

Who’s a Better Scholar: Encoder *or* Decoder?

Anonymous ACL submission

Abstract

Language modeling has seen a tremendous development over past few years, with a considerable rise in their deployment for solving domain-specific Natural Language Processing (NLP) tasks. In recent times, the fundamental building blocks of language models are essentially composed of either an encoder-based architecture or a decoder-based architecture or a combination of both. In the scholarly domain, the majority of use cases have explored only the utilization of encoder-only models for a variety of tasks using the pre-trained model fine-tuning approach. But the same has not yet been replicated for decoder based models in spite of the recent popularity of LLMs. To address this issue, we fine-tune both encoder-based language models and decoder-based language models on an array of traditional scholarly NLP tasks. This allows us to compare the effect of learned representations in contrast to generation-based techniques on standard scholarly benchmark datasets. We conduct extensive experiments on 10 highly popular human-annotated datasets over 6 different tasks and also study the effect of domain-specific pre-training on these tasks. We achieve SOTA over two tasks using decoder-based language models, although they prove to not being best in terms of computational costs or hallucinations.

1 Introduction

Scientific literature understanding is an important facet of Natural Language Understanding and is highly useful in the comprehension of large collections of scientific text. There has been a growing interest to explore the nuances of standard Natural Language Processing tasks in the scholarly domain and in most cases the best results have come from fine-tuning a pre-trained language model (Beltagy et al., 2019; Lahiri et al., 2024; Sadat and Caragea, 2022; He et al., 2020).

Researchers have been able to classify the emergence of language models into four different waves:

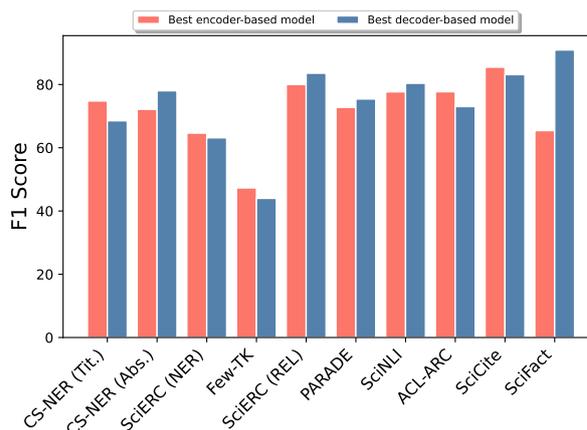


Figure 1: Comparison of the scores achieved by the best performing encoder-based and decoder-based LMs.

statistical language models, neural language models, pre-trained language models (PLMs), and large language models (LLMs) (Minaee et al., 2024). Models with tens to hundreds of billions of parameters are generally considered as LLMs and models with lesser number of parameters are referred to as PLMs. We see with LLMs the paradox of over-parametrization wherein models with greater number of parameters exhibit better performance instead of over-fitting. Decoder-based LLMs, have in fact shown to present strong emergent and reasoning capabilities (Wei et al., 2022a,b; Yao et al., 2023). The emergence of Transformer-based pre-trained language models and the subsequent popularity gained by LLMs have transformed the way we solve NLP tasks, since the language understanding capabilities of PLMs and LLMs outdo their predecessors by a large margin.

PLMs and LLMs are both categories of language models that trace their architectural roots to the original Transformer model (Vaswani et al., 2017). In theory, PLMs mainly differ from their elder siblings – the LLMs in terms of size, but may be either encoder-based or decoder-based. Encoder-based

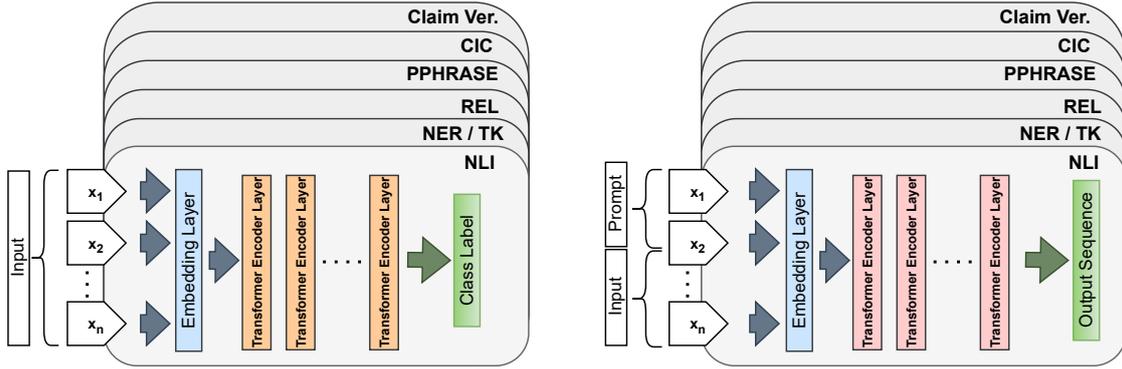


Figure 2: Fine-tuning for a Transformer encoder-based LM (left) and for a Transformer decoder-based LM (right).

models generally fall into the category of PLMs because the parameters of the available pre-trained models are in the order of millions but nowadays it is very common to find decoder-based models which have some billion parameters.

Encoder-based models like BERT (Devlin et al., 2019), although being task agnostic in nature, generally need to go through fine-tuning over a limited amount of task-specific data to achieve proficiency in that particular task. Despite the fact that LLMs possess greater emergent properties, they have been reported to be more accomplished when fine-tuned over task-specific data (Minaee et al., 2024). Moreover, simply prompting decoder-based LMs do not often produce the best results for scientific literature understanding tasks (Wadden et al., 2024).

The main objective of this paper is to create an evaluation setup that can effectively compare the ability of decoder-based LMs with that of their encoder-based LM counterparts with a special focus on scholarly tasks. To this end, we conduct extensive fine-tuning experiments on human-annotated scholarly datasets, such as Named Entity Recognition/Typed Keyphrase Recognition, Relation Classification, Natural Language Inference, Paraphrase Identification, Citation Intent Classification, and Claim Verification.

Our main contributions can be summarized as follows: a) We compare decoder-based LMs with encoder-based LMs on 10 benchmark scholarly tasks over 6 different tasks. For this purpose, we use 2 encoder-based LMs and 6 decoder-based LMs. b) We analyze the performance for each task, as well as the hallucinations generated by the models. c) We study the effect of domain-specific data in the pre-training corpus and the computational time complexity of fine-tuning these models.

2 Transformer Architecture

The original transformer architecture (Vaswani et al., 2017) consists of a combined encoder-decoder structure that is auto-regressive in nature. The encoder maps an input sequence of symbol representations (x_1, x_2, \dots, x_n) into a sequence of continuous representations $z = (z_1, z_2, \dots, z_n)$. The encoder is supposed to contain N identical layers, where each layer consists of a multi-head self-attention mechanism followed by a position-wise fully connected feed-forward network.

The decoder takes z as the input and generates an output sequence (y_1, y_2, \dots, y_m) . The decoder also consists of N identical layers, where in addition to the components of the encoder layer, there exists a new sub-layer that performs multi-head attention over the output of the encoder stack.

