the American Medical Informatics Association</i> , page ocae045.	1011 1012 1013 1014 1015
943	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	Tong Xie, Yuwei Wan, Wei Huang, Zhenyu Yin, Yixuan Liu, Shaozhou Wang, Qingyuan Linghu, Chunyu Kit, Clara Grazian, Wenjie Zhang, Imran Razzak, and Bram Hoex. 2023. Darwin series: Domain specific large language models for natural science. <i>arXiv preprint arXiv: 2308.13565</i> .	1016 1017 1018 1019 1020 1021
949	Hieu Tran, Zhichao Yang, Zonghai Yao, and Hong Yu. 2023. Bioinstruct: Instruction tuning of large language models for biomedical natural language processing. <i>JAMIA Journal of the American Medical Informatics Association</i> .	Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. 2024. WizardLM: Empowering large pre-trained language models to follow complex instructions. In <i>The Twelfth International Conference on Learning Representations</i> .	1022 1023 1024 1025 1026 1027
954	Somin Wadhwa, Silvio Amir, and Byron C. Wallace. 2023. Revisiting relation extraction in the era of large language models. <i>Proceedings of the conference. Association for Computational Linguistics. Meeting, 2023:15566–15589</i> .	An Yang, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanting Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang,	1028 1029 1030 1031 1032
959	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>arXiv preprint arXiv: 2304.08085</i> .		
965	Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. 2023b. SCIBENCH: Evaluating college-level scientific problem-solving abilities of large language models. In <i>The 3rd Workshop on Mathematical Reasoning and AI at NeurIPS'23</i> .		
972	Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A. Smith, Daniel Khashabi, and Hannaneh Hajishirzi. 2023c. Self-instruct: Aligning language		

1033	et al. 2024. Qwen2 technical report . <i>Preprint</i> , arXiv:2407.10671.	A Additional information on SciRIFF	1086
1034			
1035	Da Yin, Xiao Liu, Fan Yin, Ming Zhong, Hritik Bansal, Jiawei Han, and Kai-Wei Chang. 2023. Dynosaur: A dynamic growth paradigm for instruction-tuning data curation. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 4031–4047, Singapore. Association for Computational Linguistics.	A.1 SciRIFF task list	1087
1036		The full list of SciRIFF tasks is visualized in Figure 4. Detailed information on all tasks—including citations, URLs to source websites, and licensing information where available—will be provided in our dataset card. Where convenient, we use datasets as preprocessed by the BigBio resource (https://huggingface.co/bigbio); details are in the dataset card.	1088
1037			1089
1038			1090
1039			1091
1040			1092
1041			1093
1042	Botao Yu, Frazier N. Baker, Ziqi Chen, Xia Ning, and Huan Sun. 2024. Llamol: Advancing large language models for chemistry with a large-scale, comprehensive, high-quality instruction tuning dataset. <i>arXiv preprint arXiv: 2402.09391</i> .	A.2 Task length distribution	1096
1043		Figure 5 shows the distribution of input and output lengths for demonstrations in SciRIFF.	1097
1044		Table 5 compares SciRIFF with selected instruction-following datasets, including canonical collections commonly used for general fine-tuning and selected datasets specialized in scientific domains. Our dataset features longer input contexts than existing resources.	1098
1045			1099
1046			1100
1047	Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhui Chen. 2024a. MAMmoTH: Building math generalist models through hybrid instruction tuning. In <i>The Twelfth International Conference on Learning Representations</i> .	A.3 Truncation strategy	1105
1048		In §2.2, we mention that when an instance exceeds the maximum context length for a given version of SciRIFF, we truncate where possible and discard otherwise. In particular, we truncate for tasks (like question answering) where the task output can be localized to particular passages in the input document by randomly removing irrelevant passages until the document fits in the desired context. For tasks like summarization, where the task output cannot easily be localized, we simply discard examples that are longer than the context window.	1106
1049			1107
1050			1108
1051			1109
1052			1110
1053	Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhui Chen. 2024b. Mammoth2: Scaling instructions from the web. <i>arXiv preprint arXiv: 2405.03548</i> .	B Instruction template creation	1117
1054		Instruction templates are written in (Pallets, 2024), Guidelines and “best practices” for prompt-writing will be available at our GitHub repository. Each prompt was double-checked by an additional paper author for clarity and correctness.	1118
1055			1119
1056	Dan Zhang, Ziniu Hu, Sining Zhou, Zhengxiao Du, Kaiyu Yang, Zihan Wang, Yisong Yue, Yuxiao Dong, and Jie Tang. 2024a. Sciglm: Training scientific language models with self-reflective instruction annotation and tuning. <i>arXiv preprint arXiv: 2401.07950</i> .		1120
1057			1121
1058			1122
1059			1123
1060			1124
1061	Di Zhang, Wei Liu, Qian Tan, Jingdan Chen, Hang Yan, Yuliang Yan, Jiatong Li, Weiran Huang, Xiangyu Yue, Wanli Ouyang, Dongzhan Zhou, Shufei Zhang, Mao Su, Han-Sen Zhong, and Yuqiang Li. 2024b. Chemllm: A chemical large language model. <i>arXiv preprint arXiv: 2402.06852</i> .	C Sample template	1123
1062		In this section, we provide examples of our expert-written templates that demonstrate the complexity and precision required for scientific literature understanding tasks, described in §1 and §2.1. These templates are carefully designed to elicit structured outputs while requiring sophisticated capabilities such as information extraction with attribution, multi-step reasoning, and adherence to specific output	1124
1063			1125
1064			1126
1065			1127
1066			1128
1067	Shengyu Zhang, Linfeng Dong, Xiaoya Li, Sen Zhang, Xiaofei Sun, Shuhe Wang, Jiwei Li, Runyi Hu, Tianwei Zhang, Fei Wu, and Guoyin Wang. 2023a. Instruction tuning for large language models: A survey. <i>arXiv preprint arXiv: 2308.10792</i> .		1129
1068			1130
1069			1131
1070			1132
1071			1133
1072	Tianshu Zhang, Xiang Yue, Yifei Li, and Huan Sun. 2023b. Tablellama: Towards open large generalist models for tables. <i>arXiv preprint arXiv: 2311.09206</i> .		1134
1073			1135
1074			1136
1075	Xinlu Zhang, Chenxin Tian, Xianjun Yang, Lichang Chen, Zekun Li, and Linda Ruth Petzold. 2023c. Alpacare:instruction-tuned large language models for medical application. <i>arXiv preprint arXiv: 2310.14558</i> .		1137
1076			1138
1077			1139
1078			1140
1079			1141
1080	Chunting Zhou, Pengfei Liu, Puxin Xu, Srini Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, L. Yu, Susan Zhang, Gargi Ghosh, M. Lewis, Luke Zettlemoyer, and Omer Levy. 2023. Lima: Less is more for alignment. In <i>Advances in Neural Information Processing Systems</i> (NeurIPS).		1142
1081			1143
1082			1144
1083			1145
1084			1146
1085			1147

