

# AI-FIGURES: A Fine-grained Task-oriented Dataset for Multimodal Scientific Literature Understanding

Anonymous ACL submission

## Abstract

Diagrams and figures are a powerful medium of communication in scientific research. There is a recent spark in interest in the development of Machine Learning-driven applications involving scientific figures such as multimodal question answering, multimodal document retrieval, text-to-image generation, or image captioning. Challenging tasks in this domain may be dependent only on a specific category of scientific figures. But there are no datasets in prior literature that provide a domain-specific, broad classification of scientific figures. To fill this gap, we introduce AI-FIGURES, a large-scale dataset containing scientific figure-caption pairs which are classified into 8 different categories. We create this dataset by leveraging the idea of image segmentation and classification using the YOLO model. Our automated data acquisition pipeline can also be implemented on other datasets also to classify their figures. We benchmark our dataset for various tasks such as figure captioning, text-to-figure generation, scholarly multimodal question answering, and multimodal document retrieval using various vision-based models. We show that there is a significant increase in a model’s inference capabilities when we finetune it on targeted classes of our dataset. Our dataset and code will be made public upon acceptance.

## 1 Introduction

Images create a visual imprint on our brain that is immediately able to trigger the human perceptual system to process the simultaneous conceptual representation. Recognizing the intrinsic importance of figures, recent research endeavors have underscored the necessity of developing robust systems capable of extracting and interpreting these visual elements.

Images serve as vital elements in conveying crucial aspects of scholarly content, such as methodological explanations, experimental results, and

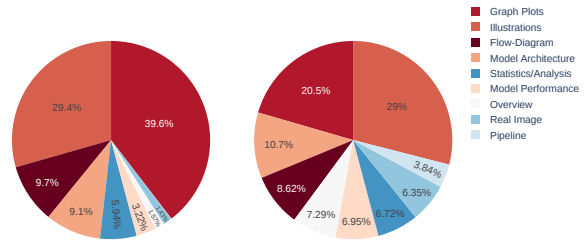


Figure 1: The class-wise distribution in the human annotated and the distantly-supervised AI-FIGURES dataset.

comparative analyses. Despite the huge availability of visual data in the scholarly domain, the advancement of automated multimodal scholarly generative tasks like figure captioning, text-to-figure, scholarly question answering and multimodal scholarly search is in a much more nascent stage than their general domain counterparts like text-to image generation (Xu et al., 2018; Ramesh et al., 2021; Saharia et al., 2022; Ramesh et al., 2022; Rombach et al., 2022; Esser et al., 2024), multimodal document summarization (Jangra et al., 2023) and multimodal question answering (Masry et al., 2022; Yue et al., 2024; Lu et al., 2024).

Scientific figures are much more nuanced than general images in terms of their content, where every minute detail holds a great deal of importance. They are also varied in terms of their content. Figures in scientific documents may range from pipeline representations to graphs to plots to real-world images, and the list goes on and on. Moreover, the diversity in scientific figures is such that each downstream task requires a different class of figures. For example, mixing plots with other types of figures could be problematic for the text-to-figure generation task (Rodriguez et al., 2023b) because plots generally require additional quantitative data not found in captions. For text-to-image synthesis, (Rodriguez et al., 2023b) filter figures

071 by searching for keywords, such as “architecture”,  
072 “model diagram” or “pipeline,” in their captions.

073 A fine-grained classification of scientific figures  
074 will help in targeted finetuning, i.e., not all down-  
075 stream tasks require fine tuning on all figures, and  
076 therefore, we may finetune on only the types of fig-  
077 ures that we require for that particular task. It will  
078 also help in scientific image search. For example,  
079 if a person wants to search for the “auto-encoder ar-  
080 chitectures in Deep Learning”, the search systems  
081 can easily narrow down the results under the archi-  
082 tecture class and then find the matching results.

083 But, in most existing datasets, there is no clas-  
084 sification indicating the type of figure that would  
085 enable researchers to pursue such tasks. There-  
086 fore, to address this gap, we introduce AI-Figures,  
087 a large fine-grained dataset obtained through a  
088 YOLO-based distantly supervised pipeline from  
089 Artificial Intelligence research papers. Our dataset  
090 is curated in such a way that it can be useful for  
091 tasks as diverse as captioning, visual question an-  
092 swering, text-to-figure, and multimodal document  
093 retrieval. Our dataset has 9 different categories for  
094 representing various kinds of scientific figures that  
095 are particularly common in Artificial Intelligence  
096 research papers.