Most recent language models follow variants of this architecture, with small changes like the activation function, or the positional embedding technique or the tokenization procedure. With the Transformer being the basic building block, language models may contain the encoder only or the decoder only or may contain both the encoder and the decoder. Pre-training of encoder models involve various language modelling objectives like masked language modelling while decoder-based models generally use the autoregressive next token prediction objective.

Figure 2 shows the fine-tuning approach followed by encoder-based LMs as well as decoder-based LMs. For encoder-based LMs, the input is tokenized and fed into the encoder blocks to generate their token representations which are then passed through an output layer. Decoder-based LMs provide a sequential output when provided with a instruction and the input.

3 Tasks

We consider 6 tasks for our experiments, each of which is briefly described here. The details of the datasets shown in Figure 3 are in the Appendix.

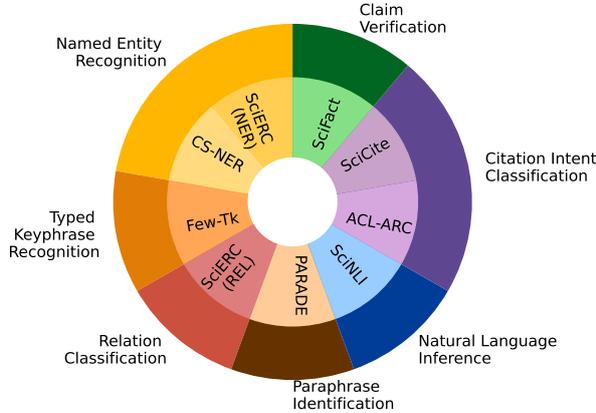


Figure 3: Tasks and Datasets

3.1 NER/TK: Named Entity Recognition/ Typed Keyphrase Recognition

Named Entity Recognition (NER) is the Information Extraction (IE) task of identifying references to rigid designators (Nadeau and Sekine, 2007). Recently (Lahiri et al., 2024) presented a broader definition for this task in the scientific domain and termed it as Typed Keyphrase Recognition.

Definition: The input is a sequence of tokens $x = (x_1, x_2, \dots, x_n)$, from which we derive a set $S = \{s_1, \dots, s_p\}$, which represents a set of semantically meaningful within-sentence contiguous sequence spans each of which is assigned a label from the set $Y = \{y_1, y_2, \dots, y_m\}$. The elements in set S may contain words, phrases or other syntactic units from the given text sequence x . Therefore, the final output can be construed as $Z = \{(s_i, y_j) : i \in 1, \dots, p; j \in 1, \dots, m; s_i \in S; y_j \in Y\}$.

3.2 REL: Relation Classification

Relation Classification is also an Information Extraction task, wherein the objective is to predict the relationship type between a given ordered pair of spans within a sentence.

Definition: The input is a sequence of tokens $x = (x_1, x_2, \dots, x_n)$ and two entities (spans), $s_A = (x_i, \dots, x_j)$ and $s_B = (x_u, \dots, x_v)$, the expected output is a triple (s_A, s_B, r) , where $r \in R$ such that R is a pre-defined set of relation labels.

3.3 PPHRASE: Paraphrase Recognition

Sentences or phrases conveying identical meaning but with the use of different wording are called paraphrases. Automated paraphrase recognition mechanisms are useful in many NLP tasks like textual entailment, machine reading, question answering, information extraction, and machine translation (Bhagat and Hovy, 2013). For the scholarly paraphrase identification task, the ability of the model to demonstrate specialized domain knowledge is tested (He et al., 2020).

Definition: A pair of sentences (s_1, s_2) are to be classified as paraphrases or non-paraphrases.

3.4 NLI: Natural Language Inference

Natural Language Inference (NLI), also known as Textual Entailment (Bowman et al., 2015; Sadat and Caragea, 2022), is the task of identifying whether there is an entailment or a contradiction between a pair of sentences or whether they are independent of each other. NLI for the scientific domain is relatively new and also quite challenging due to the difference in the vocabulary and sentence structure in comparison to the general domain.

Definition: Given a pair of sentences (s_1, s_2) , the task is to assign a label $y \in Y$ which indicates the semantic relatedness of the latter to the former.

3.5 CIC: Citation Intent Classification

Citations form an important part of scientific documents. The kind of purpose the citation serves in the scholarly document is known as its citation intent (Roman et al., 2021). Citation intents are useful in tasks like the measurement of scientific impact (Cohan et al., 2019) and the temporal study of scientific concepts (Jurgens et al., 2018).

Definition: The input is a citation sentence x and the aim is to assign a class label $y \in Y$, where Y is the set of citation intents.

3.6 CLAIM: Claim Verification

This task intends to assess the truthfulness of a claim (Vlachos and Riedel, 2014), which is important in the scientific domain due to the possibility of a far-reaching impact of a decision taken based on some scientific misinformation. We follow the simplified setting of (Vladika and Matthes, 2024) where the model is provided with golden abstracts:

Definition: Given a claim c and an evidence abstract d (each of which is a sequences of tokens), the task is to find whether c supports or refutes the abstract d .

Model	CS-NER (Titles)				CS-NER (Abstracts)			
	Precision	Recall	F1	H	Precision	Recall	F1	H
BERT	72.83	76.81	74.77	0	69.38	71.32	70.33	0
SciBERT	72.98	76.66	74.78	0	72.97	71.35	72.14	0
LLaMA-7B	66.00	70.38	68.12	1	83.29	68.18	74.98	0
LLaMA-13B	65.72	70.50	68.03	3	82.64	69.03	75.22	0
LLaMA-70B	66.41	70.61	68.45	3	90.00	62.92	74.06	0
SciLitLLM-7B	67.33	69.35	68.32	0	86.42	70.79	77.83	0
Tülu-2-dpo-7B	66.47	65.74	66.10	1	79.85	70.70	75.00	0
Tülu-2-dpo-70B	67.25	69.83	68.52	3	88.35	69.82	78.00	0

Table 1: Results for fine-tuning encoder-based LMs and instruction-tuning decoder-based LMs on CS-NER (Titles) and CS-NER (Abstracts) for Named Entity Recognition. H stands for Hallucinated Tags i.e. the tags which LLMs have generated, but are not part of the dataset’s annotation schema.

4 Experimental Setup

4.1 Encoder-based Language Models

We use the BERT (Devlin et al., 2019) and SciBERT (Beltagy et al., 2019) model checkpoints as the encoder-based LMs in our experiments. More details about these models are present in the Appendix B. The experimental details for fine-tuning encoder-based LMs are as follows:

NER/TK: We train the uncased versions of BERT and SciBERT by passing their output through a linear classifier and training using the cross-entropy loss for 20 epochs. The maximum sequence length considered is 256.

REL: This task is formulated for encoder-based LMs as a special case of text classification: the given entities are delineated with special tokens and the model learns to predict the relation between these entities (Beltagy et al., 2019).

PPHRASE: We fine-tune BERT and SciBERT by considering this task as a text classification task as was done for the original PARADE dataset (He et al., 2020). We fine-tune the backbone PLMs for 5 epochs using a learning rate of $2e - 5$.