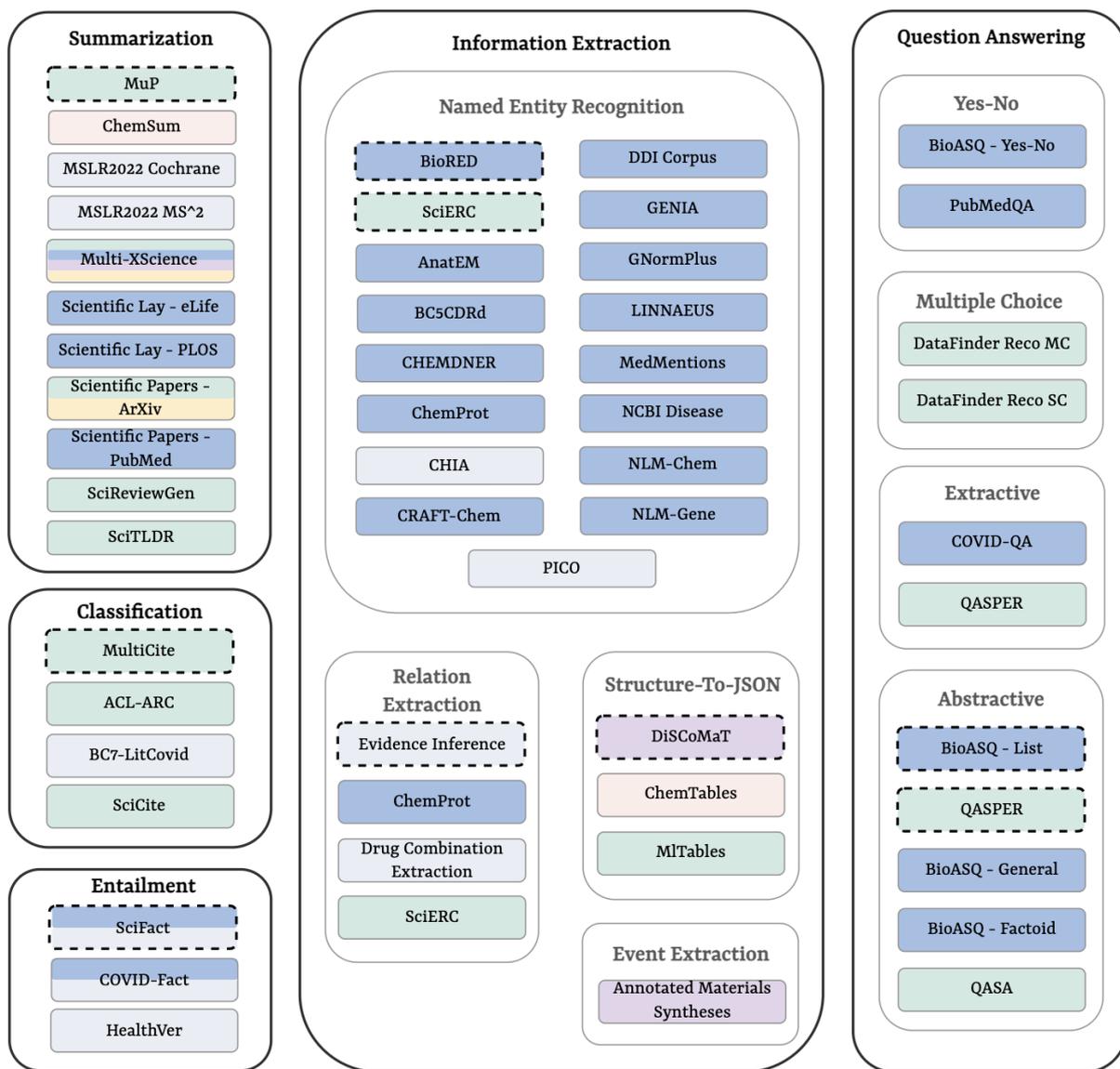


Figure 4: Overview of SciRIFF dataset. Dashed black lines indicate that a task is included in SciRIFF-EVAL and held out during model training. Scientific domains are colored as follows: ■ Biomedicine; ■ AI; ■ Clinical Medicine; ■ Chemistry; ■ Materials Science; ■ Miscellaneous.