097 We evaluate a wide spectrum of pre-trained foun-  
098 dational models on our proposed dataset for a di-  
099 versity of vision-to-text and text-to-vision tasks.  
100 Our experiments demonstrate the challenging na-  
101 ture as well as the effectiveness of our dataset. The  
102 challenging nature is exhibited by the low results  
103 obtained by state-of-the-art Large Vision Language  
104 Models (LVLMs) on the standard “figure caption-  
105 ing” and the relatively new “text-to-figure” tasks.  
106 We show the effectiveness of our dataset as a train-  
107 ing resource which can improve the scientific liter-  
108 ature understanding capability of LVLMs.

109 In summary, our contributions are the following:  
110 (a) We introduce AI-FIGURES, a large multimodal  
111 dataset containing 133, 749 figures that are classi-  
112 fied into 8 figure types.  
113 (b) We test the limitations of pre-trained LVLMs  
114 in scientific literature understanding tasks like fig-  
115 ure captioning text-to-figure generation, scholarly  
116 multimodal question-answering and multimodal  
117 document retrieval.  
118 (c) We demonstrate that training on targeted classes  
119 of our dataset can lead to performance improve-  
120 ment on three of the above tasks.

## 2 AI-FIGURES 121

122 We present our pipeline for the large-scale ex-  
123 traction and classification of figures and captions,  
124 which leverages the idea that for each document  
125 page, we only need to segment the figure area  
126 and the small chunk of text (caption) that is most  
127 adjacent to it. Our dataset curation pipeline con-  
128 sists of the following steps: Data Collection and  
129 Pre-Processing, Human Annotation, YOLO model  
130 Training, Data Creation using Distant Supervision,  
131 and Dataset Refinement. Figure 2 shows the entire  
132 pipeline used to create AI-FIGURES.

### 2.1 Dataset Schema 133

134 We design our schema keeping in mind the  
135 taxonomy of figures of research papers in the  
136 Computer Science domain, especially the Artificial  
137 Intelligence sub-domain. The 9 figure categories,  
138 along with the caption category, are described as  
139 follows:

140 The **Pipeline** class comprises figures showing  
141 high-level glances at the structural and functional  
142 aspects of the proposed technique or the step by  
143 step workflow showing the ideation of a topic.

144 The **Model Architecture** class shows figures with  
145 a detailed probe into a model structure, like the  
146 CNN or LSTM model architecture.

147 The **Auxiliary Diagrams** category consists of  
148 abstract schematic representations and general  
149 figures related to the paper topic.

150 The **Illustrations** category represents the visual  
151 depiction of an idea or feature that act as an  
152 explanatory visual aids.

153 The **Real Image** category comprises real-world  
154 photographic images, which may be either  
155 instances from a dataset used in the research paper  
156 or any other image in the open-domain.

157 The **Graph Plots** class shows non-performance  
158 and non-statistical plotting like line charts, bar  
159 charts, scatter plots, and histograms.

160 The **Model Performance Plots** class represents  
161 plots of baselines, plots of variations of different  
162 metrics with training.

163 The **Spatial Charts** category contains figures  
164 showing analytical visualizations or statistical  
165 variations like confusion matrices or regression  
166 analysis plots.

167 The **Captions** are also classified into a separate  
168 class, such that it is easier to associate each figure  
169 with its corresponding textual mention.

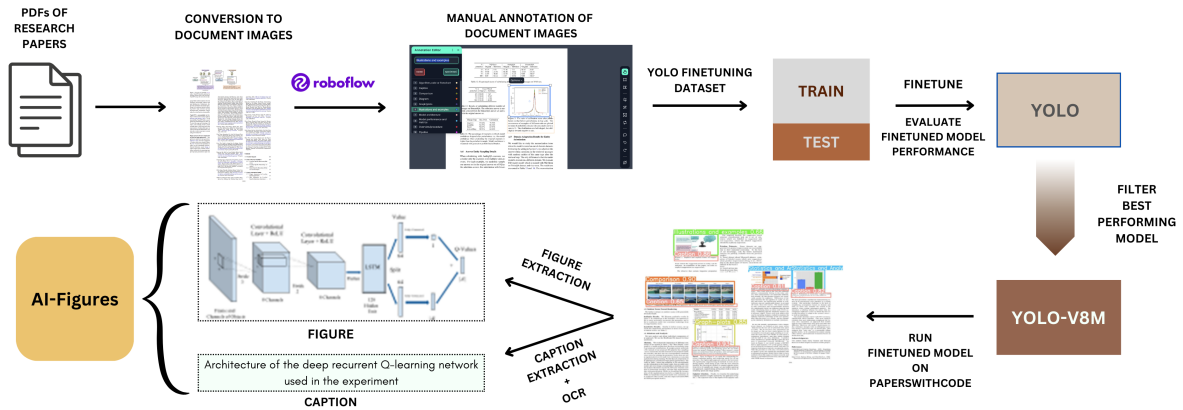


Figure 2: The construction pipeline used to curate the AI-FIGURES dataset.