NLI: The pair of sentences provided as input are concatenated separated by a [SEP] token between them. A softmax layer is used to predict the output class from the [CLS] token embedding. Each backbone model is trained for 5 epochs and the maximum input length is set at 300. We use the cased versions of the BERT and SciBERT models keeping in line with the original paper (Sadat and Caragea, 2022).

CIC: It is treated as a simple text classification problem given the citation sentence, as in (Beltagy et al., 2019). Therefore, the BERT vector is given as input into a linear classification layer. The learn-

ing rate is taken as $2e - 5$ and the model is trained for 5 epochs.

CLAIM: We model the claim verification task as a two-class classification problem, such that given the claim-evidence pair, the model predicts whether the claim supports or contradicts the evidence.

4.2 Transformer-decoder based models

We use the 7B, 13B and the 70B model variants of LLaMA-2 (Touvron et al., 2023b), SciLitLLM-7B¹ (Li et al., 2024) and 7B and 70B variants of Tülu-2 (Iverson et al., 2023) as the decoder-based LMs in our experiments. Details about these models and the prompts are described in Appendix B and D, respectively.

We instruction-tune the decoder-based LMs using QLoRA (Dettmers et al., 2023), which is an efficient approach for fine-tuning LLMs using relatively less GPU memory. QLoRA uses 4-bit NormalFloat, Double Quantization and Paged Optimizers on the Low-rank Adapter (LoRA) fine-tuning approach (Hu et al., 2022), which makes it possible to fine-tune even 70B parameter models in a 80GB-A100 GPU with minimal performance degradation. We fix both the source length and the target length to 512 for better comprehension. The learning rate is kept at $2e - 4$, and we fine-tune each model for 1, 875 steps.

5 Results

5.1 Named Entity Recognition

Table 3 and Table 1 shows the results obtained for the SciERC (Luan et al., 2018) as well as both the CS-NER (Abstracts) and CS-NER (Abstracts) (D’Souza and Auer, 2022) datasets. Apart from

¹<https://huggingface.co/Uni-SMART/SciLitLLM>

Model	Cmp.	Cnj.	Evl.-for	Ft.-of	Hyp.-of	Pt.-of	Used-for	F1	H
BERT	-	-	-	-	-	-	-	78.71*	0
SciBERT	-	-	-	-	-	-	-	79.97*	0
LLaMA-7B	87.32	94.4	87.01	71.54	94.03	68.38	93.67	74.54	2
LLaMA-13B	88.31	94.02	89.73	64.08	90	64.35	94.34	83.55	0
LLaMA-70B	88.57	93.02	86.34	66.67	84.93	37.97	93.66	78.74	0
SciLitLLM-7B	87.32	94.82	89.13	64.91	92.09	61.95	93.95	73.02	1
Tülu-2-dpo-7B	88.57	92.86	84.21	60.00	82.64	60	92.84	80.16	0
Tülu-2-dpo-70B	87.18	93.06	83.17	62.50	90.91	66.07	93.83	72.09	3

Table 2: Results for fine-tuning encoder-based LMs and instruction-tuning decoder-based LMs on SCIERC for Relation Classification. H stands for Hallucinated Tags, i.e., the tags which LLMs have generated, but are not part of the dataset’s annotation schema. The * denotes that the results are obtained from the original paper.

Model	P	R	F1	H
BERT	59.71	65.95	62.67	0
SciBERT	62.24	67.2	64.62	0
LLaMA-7B	58.57	61.83	60.16	4
LLaMA-13B	57.94	62.26	60.02	0
LLaMA-70B	61.42	64.95	63.14	4
SciLitLLM-7B	58.39	60.67	59.51	1
Tülu-2-dpo-7B	59.95	61.9	60.91	2
Tülu-2-dpo-70B	60.81	60.55	60.68	3

Table 3: Results for fine-tuning encoder-based LMs and instruction-tuning decoder-based LMs on SCIERC for Named Entity Recognition. H stands for Hallucinated Tags, i.e., the tags which LLMs have generated, but are not part of the dataset’s annotation schema.

Model	P	R	F1	H
BERT	40.59	45.05	42.66	0
SciBERT	46.87	47.82	47.29	0
LLaMA-7B	39.54	40.17	39.86	5
LLaMA-13B	40.51	46.12	43.13	8
LLaMA-70B	40.4	44.38	42.29	5
SciLitLLM-7B	41.47	44.96	43.15	16
Tülu-2-dpo-7B	38.36	41.48	39.86	15
Tülu-2-dpo-70B	42.55	45.54	43.99	5

Table 4: Results for fine-tuning encoder-based LMs and instruction-tuning decoder-based LMs on Few-TK for Typed Keyphrase Recognition. H stands for Hallucinated Tags, i.e., the tags which LLMs have generated, but are not part of the dataset’s annotation schema.

CS-NER (Abstracts), encoder-based LMs generally perform better than their decoder based counterparts for the NER task. Domain-specific pre-training in models like SciBERT, Tülu-2, and SciLitLLM help boost performance.

For the NER task, the generative decoder-based LMs, despite having the class names specified in the prompt, hallucinate new labels such as Objective, Scenario, Author, Profession, User, and Drug among others. We see that for CS-NER (Abstracts), none of the models hallucinate, which is perhaps due to the fact that it consists of only two classes.

5.2 Typed Keyphrase Recognition

Table 4 shows the results on the Few-TK dataset (Lahiri et al., 2024). Similar to the results for NER, here too we see that SciBERT outperforms all other models, although the results are generally low for this dataset. This is due to large number of classes, which is 38, in this dataset, that is much higher than that of other datasets in this domain. This shows that simple vanilla fine-tuning or instruction-tuning

may not be enough for more complex multi-label tasks such as these as they require significantly higher reasoning capabilities. We also see that due to the larger number of classes into which the keyphrases are to be divided, the number of hallucinations for this dataset are also much larger.

5.3 Relation Classification

Table 2 shows the results for relation classification on the SCIERC dataset and also includes the F1 scores for each class – Compare, Conjunction, Evaluate-for, Feature-of, Hyponym-Of, Part-of and Used-for. LLaMA-13B is found to be the best performing model for this task, which to the best of our knowledge is also the SOTA for relation classification on this dataset. The LLaMA-7B also performs well over the different classes in this task, but its overall performance dips due to the two hallucinated labels that it generates. Some of the hallucinated labels from generative decoder-based LMs are Induced-from, Sum-of and Weighted-sum, in the very rare cases where they hallucinate.

Model	Paraphrase	Non-paraphrase	Accuracy	Precision	Recall	F1
BERT	72.21	73.28	72.78	72.88	72.83	72.74
SciBERT	71.77	73.63	72.59	72.54	72.55	72.54
LLaMA-7B	73.69	72.18	72.96	73.39	73.20	72.93
LLaMA-13B	73.13	71.24	72.22	72.72	72.49	72.19
LLaMA-70B	73.30	77.30	75.46	75.58	75.25	75.30
SciLitLLM-7B	73.15	77.65	75.61	75.82	75.36	75.40
Tülu-2-dpo-7B	65.93	77.27	72.73	75.20	72.02	71.60
Tülu-2-dpo-70B	63.83	76.86	71.78	74.86	70.98	70.35

Table 5: Results for fine-tuning encoder-based LMs and instruction-tuning decoder-based LMs on PARADE for paraphrase recognition. We report the overall precision, recall, macro F1, accuracy and the class-wise macro F1.