Name	# Instances	Domain	Avg. Length
<i>General Domain</i>			
Flan V2 (Chung et al., 2024)	15M	General	355.6 / 31.2
SuperNI (Wang et al., 2022)	97K	General	291.1 / 38.7
TULU V2 MIX (Iverson et al., 2023b)	326K	General	353.3 / 696.9
<i>Scientific Domain</i>			
BoX (Parmar et al., 2022)	141K	Biomed	X*
SciInstruct (Zhang et al., 2024a)	254K	Math, PH, Chem, FP	88.4 / 265.6
Mol-Instructions (Fang et al., 2024)	2.04M	Biomolecular	126.3 / 112.9
MathInstruct (Yue et al., 2024a)	262K	Math	82.5 / 174.0
MedInstruct-52K (Zhang et al., 2023c)	52K	Medical	148.2 / 96.9
LlaSMol (Yu et al., 2024)	3.29M	Chem	81.9 / 53.0
SciRIFF (Our work)	137K	AI, Biomed, Clinical, Chem, MatSci	1242.9 / 139.6

Table 5: Comparison with selected instruction-following datasets. We use the following abbreviations: PH – Physics; FP – Formal Proof; MatSci – Materials Science. We report average token counts for input/output using Llama 2 tokenizer using up to 200k subsamples from each dataset. * BoX dataset is not readily available.

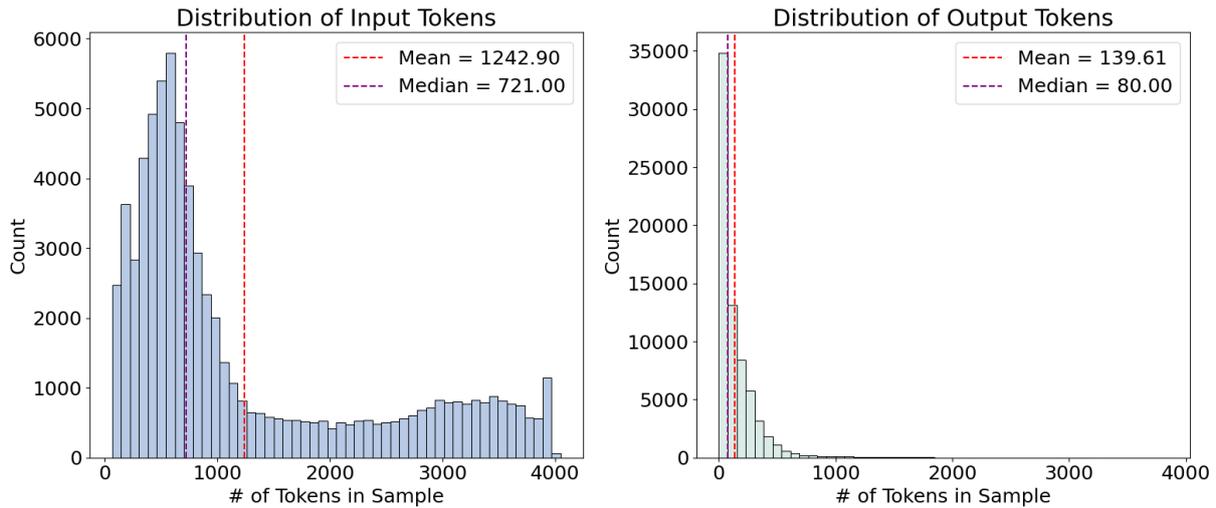


Figure 5: Distribution of input (left) and output (right) token lengths over SciRIFF training instances.

schemas. The templates shown –QASPER (QA, Figure 6), SciERC (IE, Figure 7), HealthVer (Fact-checking, Figure 8), DiSCoMaT (IE over tabular data, Figure 9), and DataFinder Reco MC (Multiple Choice QA, Figure 10) – demonstrates how our instruction format guides models to perform challenging tasks like answering questions with evidence attribution, extracting nested entity relationships, and verifying scientific claims with supporting rationales.⁹

D Training Details

For instruction-tuning, our training hyperparameters were as follows:

- Precision: BFloat16
- Epochs: 5
- Weight decay: 0
- Warmup ratio: 0.03
- Learning rate: $2e-5$ ($1e-5$ for 70B)
- Max. seq. length: 4,096
- Effective batch size: 128

All of our models were trained on v3-128 TPUs on the Google TPU Research Cloud.

E Evaluation details

The following pages show full input / output examples for all SciRIFF-EVAL tasks, along with details on metric calculations. This information will be available on our project GitHub page. We use

⁹Our preliminary experiments showed that even strong proprietary models like GPT-4o struggled to reliably generate such structured outputs without explicit templates. This observation motivated our decision to use expert-written templates.

gpt-4o-2024-08-06 model for tasks using an LLM judge as evaluation.

1159
1160

QASPER

You will be shown sections from a scientific research paper, together with a question about the paper. This is an extractive question-answering task, where you must find and extract relevant text spans directly from the paper to answer the question. Your response should strictly be a 'json' object with two fields:

- "answer": An array of strings extracted directly from the paper which, collectively, answer the question.
- "evidence": An array of strings. Each should be an excerpt from the paper, in which one or more of the extracted answers can be found.

For example, for the question "What baselines did the authors compare against?", a sample response might be:

```
{
  "answer": ["BERT", "RoBERTa"],
  "evidence": ["In our experiments, we compare the performance of our model against BERT and RoBERTa."]
}
```

Do not include any text in your response other than the json.

If the question is unanswerable given the provided excerpts, respond with the single word "null".