## 2.2 Manual Annotation (AI-FIGURES-HUMAN)

The first step in our dataset creation pipeline is to annotate a small set of figures and their corresponding captions from a set of documents. This is a human-annotated dataset comprising of a corpus of figures in the Artificial Intelligence/Machine Learning domain paired with their textual contexts, i.e., figure captions. Human annotations were performed on a set of 200 research documents, with 100 each from ACL Anthology<sup>1</sup> and CVPR<sup>2</sup>. We labeled bounding boxes for figure and caption regions for each document page and subsequently assign each figure to one of the pre-defined classes. The Roboflow Annotate<sup>3</sup> platform was used to assist annotators to mark the bounding regions. This platform facilitated dataset pre-processing, division into train, validation, and test sets. AI-FIGURES-HUMAN contains 4975 figures split into training (3790), validation (803) and test (382) sets. While annotating, two more classes (Overview-procedure and Algorithm/Code/Flowchart) were included but were later left out as described in the Section 2.5. The annotation guidelines provided to the human annotators can be found in the Appendix.

## 2.3 YOLO model Training

YOLO (You Only Look Once) (Redmon et al., 2015) is a hugely popular fast object detection and image segmentation model. In YOLO, object detection is reformulated as a regression problem from image pixels to bounding box coordinates and class

<sup>1</sup><https://aclanthology.org/>

<sup>2</sup><https://openaccess.thecvf.com/>

<sup>3</sup><https://roboflow.com/annotate>

probabilities.

Model	mAP	mAP 0.5-0.95	P	R
YOLOv5s	0.506	0.434	0.461	0.607
YOLOv5m	0.497	0.449	<b>0.471</b>	0.558
YOLOv5l	0.51	0.458	0.434	0.644
YOLOv8s	0.49	0.441	0.424	0.62
YOLOv8m	<b>0.515</b>	<b>0.462</b>	0.445	0.667
YOLOv8l	0.505	0.456	0.42	<b>0.695</b>

Table 1: Results of YOLO on AI-FIGURES (Human). P represents Precision while R represents Recall

We train several versions of YOLO models, including YOLOv5s, YOLOv5m, YOLOv5l, YOLOv8s, YOLOv8m, and YOLOv8l on AI-FIGURES-HUMAN. Based on the mean Average Precision (mAP) scores on the test set of AI-FIGURES-HUMAN, as shown in Table 1, we select YOLOv8m for creation of the larger AI-FIGURES corpus.

## 2.4 AI-FIGURES data creation

### 2.4.1 Data Collection and Pre-Processing

We use the URLs present in the *PapersWithCode*<sup>4</sup> repository to curate a corpus of open-access research papers as PDF documents. This open-sourced corpus covers a wide variety of research papers in the AI-ML domain across multiple conferences and journals. Each document page is converted into an image using the *pdf2image*<sup>5</sup> library.

### 2.4.2 Data Creation using Distant Supervision

We metamorphose the figure and caption extraction task into a segmentation task. The figures and

<sup>4</sup><https://paperswithcode.com/>

<sup>5</sup><https://pypi.org/project/pdf2image/>

captions in the page are only objects in the image. This allows us to assign separate classes to each figure. Captions, which are present alongside figures, are naturally treated as objects as well and are classified into a separate class. We use the best YOLO model from Table 1 (YOLOv8m) to extract figures from the document page images that were extracted in the previous step. Subsequently, we run an OCR (Optical character recognition) model<sup>6</sup> over the Captions class to convert them to textual form.

## 2.5 Dataset Refinement Process

After manual pilot assessment of the extracted figures and captions, three issues were revealed with our proposed approach.

Firstly, if the document page image contains two or more figures belonging to the same category, the YOLO model extracts only the last extracted figure, although it detects both the figures. However, the model extracts all the captions in the input page image. This leads to a mismatch in the number of captions and images. To circumvent this problem, we map the selected figure’s bounding box co-ordinate to the bounding box coordinates of the closest caption by calculating the Euclidean distance between the centers of the bounding boxes.

Secondly, if a detected figure crop has been classified into multiple classes at the same time with varying confidence scores, then the YOLO model allocates the figure to both classes. To remove such ambiguity, we first detect multiple class assignments based on the maximal overlap of bounding box co-ordinates. We then assign the class with the greatest confidence score to the figure.

Finally, we remove the Algorithm/Code/Flowchart class from the dataset due to the high occurrence of hallucinations in this category. The frequent hallucinations arise because the model often confuses an Algorithm/Code image with a regular text snippet. We also combine the Overview/Procedure and the Pipeline into a single unified Pipeline class because both these classes contain similar figures. We remove all figures that have captions shorter than 5 words. Also, phrases like Figure x./Figure x./Fig. x./ Fig. x is deleted from the beginning of each caption.