Model	Contrasting	Reasoning	Entailment	Neutral	F1	Accuracy
BERT	77.17	71.25	74.37	74.01	74.20	74.27
SciBERT	79.69	74.35	74.35	76.46	77.68	77.67
LLaMA-7B	78.22	69.53	73.53	61.05	70.58	71.10
LLaMA-13B	82.92	74.93	77.60	71.71	76.79	76.98
LLaMA-70B	86.17	74.45	77.77	64.51	75.73	76.50
SciLitLLM-7B	82.54	76.52	77.06	69.77	76.47	76.80
Tülu-2-dpo-7B	79.82	71.03	74.87	63.86	72.39	72.85
Tülu-2-dpo-70B	87.24	78.22	79.20	76.23	80.22	80.37

Table 6: Results for fine-tuning encoder-based LMs and instruction-tuning decoder-based LMs on SciNLI for Natural Language Inference. We report the overall macro F1, accuracy and the class-wise macro F1.

5.4 Paraphrase Recognition

Table 5 shows the results for the task of paraphrase recognition. Although the results achieved by each of the models are very close to each other, decoder-based LMs hold a slight edge in performance over encoder-based LMs, with the SciLitLLM-7B being the best performing model by outperforming even the 70B models.

5.5 Natural Language Inference

Table 6 shows the results for scientific Natural Language Inference. The Tülu-2-dpo-70B model shows superior performance among the tested models and also achieves the SOTA performance on this dataset (Sadat and Caragea, 2024).

5.6 Citation Intent Classification

Table 7 and Table 8 shows the result for Citation Intent Classification on the ACL-ARC (Jurgens et al., 2018) and SciCite (Cohan et al., 2019) datasets, respectively. We see that for both the datasets SciBERT shows better performance. Only for F1 scores of two classes of the ACL-ARC dataset and the overall accuracy score, other language models are able to perform better than SciBERT. LLaMA-70B and Tülu-2-dpo-70B – both 70B LLMs clock al-

most about the same overall F1 score, whereas the two 7B models show some hallucinations like Repeats and Inspired.

5.7 Claim Verification

Table 9 shows the result for Claim Verification on the SCIFACT dataset (Wadden et al., 2020). This is the only task where we find that a large language model i.e. the Tülu-2-dpo-70B model is the best performing model on all metrics and is also separated from the encoder-based LMs by a huge margin.

6 Performance Analysis

We find that encoder-based LMs offer stiff competition to their decoder-based counterparts even though the encoder-based LMs are quite smaller in size and trained on much less data. Decoder-based LMs perform well in those tasks where the number of labels or classification heads are less than or equal to 3. Among the tasks considered, decoder-based LMs have been found to work well in tasks like Paraphrase Recognition, Natural Language Inference and Claim Verification.

(Wadden et al., 2024) reports the F1 score in the SciERC using GPT-4 to be 42.2 and using their

Model	Bckg.	Comp.	Extends	Future	Motiv.	Uses	Accuracy	F1	H
BERT	84.12	59.15	44.81	21.67	00.00	64.91	45.78	70.74	0
SciBERT	87.67	73.76	73.13	76.26	41.79	78.42	74.96	77.70	0
LLaMA-7B	84.62	60.00	61.54	50.00	71.43	84.44	77.70	58.86	2
LLaMA-13B	86.09	68.18	50.00	66.67	40.00	80.77	78.42	65.29	0
LLaMA-70B	84.97	63.41	72.73	80.00	26.67	79.17	76.98	67.82	0
SciLitLLM-7B	84.00	60.47	61.54	72.73	36.36	76.00	75.54	65.18	0
Tülu-2-dpo-7B	84.93	60.00	46.15	72.73	44.44	77.55	74.82	55.12	1
Tülu-2-dpo-70B	84.97	61.90	80.00	72.73	53.33	85.11	79.14	73.01	0

Table 7: Results for fine-tuning encoder-based LMs and instruction-tuning decoder-based LMs on ACL-ARC for Citation Intent Classification. We report the overall macro F1, accuracy and the class-wise macro F1. H stands for Hallucinated Tags i.e. the tags which LMs have generated, but are not part of the dataset’s annotation schema.

Model	Background	Method	Result	Accuracy	F1
BERT	88.28	85.28	80.6	86.17	84.72
SciBERT	88.51	86.33	81.53	86.75	85.46
LLaMA-7B	85.85	81.44	77.96	83.37	81.75
LLaMA-13B	85.31	80.28	77.12	82.56	80.90
LLaMA-70B	86.83	82.58	79.92	84.55	83.11
SciLitLLM-7B	86.10	81.02	79.06	83.48	82.06
Tülu-2-dpo-7B	86.54	82.41	76.73	83.80	81.89
Tülu-2-dpo-70B	86.19	83.09	80.00	84.23	83.10

Table 8: Results for fine-tuning encoder-based LMs and instruction-tuning decoder-based LMs on SciCite for Citation Intent Classification. We report the overall macro F1, accuracy and the class-wise macro F1.

own SciTÜLU 70B model to be 35.9. Therefore, we see that fine-tuning decoder-based LMs gives far better results than the simply prompting.

We see that many of the decoder-based LMs hallucinate when there are too many labels for classification. Hallucinations are a major reason for the overall decrease in performance of decoder-based LMs in many tasks. We postulate that the pre-training of large generative models plays a major part in such hallucinations, where in spite of the classes being mentioned in the training prompt, the model in a few exceptional cases generates data which is meaningful but does not pertain to the constrained framework of the given task.

On the bright side, our experiments on decoder-based LMs have led to achieving SOTA performance on two tasks – Relation Classification and Natural Language Inference.

6.1 Computational Time Complexity

Encoder-based LMs take much lower time for both training and inferencing than decoder-based LMs, which require anywhere about 4 to 26 A100 GPU hours per dataset only for the training part. Apart from this, the inferencing stage is also a time-

consuming process with datasets like CS-NER which have large amounts of test data requiring more than 12 hours on an A100 GPU. In comparison, encoder-based LMs require at most 5-6 hours for the completion of both the training and inferencing stages. SciLitLLM (Li et al., 2024) takes an inordinately large amount of time for the inferencing phase in spite of its model size.