Paper: {{paper}}

Question: {{question}}

|||

```
{% if unanswerable %} null
{% else %}
{{ {"answer": answer, "evidence": evidence} | tojson }}
{% endif %}
```

Figure 6: Canonical template for QASPER task in Figure 4. See §C for description.

SciERC

You will be shown an abstract from a computer science research paper. Given this abstract, your task is to extract all unique entities with the following types:

- "Task": Applications, problems to solve, systems to construct. Examples include "information extraction", "machine reading system", "image segmentation".
- "Method": : Methods, models, systems to use, or tools, components of a system, frameworks. Examples include "language model", "CORENLP", "POS parser".
- "Metric": Metrics, measures, or entities that can express quality of a system / method. Examples include "F1", "BLEU", "Precision", "time complexity".
- "Material": Data, datasets, resources, Corpus, Knowledge base. Examples include "image data", "speech data", "stereo images", "CoNLL", "Wikipedia".
- "OtherScientificTerm": Phrases that are a scientific terms but do not fall into any of the above classes. Examples include "physical or geometric constraints", "qualitative prior knowledge", "tree kernel", "noise".
- "Generic": General terms or pronouns that may refer to an entity but are not themselves informative, often used as connection words. Examples include "model", "approach", "them".

Please return the output as a JSON object of the format: {"type1": ["example_entity", ...], "type2": ["example_entity", ...]}. The keys should be entity types and values should be lists of extracted entities belonging to the corresponding type. Entity types with no matching entities should be assigned an empty array [].

For instance, the output might look like: {"Task": ["speech recognition", ...], "Method": ["Conditional random field"], "Material": [], ...}.

Only output the JSON object and do not include any additional text.

Abstract:

```
{{ org_text }}
```

|||

```
{{ ner_dict | tojson }}
```

Figure 7: Canonical template for SciERC task in Figure 4. See §C for description.

HealthVer

You will be shown a claim about public health and the abstract of a biomedical research paper. Each sentence from the abstract will be on a separate line. Your task is to return a JSON object with two fields:

- "verdict": The fact-checking verdict. If the information in the abstract supports the claim, write "SUPPORT". If the abstract contradicts the claim, write "CONTRADICT". If the abstract does not provide enough information to arrive at a verdict, write "NEI" (for "not enough information").
- "evidence": An array of sentences providing evidence for the verdict. Please copy all relevant sentences verbatim from the abstract. If the verdict was "NEI", then return an empty array.

For instance, if the model were given the claim "wearing masks can prevent the spread of COVID", the output might be:

```
{  
  "verdict": "SUPPORT",  
  "evidence": ["Our findings indicate that mass mask-wearing reduces the transmission rate for COVID-19."]  
}
```

Claim: {{ claim }}

Abstract:

{{ abstract_with_newlines }}

|||

{{ output_json_with_sentences }}

Figure 8: Canonical template for HealthVer task in Figure 4. See §C for description.

DiSCoMaT

{{ table_code_text }}

You are provided with the table above from a materials science paper. Here are JSON templates for two types of numeric cells: "Other" and "Glass_Compound_Amount":

```
{"value": "xx", "type": "Other"}  
{"value": "xx", "type": "Glass_Compound_Amount", "constituent": "xx", "unit": "xx", "material": "xx"}
```

Please describe all numeric cells in the above table following the JSON templates (proceeding by row in a left-right, top-down direction). For each cell, output one JSON description per line. For any unanswerable attributes in the templates, set their value to the placeholder "xx".

Cell Description:

|||

{{ json_records }}

Figure 9: Canonical template for DiSCoMaT task in Figure 4. See §C for description.

DataFinder Reco MC

You are provided with a research question, keyphrases about the question, a description of candidate datasets and dataset options. Read the description of popular datasets provided below and select the ones that can be used to validate the following research question. Use your knowledge of machine learning datasets to make the best judgement.

Your response should be formatted as a 'json' array. For instance, for the query "Semi supervised image classification", a sample response might be: ["CIFAR-10", "CIFAR-100"]. Do not include any extra text in the response other than the answer array.

Query: {{ query }}

Keyphrases: {{ keyphrase_query }}

Dataset description:
{{ context }}

Options:- {{ options }}

|||

```
{%- set ans_list = answer.split(", ") %}  
{{ ans_list | tojson }}
```

Figure 10: Canonical template for DataFinder Reco MC (QA-multiple choice) task in Figure 4. See §C for description.

Evaluation tasks

This doc has a list of all evaluation tasks, including input / output examples and evaluation metrics.

Table of contents

- [BioASQ](#): question answering
- [BioRED](#): named entity recognition
- [Discomat](#): table extraction
- [Evidence inference](#): evidence tuple extraction
- [Multicite](#): citation intent classification
- [MUP](#): summarization
- [Qasper](#): paper question answering
- [SciERC](#): named entity recognition
- [SciFact](#): claim verification

BioASQ

- Task input: A collection of biomedical research excerpts and a question answerable from the excerpts.
- Task output: A list of answers to the question.
- Metrics: Compare predicted vs. reference answers using exact-match F1.

Input

Below are a collection of excerpts from biomedical research articles. Excerpts are separated by newlines. Your task is to answer a question based these excerpts. Your response should be formatted as a `json` array.

For instance, given excerpts from articles studying breast cancer, and the question "what are some common genes associated with breast cancer?", an answer might be formatted like: ["BRCA1", "BRCA2", "TP53", ...]. Only include answers that are mentioned in the provided excerpts. The array should always have at least one answer; do not respond with an empty array []. Do not include any text in your response other than the answer array.

Context: sensitization, behavioral changes, and low body mass index (BMI). One possible cellular target that may mediate some of these findings is the hypocretin/orexin neurons. This neuronal system plays a role in regulating wakefulness/sleep cycles, pain perception, and appetite. Food intake, in contrast, receives circadian modulation through hormones such as leptin, ghrelin, insulin and orexin. A low level of hypocretin-1/orexin-A in the cerebrospinal fluid is sufficient to diagnose narcolepsy type 1, being a highly specific and sensitive biomarker, and the irreversible loss of hypocretin neurons is responsible for the main symptoms of the disease: Orexins, or hypocretins, are excitatory neuropeptides involved in the regulation of feeding behavior and the sleep and wakefulness states.

[Lines omitted for space]

Orexin A (OXA) and orexin B (OXB) are recently discovered neuropeptides that appear to play a role in various distinct functions such as arousal and the sleep-wake cycle as well as on appetite and regulation of feeding and energy homeostasis. Orexins were first described as neuropeptides expressed by a sp Orexin/hypocretin neurons located in the lateral hypothalamus play a critical role in the maintenance of arousal and contribute to the regulation of multiple homeostatic and behavioral processes.

Question: What processes do orexin/hypocretin neurons regulate?

Output

```
[
  "sleep",
  "appetite",
  "wakefulness",
  "pain",
  "reward",
  "energy homeostasis",
  "goal-directed behaviors",
  "Arousal",
  "addiction"
]
```

BioRed

- Task input: Abstract of a biomedical research article.
- Task output: All entities in the article of the following types:
 - cell line
 - chemical
 - disease
 - gene
 - gene variant
 - species
- Metrics: Compare predicted vs. reference entities using exact-match F1.

Input

You will be shown an abstract from a biomedical research paper. Given this abstract, your task is to extract all unique entities of the following types: ["Chemical", "Variant", "Gene", "CellLine", "Disease", "Species"].

Please return the output as a JSON object of the format: {"CellLine": ["hRPTEC", ...], "Chemical": ["Glucose", ...], "Disease": ["Diabetes", ...], "Gene": ["HNF-6", ...], "Species": ["Patients", ...], "Variant": ["Pro75Ala", ...]}. The keys should be entity types and values should be lists of extracted entities belonging to the corresponding type. If you cannot find entities belonging to a specific type, the value should be [].

Only output the JSON object and do not include any additional text.

Abstract:

Fatal carbamazepine induced fulminant eosinophilic (hypersensitivity) myocarditis: emphasis on anatomical and histological characteristics, mechanisms and genetics of drug hypersensitivity and differential diagnosis. The most severe adverse reactions to carbamazepine have been observed in the haemopoietic system, the liver and the cardiovascular system. A frequently fatal, although exceptionally rare side effect of carbamazepine is necrotizing eosinophilic (hypersensitivity) myocarditis. We report a case of hypersensitivity myocarditis secondary to administration of carbamazepine. Acute hypersensitivity myocarditis was not suspected clinically, and the diagnosis was made post-mortem. Histology revealed diffuse infiltration of the myocardium by eosinophils and lymphocytes with myocyte damage. Clinically, death was due to cardiogenic shock. To best of our knowledge this is the second case of fatal carbamazepine induced myocarditis reported in English literature.

Output

```
{
  "CellLine": [],
  "Chemical": ["carbamazepine"],
  "Disease": [
    "hypersensitivity",
    "death",
    "myocarditis",
    "cardiogenic shock",
    "drug hypersensitivity"
  ],
  "Gene": [],
  "Species": [],
  "Variant": []
}
```

Discomat

- Task input: A passage from a research paper including a table.
- Task output: The table, with each cell as a `json` line.
- Metrics: BLEU score between predicted and gold reference. Manual inspection showed that BLEU was pretty reliable for this task.

Input

Sample no.	Ph, volume percent of crystals	Activation energy (kJ/mol) of the scale factor for normalised frequency	$\ln(t_0, s)$ of the scale factor	G unrelaxed shear modulus (GPa)
Glas 0	0			137+-18
-50.15		24.3		
Glas 1	17			129+-13

	-47.68		23.8	
Glas 3	22		126+-16	
	-45.72		24.7	
Glas 5	27		117+-10	
	-42.25		25.0	

Caption: Activation energies of shear stress relaxation and unrelaxed shear modulus of disilicate lithium glasses

You are provided with a table from a material science paper. Here are JSON templates for two types of numeric cells: "Other" and "Glass_Compound_Amount": {"value": "xx", "type": "Other"} {"value": "xx", "type": "Glass_Compound_Amount", "constituent": "xx", "unit": "xx", "material": "xx"}

Please describe all numeric cells in the above table following the JSON templates (proceeding by row in a left-right, top-down direction). For each cell, output one JSON description per line. For any unanswerable attributes in the templates, set their value to the placeholder "xx".

Cell Description:

Output

```
{ "value": "0", "type": "Other" }
{ "value": "137", "type": "Other" }
{ "value": "24.3", "type": "Other" }
{ "value": "17", "type": "Other" }
{ "value": "129", "type": "Other" }
{ "value": "23.8", "type": "Other" }
{ "value": "22", "type": "Other" }
{ "value": "126", "type": "Other" }
{ "value": "24.7", "type": "Other" }
{ "value": "27", "type": "Other" }
{ "value": "117", "type": "Other" }
{ "value": "25.0", "type": "Other" }
```

Evidence Inference

- Task input: Abstract of a clinical trial report.
- Task output: List of all ICO (**intervention / comparator / outcome**) tuples, together with the effect of the intervention on the outcome and the textual evidence of this effect.
- Metrics: "Fuzzy" F1. Given a prediction and a reference tuple, compute the token overlap for each tuple item. If token overlaps for all fields exceed 0.3, the predicted tuple is judged as a match to the reference.