6.2 Effect of using domain-specific pre-trained models

We see across all tasks that language models that have been pre-trained on scholarly data perform better than those trained on general domain data. We observe this trend both in the case of encoder-based models (SciBERT) and decoder-based models (SciLitLLM and Tülu-2). But, we notice an interesting scenario in the case of Tülu-2: SCI-ERC (one of our NER and relation classification datasets) is included within its pre-training data and even after explicitly fine-tuning on the same data, we do not obtain an improvement in the results. Yet, although SciFact occurs in Tülu-2 pre-training corpus, hallucinations do not occur during claim verification on SciFact. Therefore, we again con-

Model	Support	Contradict	Accuracy	Precision	Recall	F1
BERT	77.14	00.52	62.82	34.15	49.21	38.83
SciBERT	80.22	53.15	69.82	66.89	65.15	65.41
LLaMA-7B	81.87	51.89	73.67	74.64	66.20	66.88
LLaMA-13B	85.59	71.11	80.77	79.90	77.46	78.35
LLaMA-70B	90.20	79.26	86.69	87.86	83.16	84.73
SciLitLLM-7B	85.27	69.68	80.18	79.47	76.46	77.48
Tülu-2-dpo-7B	83.41	67.83	78.11	76.55	75.02	75.62
Tülu-2-dpo-70B	93.08	88.72	91.42	90.25	91.86	90.9

Table 9: Results for fine-tuning encoder-based LMs and instruction-tuning decoder-based LMs on SciFact for Claim Verification. We report the overall precision, recall, macro F1, accuracy and the class-wise macro F1.

clude that hallucinations play a large role in the performance of decoder-based models.

6.3 Experimental Setup Analysis

We do not opt for multi-task fine-tuning of LLMs as we have chosen a diverse range of tasks and therefore, there is a high possibility of negative transfer even though multi-task fine-tuning is a viable option sometimes while dealing with related tasks (Karimi Mahabadi et al., 2021).

We choose BERT (Devlin et al., 2019) over other variants of Transformer encoder based model variants because other architecturally similar models do not show any drastic improvement in performance over BERT and also because of the popularity of BERT on standard NLP tasks. We do not use the SCITÜLU (Wadden et al., 2024) checkpoints for our experiments as most of the datasets overlap with their training data and this would not have been suitable for our experiments.

7 Related Work

A series of instruction-tuned models have been built on LLaMA (Touvron et al., 2023a) and LLaMA-2 (Touvron et al., 2023b) including Code LLaMA (Rozière et al., 2024), Gorilla (Patil et al., 2023), Giraffe (Pal et al., 2023), Vigogne (Huang, 2023), Tülu (Wang et al., 2023), Tülu-2 (Iverson et al., 2023), Long LLaMA (Tworkowski et al., 2023), and Stable Beluga2 (Mahan et al.).

Galactica (Taylor et al., 2022), DARWIN (Xie et al., 2023), SCITÜLU (Wadden et al., 2024) and SciLitLLM (Li et al., 2024) are some recently developed LLMs that have scientific knowledge injected into them and are able to perform better than general-domain LLMs on scientific tasks.

(AI4Science and Quantum, 2023) explores the performance of GPT-4 on a range of scientific

domains, SCIBENCH (Wang et al., 2024) is a benchmark for examining the reasoning capabilities of LLMs, SciEval (Sun et al., 2024) contains 18,000 objective and subjective questions for evaluating the scientific reasoning capabilities of LLMs. Domain-specific evaluation of LLMs has been carried out in areas like chemistry (Castro Nascimento and Pimentel, 2023) (Guo et al., 2024), molecular discovery (Janakarajan et al., 2024), biomedicine (Jahan et al., 2024), biological protocol planning (O’Donoghue et al., 2023) and material science (Jablonka et al., 2023). These studies mainly examine only the zero-shot, few-shot and chain-of-thought inferencing capabilities of LLMs, whereas our study highlights the difference of fine-tuning encoder-based LMs with decoder-based LMs. With respect to scientific literature understanding, perhaps the closest work to ours is the SCIRIFF (Wadden et al., 2024), which creates an instruction-tuning dataset for scientific literature understanding and fine-tunes the TÜLU V2 checkpoint on the dataset to finally create a set of models called SCITÜLU. In contrast, our work is more aligned towards the evaluation of decoder-based LMs and encoder-based LMs.

8 Conclusion

We fine-tune and examine 2 encoder-based language models and 6 decoder-based language models on 10 benchmark scholarly datasets over a span of 6 tasks. We observe that there is no clear winner among these two groups of models. In the case of decoder-based language models, we find that there is a huge dissimilarity between the performance achieved and the computational costs involved. We also report the usefulness of fine-tuning and using domain-specific large language models.

502 Limitations

503 We do not test over different prompt templates
504 due to computational costs. Moreover, using more
505 prompt engineering and using more latest decoder-
506 based language models can be tested for these
507 tasks.

508 References

509 Microsoft Research AI4Science and Microsoft Azure
510 Quantum. 2023. [The impact of large language mod-
511 els on scientific discovery: a preliminary study using
512 gpt-4](#). *Preprint*, arXiv:2311.07361.

513 Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. [SciB-
514 ERT: A pretrained language model for scientific text](#).
515 In *Proceedings of the 2019 Conference on Empirical
516 Methods in Natural Language Processing and the
517 9th International Joint Conference on Natural Lan-
518 guage Processing (EMNLP-IJCNLP)*, pages 3615–
519 3620, Hong Kong, China. Association for Computa-
520 tional Linguistics.

521 Rahul Bhagat and Eduard Hovy. 2013. [What Is a Para-
522 phrase?](#) *Computational Linguistics*, 39(3):463–472.

523 Samuel R. Bowman, Gabor Angeli, Christopher Potts,
524 and Christopher D. Manning. 2015. [A large anno-
525 tated corpus for learning natural language inference](#).
526 In *Proceedings of the 2015 Conference on Empiri-
527 cal Methods in Natural Language Processing*, pages
528 632–642, Lisbon, Portugal. Association for Computa-
529 tional Linguistics.

530 Isabel Cachola, Kyle Lo, Arman Cohan, and Daniel
531 Weld. 2020. [TLDR: Extreme summarization of sci-
532 entific documents](#). In *Findings of the Association
533 for Computational Linguistics: EMNLP 2020*, pages
534 4766–4777, Online. Association for Computational
535 Linguistics.

536 Cayque Monteiro Castro Nascimento and André Silva
537 Pimentel. 2023. [Do large language models un-
538 derstand chemistry? a conversation with chatgpt](#).
539 *Journal of Chemical Information and Modeling*,
540 63(6):1649–1655. PMID: 36926868.

541 Arman Cohan, Waleed Ammar, Madeleine van Zuylen,
542 and Field Cady. 2019. [Structural scaffolds for ci-
543 tation intent classification in scientific publications](#).
544 In *Proceedings of the 2019 Conference of the North
545 American Chapter of the Association for Computa-
546 tional Linguistics: Human Language Technologies,
547 Volume 1 (Long and Short Papers)*, pages 3586–3596,
548 Minneapolis, Minnesota. Association for Computa-
549 tional Linguistics.

550 Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan,
551 Noah A. Smith, and Matt Gardner. 2021. [A dataset
552 of information-seeking questions and answers an-
553 chored in research papers](#). In *Proceedings of the
554 2021 Conference of the North American Chapter of*

*the Association for Computational Linguistics: Hu-
555 man Language Technologies*, pages 4599–4610, On-
556 line. Association for Computational Linguistics. 557

558 Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and
559 Luke Zettlemoyer. 2023. [Qlora: Efficient finetuning
560 of quantized llms](#). In *Advances in Neural Information
561 Processing Systems*, volume 36, pages 10088–10115.
562 Curran Associates, Inc.