Input

You will be shown the abstract of a medical clinical trial report. Your task is to extract all the findings from this report into a JSON array. Each finding should contain the following five elements:

- Intervention: The medical intervention being tested. This should be a text span copied from the

input passage.

- Comparator: The baseline against which the intervention is being evaluated. This should be a text span copied from the input passage. If no comparator is reported, set to `null`.
- Outcome: The medical outcome whose effect is being measured. This should be a text span copied from the input passage.
- Effect: The effect of the intervention on the outcome, relative to the comparator. The effect should be one of the following three values: ("significantly increased", "significantly decreased", "no significant difference").
- Evidence: The evidence for the effect. This should be a text span copied from the input passage.

Please format your output as a JSON array. Each entry in the output should be an array containing the 5 elements listed above, in the following order: [<intervention>, <comparator>, <outcome>, <effect>, <evidence>].

For example, an output with two findings might read: [["aspirin", "placebo", "headache severity", "significantly decreased", "Mean headache severity was significantly decreased in the aspirin group compared to the placebo group (p < 0.05)."], ["aspirin", "placebo", "weight loss", "no significant difference", "We did not observe any difference in weight loss between the group given aspirin relative to the control group"]]

There are 3 finding(s) in the abstract below. Please extract them. Output only the JSON array with these 3 findings. Do not include any additional text.

Abstract: ABSTRACT.OBJECTIVES: To compare the efficacy and safety of SB4 (an etanercept biosimilar) with reference product etanercept (ETN) in patients with moderate to severe rheumatoid arthritis (RA) despite methotrexate (MTX) therapy.

ABSTRACT.METHODS: This is a phase III, randomised, double-blind, parallel-group, multicentre study with a 24-week primary endpoint. Patients with moderate to severe RA despite MTX treatment were randomised to receive weekly dose of 50 mg of subcutaneous SB4 or ETN. The primary endpoint was the American College of Rheumatology 20% (ACR20) response at week 24. Other efficacy endpoints as well as safety, immunogenicity and pharmacokinetic parameters were also measured.

ABSTRACT.RESULTS: 596 patients were randomised to either SB4 (N=299) or ETN (N=297). The ACR20 response rate at week 24 in the per-protocol set was 78.1% for SB4 and 80.3% for ETN. The 95% CI of the adjusted treatment difference was -9.41% to 4.98%, which is completely contained within the predefined equivalence margin of -15% to 15%, indicating therapeutic equivalence between SB4 and ETN. Other efficacy endpoints and pharmacokinetic endpoints were comparable. The incidence of treatment-emergent adverse events was comparable (55.2% vs 58.2%), and the incidence of antidrug antibody development up to week 24 was lower in SB4 compared with ETN (0.7% vs 13.1%).

ABSTRACT.CONCLUSIONS: SB4 was shown to be equivalent with ETN in terms of efficacy at week 24. SB4 was well tolerated with a lower immunogenicity profile. The safety profile of SB4 was comparable with that of ETN.

ABSTRACT.TRIAL REGISTRATION NUMBERS: NCT01895309, EudraCT 2012-005026-30.

Findings:

Output

```
[
  [
    "SB4 (an etanercept biosimilar)",
    "etanercept (ETN)",
    "therapeutic equivalence – Patients with moderate to severe rheumatoid arthritis (RA) despite methotrexate (MTX) treatment",
    "no significant difference",
    "The 95% CI of the adjusted treatment difference was \u22129.41% to 4.98%, which is completely contained within the predefined equivalence margin of \u221215% to 15%, indicating therapeutic equivalence between SB4 and ETN."
  ],
  [
    "SB4 (an etanercept biosimilar)",
    "etanercept (ETN)",
    "incidence of antidrug antibody development up to week 24 – Patients with moderate to severe rheumatoid arthritis (RA) despite methotrexate (MTX) treatment",
    "significantly decreased",
    "the incidence of antidrug antibody development up to week 24 was lower in SB4 compared with ETN (0.7% vs 13.1%)."
  ],
  [
    "SB4 (an etanercept biosimilar)",
    "etanercept (ETN)",
    "incidence of treatment-emergent adverse events – Patients with moderate to severe rheumatoid arthritis (RA) despite methotrexate (MTX) treatment",
    "no significant difference",
    "The incidence of treatment-emergent adverse events was comparable (55.2% vs 58.2%)"
  ]
]
```

Multicite

- Task Input: A citation sentence from a research paper.
- Task output: A list of intents for the citation sentence.
- Metrics: Compare predicted vs. reference intents using exact-match F1.

Input

Your task is to classify the citation intent within the following provided text from a computational linguistics research paper. The cited work is demarcated by "<cite>" and "</cite>". Determine the purpose of the cited work by selecting from the listed categories:

- Background: The cited paper underpins the subject matter.
- Motivation: The cited paper inspires or provides a rationale for the current research.
- Uses: The current work utilizes concepts or tools from the cited paper.
- Extends: The current work advances ideas or methods from the cited paper.

- Similarities: The current work identifies commonalities¹⁶⁷ with the cited paper.
- Differences: The current work delineates its distinction from the cited paper.
- FutureWork: The cited paper is acknowledged as groundwork for prospective research.

Indicate the intents by listing them in a `json` array, e.g. ["Background", "Uses"]. More than one intent may be applicable. Do not include any extraneous text in your response.

Context with Citation: In addition to that, we implemented semi-supervised classification by training in the positive samples of the <cite>[9]</cite> dataset and training in only the lexicon as negative samples.

Output

```
["Similarities", "Uses"]
```

MUP

- Task input: Full text of a machine learning paper.
- Task output: Short paper summary that a reviewer might write as part of a paper review.
- Metrics: Use GPT-3.5 to judge similarity of generated summary to human reference on 1-5 scale. Based on manual inspection, this was higher-quality than automated metrics like ROUGE.

Input

You will be presented with the title and body text of a computer science research paper. Please write a summary of the work that would be informative for a peer reviewer assessing its quality. Your summary should be 3 sentences long. In your response, include only the summary and no additional text.

Paper title: Reinforcement Learning with Efficient Active Feature Acquisition

Paper body: 1 INTRODUCTION . Recently , machine learning models for automated sequential decision making have shown remarkable success across many application areas , such as visual recognition (Mathe et al. , 2016 ; Das et al. , 2017) , robotics control (Finn et al. , 2016 ; Zhang et al. , 2018) , medical diagnosis (Ling et al. , 2017 ; Peng et al. , 2018) and computer games (Mnih et al. , 2015 ; Silver et al. , 2016) . One fundamental reason that drives the success of such models and enables them to outperform classical algorithms is the availability of large amounts of training data . Typically such training data is either fully observed or the features stem from an action-independent observation model (which clearly can depend on the state of the system) . However , the fundamental assumption that the same features are always readily available during deployment could not hold in many real-world applications . For instance , consider a medical support system for monitoring and treating patients during their stay at hospital which was trained on rich historical medical data . To provide the best possible treatment , the system might need to

perform several measurements of the patient over time , while some of them could be costly or even pose a health risk . Therefore , during deployment , it is more ideal that the system could function with minimal features while during training more features might have been available . In such cases , we are interested in decision making models that actively take the measurement process , i.e. , feature acquisition , into account and only acquire the information relevant for making a decision .

In this paper , we consider the challenging problem of learning effective policies when the cost of information acquisition can not be neglected . To be successful , we need to learn policies which acquires the information required for solving a task in the cheapest way possible . [Truncated for space].

3-sentence paper summary:

Output

In this paper the authors propose an approach for simultaneously learning how to explore more efficiently in POMDPs via targeted feature acquisition, and learning a reward-maximizing control policy, balancing the cost of feature acquisition with the expected reward. Learning is done via a VAE framework which combines a belief inference model and an observation decoder, with a key innovation being that inference is done as a sequential process. Results comparing this approach to other variational inference approaches show the proposed framework reaches better performance with lower cost (particularly, number of acquired features).

Qasper

- Task input: The full text of an NLP research paper, and a question answerable from the paper body (but not the abstract).
- Task output: An answer to the question, accompanied by the extracts from the paper body supplying the answer.
- Metrics: We compute metrics for both the answer and the evidence.
 - Answer: GPT-3.5 judge of similarity of model answer to human reference (1-5 scale).
 - Evidence: Token F1 overlap with gold evidence.

Input

You will be shown sections from a scientific research paper, together with a question about the paper. Paragraphs in the paper are separated by newlines. Your task is to answer the question based on the contents of the paper.

Paper:

Named Entity Disambiguation for Noisy Text

We address the task of Named Entity Disambiguation (NED) for noisy text. We present WikilinksNED, a

large-scale NED dataset of text fragments from the web, which is significantly noisier and more challenging than existing news-based datasets. To capture the limited and noisy local context surrounding each mention, we design a neural model and train it with a novel method for sampling informative negative examples. We also describe a new way of initializing word and entity embeddings

that significantly improves performance. Our model significantly outperforms existing state-of-the-art methods on WikilinksNED while achieving comparable performance on a smaller newswire dataset.

The WikilinksNED Dataset: Entity Mentions in the Web We introduce WikilinksNED, a large-scale NED dataset based on text fragments from the web. Our dataset is derived from the Wikilinks corpus BIBREF14, which was constructed by crawling the web and collecting hyperlinks (mentions) linking to Wikipedia concepts (entities) and their surrounding text (context).

Wikilinks

contains 40 million mentions covering 3 million entities, collected from over 10 million web pages.

Wikilinks can be seen as a large-scale, naturally-occurring, crowd-sourced dataset where thousands of human annotators provide ground truths for mentions of interest. This means that the dataset contains various kinds of noise, especially due to incoherent contexts. The contextual noise presents an interesting test-case that supplements existing datasets that are sourced from mostly coherent and well-formed text.

[Truncated for space]

Question: How was a quality control performed so that the text is noisy but the annotations are accurate?

To answer the question, format your response as a `json` object with two fields:

"answer": A string providing a succinct answer to the question, in your own words. "evidence": An array of strings. Each entry should be a full paragraph from the paper. Together, the evidence should serve as a justification for the answer.

For instance, for the question "What baselines did the authors compare against?", a sample response might be:

```
{ "answer": "BERT and RoBERTa." "evidence": ["We compare our approach against two baselines. In Table 1, we compare against BERT. In Table 2, we compare against RoBERTa. Our findings indicate that our approach improves over both baselines..."] }
```

The "answer" field should be roughly 190 characters in length.

Do not include any text in your response other than the json. If the question is unanswerable given the provided excerpts, respond with the single word "null".

To repeat, the question is: How was a quality control performed so that the text is noisy but the annotations are accurate?