563 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
564 Kristina Toutanova. 2019. [BERT: Pre-training of
565 deep bidirectional transformers for language under-
566 standing](#). In *Proceedings of the 2019 Conference of
567 the North American Chapter of the Association for
568 Computational Linguistics: Human Language Tech-
569 nologies, Volume 1 (Long and Short Papers)*, pages
570 4171–4186, Minneapolis, Minnesota. Association for
571 Computational Linguistics.

572 Jennifer D’Souza and Sören Auer. 2022. Computer
573 science named entity recognition in the open research
574 knowledge graph. *arXiv preprint arXiv:2203.14579*.

575 Taicheng Guo, Kehan Guo, Bozhao Nan, Zhenwen
576 Liang, Zhichun Guo, Nitesh V. Chawla, Olaf Wiest,
577 and Xiangliang Zhang. 2024. [What can large lan-
578 guage models do in chemistry? a comprehensive
579 benchmark on eight tasks](#). In *Proceedings of the 37th
580 International Conference on Neural Information Pro-
581 cessing Systems, NIPS ’23*, Red Hook, NY, USA.
582 Curran Associates Inc.

583 Yun He, Zhuoer Wang, Yin Zhang, Ruihong Huang, and
584 James Caverlee. 2020. [PARADE: A New Dataset
585 for Paraphrase Identification Requiring Computer
586 Science Domain Knowledge](#). In *Proceedings of the
587 2020 Conference on Empirical Methods in Natural
588 Language Processing (EMNLP)*, pages 7572–7582,
589 Online. Association for Computational Linguistics.

590 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan
591 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and
592 Weizhu Chen. 2022. [LoRA: Low-rank adaptation of
593 large language models](#). In *International Conference
594 on Learning Representations*.

595 Bofeng Huang. 2023. [Vigogne: French instruction-
596 following and chat models](#). [https://github.com/
597 bofenghuang/vigogne](https://github.com/bofenghuang/vigogne).

598 Hamish Ivison, Yizhong Wang, Valentina Pyatkin,
599 Nathan Lambert, Matthew Peters, Pradeep Dasigi,
600 Joel Jang, David Wadden, Noah A. Smith, Iz Belt-
601 agy, and Hannaneh Hajishirzi. 2023. [Camels in a
602 changing climate: Enhancing lm adaptation with tulu
603 2](#). *Preprint*, arXiv:2311.10702.

604 Kevin Maik Jablonka, Qianxiang Ai, Alexander Al-
605 Feghali, Shruti Badhwar, Joshua D. Bocarsly, An-
606 dres M. Bran, Stefan Bringuier, L. Catherine Brinson,
607 Kamal Choudhary, Defne Circi, Sam Cox, Wibe A.
608 de Jong, Matthew L. Evans, Nicolas Gastellu,
609 Jerome Genzling, María Victoria Gil, Ankur K.
610 Gupta, Zhi Hong, Alishba Imran, Sabine Kruschwitz,
611 Anne Labarre, Jakub Lála, Tao Liu, Steven Ma,

612	Sauradeep Majumdar, Garrett W. Merz, Nicolas Moitessier, Elias Moubarak, Beatriz Mouriño, Brenden Pelkie, Michael Pieler, Mayk Caldas Ramos, Bojana Ranković, Samuel G. Rodrigues, Jacob N. Sanders, Philippe Schwaller, Marcus Schwarting, Jiale Shi, Berend Smit, Ben E. Smith, Joren Van Herck, Christoph Völker, Logan Ward, Sean Warren, Benjamin Weiser, Sylvester Zhang, Xiaoqi Zhang, Ghezal Ahmad Zia, Aristana Scourtas, K. J. Schmidt, Ian Foster, Andrew D. White, and Ben Blaiszik. 2023. 14 examples of how llms can transform materials science and chemistry: a reflection on a large language model hackathon . <i>Digital Discovery</i> , 2:1233–1250.	
626	Israt Jahan, Md Tahmid Rahman Laskar, Chun Peng, and Jimmy Xiangji Huang. 2024. A comprehensive evaluation of large language models on benchmark biomedical text processing tasks . <i>Computers in Biology and Medicine</i> , 171:108189.	
631	Nikita Janakarajan, Tim Erdmann, Sarath Swaminathan, Teodoro Laino, and Jannis Born. 2024. Language Models in Molecular Discovery , pages 121–141. Springer Nature Singapore, Singapore.	
635	David Jurgens, Srijan Kumar, Raine Hoover, Dan McFarland, and Dan Jurafsky. 2018. Measuring the Evolution of a Scientific Field through Citation Frames . <i>Transactions of the Association for Computational Linguistics</i> , 6:391–406.	
640	Rabeeh Karimi Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. 2021. Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 565–576, Online. Association for Computational Linguistics.	
649	Avishek Lahiri, Pratyay Sarkar, Medha Sen, Debarshi Kumar Sanyal, and Imon Mukherjee. 2024. Few-TK: A dataset for few-shot scientific typed keyphrase recognition . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 4011–4025, Mexico City, Mexico. Association for Computational Linguistics.	
656	Sihang Li, Jin Huang, Jiaxi Zhuang, Yaorui Shi, Xiaochen Cai, Mingjun Xu, Xiang Wang, Linfeng Zhang, Guolin Ke, and Hengxing Cai. 2024. Scilitlm: How to adapt llms for scientific literature understanding . <i>Preprint</i> , arXiv:2408.15545.	
661	Yi Luan, Luheng He, Mari Ostendorf, and Hannaneh Hajishirzi. 2018. Multi-task identification of entities, relations, and coreference for scientific knowledge graph construction . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 3219–3232, Brussels, Belgium. Association for Computational Linguistics.	
668	Dakota Mahan, Ryan Carlow, Louis Castricato, Nathan Cooper, and Christian Laforte. Stable beluga models .	
	Shervin Minaee, Tomas Mikolov, Narjes Nikzad, Meysam Chenaghlu, Richard Socher, Xavier Amatriain, and Jianfeng Gao. 2024. Large language models: A survey . <i>Preprint</i> , arXiv:2402.06196.	670 671 672 673
	David Nadeau and Satoshi Sekine. 2007. A survey of named entity recognition and classification . <i>Linguisticae Investigationes</i> , 30:3–26.	674 675 676
	Odhran O’Donoghue, Aleksandar Shtedritski, John Ginger, Ralph Abboud, Ali Ghareeb, and Samuel Rodrigues. 2023. BioPlanner: Automatic evaluation of LLMs on protocol planning in biology . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 2676–2694, Singapore. Association for Computational Linguistics.	677 678 679 680 681 682 683 684
	Arka Pal, Deep Karkhanis, Manley Roberts, Samuel Dooley, Arvind Sundararajan, and Siddhartha Naidu. 2023. Giraffe: Adventures in expanding context lengths in llms . <i>Preprint</i> , arXiv:2308.10882.	685 686 687 688
	Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. 2023. Gorilla: Large language model connected with massive apis . <i>Preprint</i> , arXiv:2305.15334.	689 690 691 692
	Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2023. Direct preference optimization: Your language model is secretly a reward model . In <i>Advances in Neural Information Processing Systems</i> , volume 36, pages 53728–53741. Curran Associates, Inc.	693 694 695 696 697 698
	Muhammad Roman, Abdul Shahid, Shafiqullah Khan, Anis Koubâa, and Lisu Yu. 2021. Citation intent classification using word embedding . <i>IEEE Access</i> , 9:9982–9995.	699 700 701 702
	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2024. Code llama: Open foundation models for code . <i>Preprint</i> , arXiv:2308.12950.	703 704 705 706 707 708 709 710 711 712
	Mobashir Sadat and Cornelia Caragea. 2022. SciNLI: A corpus for natural language inference on scientific text . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7399–7409, Dublin, Ireland. Association for Computational Linguistics.	713 714 715 716 717 718
	Mobashir Sadat and Cornelia Caragea. 2024. Co-training for low resource scientific natural language inference . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2538–2550, Bangkok, Thailand. Association for Computational Linguistics.	719 720 721 722 723 724 725

726	Noam Shazeer. 2020. GLU variants improve transformer . <i>CoRR</i> , abs/2002.05202.		
727			
728	Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. 2021. Roformer: Enhanced transformer with rotary position embedding . <i>CoRR</i> , abs/2104.09864.		
729			
730			
731	Liangtai Sun, Yang Han, Zihan Zhao, Da Ma, Zhen-nan Shen, Baocai Chen, Lu Chen, and Kai Yu. 2024. Scieval: A multi-level large language model evaluation benchmark for scientific research . <i>Proceedings of the AAI Conference on Artificial Intelligence</i> , 38(17):19053–19061.		
732			
733			
734			
735			
736			
737	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>Preprint</i> , arXiv:2211.09085.		
738			
739			
740			
741			
742	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023a. Llama: Open and efficient foundation language models . <i>Preprint</i> , arXiv:2302.13971.		
743			
744			
745			
746			
747			
748			
749	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023b. Llama 2: Open foundation and fine-tuned chat models . <i>Preprint</i> , arXiv:2307.09288.		
750			
751			
752			
753			
754			
755			
756			
757			
758			
759			
760			
761			
762			
763			
764			
765			
766			
767			
768			
769			
770			
771			
772	Szymon Tworkowski, Konrad Staniszewski, Mikołaj Pacek, Yuhuai Wu, Henryk Michalewski, and Piotr Miłoś. 2023. Focused transformer: Contrastive training for context scaling . <i>Preprint</i> , arXiv:2307.03170.		
773			
774			
775			
776	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need . In <i>Advances in Neural Information Processing Systems</i> , volume 30. Curran Associates, Inc.		
777			
778			
779			
780			
781	Andreas Vlachos and Sebastian Riedel. 2014. Fact checking: Task definition and dataset construction . In <i>Proceedings of the ACL 2014 Workshop on Language Technologies and Computational Social Science</i> , pages 18–22, Baltimore, MD, USA. Association for Computational Linguistics.		783
782			784
			785
			786
			787
			788
			789
			790
			791
			792
			793
			794
			795
			796
			797
			798
			799
			800
			801
			802
			803
			804
			805
			806
			807
			808
			809
			810
			811
			812
			813
			814
			815
			816
			817
			818
			819
			820
			821
			822
			823
			824
			825
			826
			827
			828
			829
			830
			831
			832
			833
			834
			835
			836
			837
			838
			839
			840

841	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,	(three categories) (Cohan et al., 2019). SciCite con-	887
842	Tom Griffiths, Yuan Cao, and Karthik Narasimhan.	sists of 11, 020 instances and is larger than ACL-	888
843	2023. Tree of thoughts: Deliberate problem solving	ARC which contains 1, 941 data points.	889
844	with large language models . In <i>Advances in Neural</i>		
845	<i>Information Processing Systems</i> , volume 36, pages		
846	11809–11822. Curran Associates, Inc.		
847	Biao Zhang and Rico Sennrich. 2019. Root mean square	A.6 Claim Verification	890
848	layer normalization . In <i>Advances in Neural Informa-</i>	SciFACT (Wadden et al., 2020) is a dataset that is	891
849	<i>tion Processing Systems</i> , volume 32. Curran Asso-	made up of 1, 409 expert-written scientific claims	892
850	ciates, Inc.	which are verified against a corpus of 5, 183 ab-	893
		stracts. The claims in this dataset	894
851	A Dataset	B Model Checkpoints	895
852	A.1 Named Entity Recognition/ Typed	B.1 BERT	896
853	Keyphrase Recognition	BERT (Devlin et al., 2019) stands for Bidirec-	897
854	We make use of the following popular datasets for	tional Encoder Representations from Transformers.	898
855	Named Entity Recognition: SCIERC (Luan et al.,	BERT is a multi-layer bidirectional Transformer	899
856	2018), CS-NER (Abstracts) (D’Souza and Auer,	encoder model that is pre-trained on unlabelled	900
857	2022), CS-NER (Abstracts) (D’Souza and Auer,	data from the BooksCorpus and English Wikipedia	901
858	2022). For the Typed Keyphrase Extraction task,	for two different tasks: the masked language mod-	902
859	we use FEW-TK (Lahiri et al., 2024). Almost all	elling (MLM) task and the next sentence prediction	903
860	of these datasets are annotated on research paper	(NSP) task. The BERT model may be fine-tuned	904
861	abstracts or titles or both.	for several downstream tasks and this fine-tuning	905
862	A.2 Relation Classification	paradigm has found success in almost all major	906
863	We use SCIERC (Luan et al., 2018), which con-	NLP tasks.	907
864	tains about 4, 716 relations over 500 scientific doc-	B.2 SciBERT	908
865	ument abstracts.	SciBERT (Beltagy et al., 2019) is domain-specific	909
866	A.3 Paraphrase Recognition	variant of BERT that is pre-trained on scientific	910
867	PARADE (PARAphrase identification based on Do-	text. SciBERT retains the architecture as well as all	911
868	main knowledge) (He et al., 2020) is a dataset	the major characteristics of BERT except that it is	912
869	tailored for paraphrase identification consisting of	pre-trained on a corpus that consists of papers from	913
870	10, 182 pairs of definitions that describe 788 dis-	the biomedical domain and the computer science	914
871	ting entities in the Computer Science domain. Out	domain in a 82 : 18 ratio.	915
872	of these, 4, 778 are paraphrases and 5, 404 are non-	B.3 LLaMA family of models	916
873	paraphrases.	LLaMA is a family of pre-trained foundational lan-	917
874	A.4 Natural Language Inference	guage models that have been open-sourced by Meta	918
875	SciNLI (Sadat and Caragea, 2022) is a Natural	in recent times. LLaMA models incorporates the	919
876	Language Inference (NLI) dataset tailored for the	following three minor architectural changes within	920
877	scientific domain, consisting of 101,412 samples in	the original Transformer architecture (Vaswani	921
878	the training set, 2,000 samples in the validation set,	et al., 2017): (1) use of SwiGLU (Shazeer, 2020)	922
879	and 4,000 samples in the test set. In comparison	activation function instead of ReLU, (2) use of ro-	923
880	to traditional datasets, this dataset contains two	tary positional embeddings (Su et al., 2021) instead	924
881	new classes, taking the total number of classes to	of absolute positional embedding, and, (3) use of	925
882	four: "Contrasting", "Entailment", "Reasoning"	RMSNorm (Zhang and Sennrich, 2019) normaliz-	926
883	and "Neutral".	ing function instead of layer-normalization.	927
884	A.5 Citation Intent Classification	B.4 SciLitLLM	928
885	We consider two datasets for this task: ACL-ARC	SciLitLLM (Li et al., 2024) is a very recently re-	929
886	(six categories) (Jurgens et al., 2018) and SciCite	leased LLM designed for the task of scientific litera-	930
		ture understanding that has been trained using both	931