Answer JSON object:

Output

```
{
  "answer": "Profile pictures from the Twitter users' profiles.",
  "evidence": [
    "The recent advancements in deep neural networks, specifically for image analysis task, can lead to determining demographic features such as age and gender BIBREF13 . We show that by determining and integrating heterogeneous set of features from different modalities \u2013 aesthetic features from posted images (colorfulness, hue variance, sharpness, brightness, blurriness, naturalness), choice of profile picture (for gender, age, and facial expression), the screen name, the language features from both textual content and profile's description (n-gram, emotion, sentiment), and finally sociability from ego-network, and user engagement \u2013 we can reliably detect likely depressed individuals in a data set of 8,770 human-annotated Twitter users."
  ]
}
```

SciERC

- Task input: An abstract of an NLP paper.
- Task output: A list of all entities mentioned in the paper of the following types:
 - Material
 - Method
 - Metric
 - Task
 - Generic
 - Other scientific term
- Metrics: Exact-match F1.

Input

You will be shown an abstract from a computer science research paper. Given this abstract, your task

is to extract all unique entities with the following types:

- "Task": Applications, problems to solve, systems to construct. Examples include "information extraction", "machine reading system", "image segmentation".
- "Method": : Methods, models, systems to use, or tools, components of a system, frameworks. Examples include "language model", "CORENLP", "POS parser".
- "Metric": Metrics, measures, or entities that can express quality of a system / method. Examples include "F1", "BLEU", "Precision", "time complexity".
- "Material": Data, datasets, resources, Corpus, Knowledge base. Examples include "image data", "speech data", "stereo images", "CoNLL", "Wikipedia".
- "OtherScientificTerm": Phrases that are a scientific terms but do not fall into any of the above classes. Examples include "physical or geometric constraints", "qualitative prior knowledge", "tree kernel", "noise".
- "Generic": General terms or pronouns that may refer to a entity but are not themselves informative, often used as connection words. Examples include "model", "approach", "them".

Please return the output as a JSON object of the format: `{ "type1" : ["example_entity", ...], "type2" : ["example_entity", ...] }`. The keys should be entity types and values should be lists of extracted entities belonging to the corresponding type. Entity types with no matching entities should be assigned an empty array `[]`.

For instance, the output might look like: `{ "Task": ["speech recognition", ...], "Method": ["Conditional random field"], "Material": [], ...}`.

Only output the JSON object and do not include any additional text.

Abstract:

We present a syntax-based constraint for word alignment, known as the cohesion constraint. It requires disjoint English phrases to be mapped to non-overlapping intervals in the French sentence.

We evaluate the utility of this constraint in two different algorithms. The results show that it can provide a significant improvement in alignment quality.

Output

```
{
  "Generic": ["algorithms"],
  "Material": ["English phrases", "French sentence"],
  "Method": [],
  "Metric": ["alignment quality"],
  "OtherScientificTerm": ["cohesion constraint", "syntax-based constraint"],
  "Task": ["word alignment"]
}
```

SciFact

- Task input: An abstract from a biomedical research article, and a scientific claim.
- Task output:
 - A fact-checking verdict indicating whether the abstract supports or refutes the claim, or has no relevant information.
 - The evidence -- i.e. sentences from the abstract justifying the verdict.
- Metrics: We compute metrics for both the answer and the evidence.
 - Verdict: Label F1.
 - Evidence: Token F1 overlap with gold evidence.

Input

You will be shown a scientific claim, and the abstract of a biomedical research paper. Each sentence from the abstract will be on a separate line. Your task is to return a JSON object with two

fields:

- "verdict": The fact-checking verdict. If the information in the abstract supports the claim, write "SUPPORT". If the abstract contradicts the claim, write "CONTRADICT". If the abstract does not provide enough information to arrive at a verdict, write "NEI" (for "not enough information").
- "evidence": An array of sentences providing evidence for the verdict. Please copy all relevant sentences verbatim from the abstract. If the verdict was "NEI", then return an empty array.

For instance, if the model were given the claim "smoking causes cancer", the output might be { "verdict": "SUPPORT", "evidence": ["The results of our meta-analysis provide overwhelming support that cigarette smoking is a risk cause for lung cancer."] }

Your response should not include any text other than the json.

Claim: Therapeutics receiving accelerated approval encounter a lower frequency of post-marketing safety events

Abstract: Importance Postmarket safety events of novel pharmaceuticals and biologics occur when new safety risks are identified after initial regulatory approval of these therapeutics. These safety events can change how novel therapeutics are used in clinical practice and inform patient and clinician decision making. Objectives To characterize the frequency of postmarket safety events among novel therapeutics approved by the US Food and Drug Administration (FDA), and to examine whether any novel therapeutic characteristics known at the time of FDA approval were associated with increased risk. [Truncated for space] Biologics, psychiatric therapeutics, and accelerated and near-regulatory deadline approval were statistically significantly associated with higher rates of events, highlighting the need for continuous monitoring of the safety of novel therapeutics throughout their life cycle.

Output

```
{
  "verdict": "CONTRADICT",
  "evidence": [
    "In multivariable analysis, postmarket safety events were statistically significantly more frequent among biologics (incidence rate ratio [IRR] = 1.93; 95% CI, 1.06-3.52; P = .03), therapeutics indicated for the treatment of psychiatric disease (IRR = 3.78; 95% CI, 1.77-8.06; P < .001), those receiving accelerated approval (IRR = 2.20; 95% CI, 1.15-4.21; P = .02), and those with near-regulatory deadline approval (IRR = 1.90; 95% CI, 1.19-3.05; P = .008); events were statistically significantly less frequent among those with regulatory review times less than 200 days (IRR = 0.46; 95% CI, 0.24-0.87; P = .02)."
```

1173 **F Information About Use of AI Assistants**

1174 We use OpenAI ChatGPT and Anthropic Claude
1175 for grammar checking in manuscript preparation.