Corpora	Domain	Classes	Papers	Tokens	Entities
SciERC (Luan et al., 2018)	AI	5	500	60,749	8,089
CS-NER (Abstracts) (D’Souza and Auer, 2022)	AI	2	12,271	1,317,256	29,273
CS-NER (Titles) (D’Souza and Auer, 2022)	CL	7	31,044	263,143	67,270
FEW-TK (Lahiri et al., 2024)	AI	38	500	115,745	20064

Table 10: Details of standard scientific-domain Named Entity Recognition datasets and FEW-TK for Typed Keyphrase Recognition

continual pre-training (CPT) and supervised fine-tuning (SFT). This strategy is used on Qwen2.5 to obtain SciLitLLM. The CPT stage uses 73,000 textbooks and 625,000 academic papers, while the SFT stage uses SciLitIns, SciRIFF (Wadden et al., 2024) and Infinity-Instruct². We use the SciLitLLM 7B³ for our experimental purposes.

B.5 Tülu family of models

Tülu (Wang et al., 2023) is a set of models that are instruction-tuned on LLaMA (Touvron et al., 2023a) using a mixture of human-generated as well as GPT-generated data. Tülu-2 (Iverson et al., 2023) is trained on LLaMA-2 over a more updated and refined data mixture, which contains even datasets from scientific literature like SciERC (Luan et al., 2018), Qasper (Dasigi et al., 2021), SciFact (Wadden et al., 2020) and SciTLDR (Cachola et al., 2020). Tülu-2 is further trained using the direct preference optimization (DPO) algorithm (Rafailov et al., 2023).

C Hallucinated Labels

The following tables show the hallucinated labels in different decoder-based language models.

Model	SciERC (REL)
LLaMA-7B	COMBINATION-STRATEGY -OVER, WEIGHTED-SUM.
LLaMA-13B	-
LLaMA-70B	-
SciLitLLM-7B	INDUCED-FROM
Tulu-2-dpo-7B	-
Tulu-2-dpo-70B	FOR-FOR, SUM-OF, OUT-OF-NLP.

Table 11: Hallucinated Labels for Relation Extraction datasets

²<https://huggingface.co/datasets/BAAI/Infinity-Instruct>

³<https://huggingface.co/Uni-SMART/SciLitLLM>

Model	ACL-ARC
LLaMA-7B	INSPIRED, TUV
LLaMA-13B	-
LLaMA-70B	-
SciLitLLM-7B	-
Tulu-2-dpo-7B	-
Tulu-2-dpo-70B	REPEATS

Table 12: Hallucinated Labels for Citation Intent Classification datasets

D Prompt Template

Table 15 shows the prompt templates used by the generative decoder-based language models.

Model	Few-TK
LLaMA-7B	'Data Mining Information Retrieval metrics', 'Compute architecture', 'Data Mining' 'Information Retrieval dataset', 'Statistical Mathematical domain', 'Statistical Mathematical phenomenon'
LLaMA-13B	'Astronomy term', 'Astronomy term', 'Astronomy term', 'Astronomy term', 'Statistical Mathematical domain', 'Statistical Mathematical technique', 'Statistical Mathematical domain', 'Bioinformatics algorithm tool'
LLaMA-70B	'Garbage value: Tourism is the typed keyphrase identified from the given text.', 'Statistical Mathematical focus', 'Statistical Mathematical domain', 'New York City dog park', 'AI ML DL metrics'
SciLitLLM-7B	'Reference', 'Optimization algorithm tool', 'Data Mining Information Retrieval dataset', 'AI ML DL library', 'Q&A site for programmers', 'Commercial LP solver', 'Data Mining Information Retrieval dataset', 'Miscellaneous result', 'Data Mining Information Retrieval strategy', 'Statistical Mathematical focus', 'Statistical Mathematical domain', 'NLP author', 'NLP author', 'Information Retrieval focus', 'Garbage value: 600 words of type'
Tulu-2-dpo-7B	'Miscellaneous dataset', 'Miscellaneous dataset', 'Miscellaneous result', 'Statistical Mathematical focus', 'Statistical Mathematical focus', 'Data Mining Information Retrieval dataset', 'Computer vision algorithm step', 'Financial term', 'Quality metrics', 'Statistical Mathematical focus', 'Statistical Mathematical discipline', 'author', 'author', 'Information retrieval focus', 'Statistical Mathematical focus'
Tulu-2-dpo-70B	'Application term', 'Computer Vision algorithm tool', 'Data Mining Information Retrieval tool', 'Miscellaneous dataset', 'NLP framework'

Table 13: Hallucinated Labels for Typed Keyphrase Recognition dataset, Few-TK

Model	CS-NER (Titles)	SciERC (NER)
LLaMA-7B	AUTHOR	OBJECTIVE, SCENARIO, AUTHOR
LLaMA-13B	DATE	-
LLaMA-70B	AUTHOR, R, REGION	AUTHOR, HUMAN
SciLitLLM-7B	-	PROFESSION
Tulu-2-dpo-7B	DATE	FUNCTION, AUTHOR
Tulu-2-dpo-70B	DATE, REGION, DATE	USER, PLATFORM, DRUG

Table 14: Hallucinated Labels for Named Entity Recognition datasets

Task	Instruction	Input	Output
Named Entity Recognition	In the given sentence, find the named entity mentions and classify them among the following possible categories - Y	X	The entities s_i of type y_i are identified from the given text.
Typed Keyphrase Recognition	In the given sentence, find the typed keyphrase mentions and classify them among the following possible categories - Y	X	The typed keyphrases s_i of type y_i are identified from the given text.
Relation Extraction	In the given sentence, find and classify the relation between the mentioned pair of named entities, where the relation can be of the following types: Y	X	The relation between s_A and s_B is r .
Paraphrase Recognition	Paraphrases are sentences that express the same meaning by using different wording. Are the following pair of sentences paraphrases or non-paraphrases? SEP separates the two sentences.	(s_1, s_2)	The given pair of sentences are paraphrases/non-paraphrases.
Natural Language Inference	Analyze the provided pair of sentences to determine their relationship. Choose one of the following categories: Y	(s_1, s_2)	$y \in Y$
Citation Intent Classification	Given a scientific text containing a citation and the citation string, classify the intent of the citation among the following categories: Y .	X	The intent of the citation falls under the following category: $y \in Y$
Claim Verification	Given a scientific claim, evaluate the evidence to determine whether it supports or refutes the claim.	(s_1, s_2)	The given evidence supports/refutes the scientific claim.

Table 15: Table showing prompts used to instruction-tune LLMs