<sup>6</sup><https://github.com/tesseract-ocr/tesseract>

## 2.6 Human Evaluation

We construct the following manual evaluation setup to evaluate our dataset construction process. We divide 3000 figure-caption-category triplets among 8 annotators, who are provided with the extracted figure, the extracted caption, the predicted category, and the URL to the original paper from which they have been extracted.

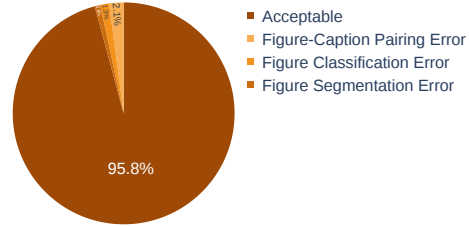


Figure 3: Results for the manual analysis of our distantly supervised dataset, AI-FIGURES

We select 8 graduate students with knowledge in Computer Science as annotators to evaluate our distantly-supervised dataset. We ask the annotators to categorize the dataset samples into the following four categories: (1) **Acceptable**, where the image segmentation is done correctly; (2) **Figure Segmentation Error**, where the figure crop is done incorrectly; (3) **Figure Classification Error**, where the model inaccurately classifies the figure into an unrelated category; (4) **Figure-Caption Pairing Error**, where the figure is paired with an incorrect caption. Figure 3 shows the aggregated results of the manual evaluation of the dataset construction pipeline by human annotators.

## 2.7 Comparison with Baseline Approaches

Model	Precision	Recall	F1	Avg. P (IOU 0.5)
PDFFigures	0.395	0.634	0.487	0.333
YOLOv8m	<b>0.445</b>	<b>0.667</b>	<b>0.534</b>	<b>0.515</b>

Table 2: Comparison of our approach with PDFFigures

**PDFFigures 2.0** (Clark and Divvala, 2015, 2016): It consists of three phases: Caption Detection using keyword search, Region Identification using paragraph grouping with classification, and Figure Assignment using a scoring function to rate the proposed regions. We test the approach of PDFFigures 2.0 with our construction pipeline on the test set of our human-annotated dataset, AI-FIGURES-HUMAN. The results are present in Table 2, where we see that our method comprehensively

Dataset	Source	Annotation	Papers	Figures	Classes
CS-150 (Clark and Divvala, 2015)	CS-conferences	Manual	150	458	✗
CS-Large (Clark and Divvala, 2016)	Semantic Scholar	Manual	346	952	✗
DocFigures (Jobin et al., 2019)	Computer vision conferences	SVM	20K	33K	28
Paper2Fig100k (Rodriguez et al., 2023b)	ArXiv	GROBID	183,427	102,453	✗
ArXivCap (Li et al., 2024b)	ArXiv	ImageMagick	572K	6.4M	✗
ACL-FIG-PILOT (Karishma et al., 2023)	ACL Anthology	Manual	-	1,671	19
ACL-FIG (Karishma et al., 2023)	ACL Anthology	-	55,760	112,052	✗
VisImages (Deng et al., 2022)	IEEE InfoVis and VAST	Manual	1,397	12,267	34
EXAMS-V (Das et al., 2024)	School exam questions	Manual	-	20,932	20
MMMU (Yue et al., 2024)	Exams, quizzes and textbooks	Manual	-	12,507	30
Science QA (Lu et al., 2022)	Elementary and high school science curricula	Manual	-	10,332	127
SCI-3000 (Darmanović et al., 2023)	DOAJ (Comp. science, medicine, physics, chemistry, and technology)	None	3000	✗	
MathVista (Lu et al., 2024)	Existing VQA, Math QA datasets	Repurposed	-	5487	5
SciBench (Wang et al., 2024b)	Textbooks: Phys, Chem, Maths	Manual	-	177	✗
VGbench (Zou et al., 2024)	Github, Kaggle	Manual	-	5845	3
CharXiv (Wang et al., 2024c)	ArXiv	SigLIP	-	2323	8
Datikz (Belouadi et al., 2024)	TEX Stack Exchange, ArXiv papers, artificial examples	Manual	-	119 789	✗
SciFIBench (Roberts et al., 2024)	SciCap	Manual	6000	1,96,000	2
MMSCI (Li et al., 2025)	Nature Articles	Directly from website	131,393	742,273	7
ChartMimic (Shi et al., 2024)	ArXiv	Manual	-	600	22
AI-FIGURES-HUMAN	PaperswithCode	Manual	200	4,844	10
AI-FIGURES	PaperswithCode	YOLO	26,969	133,749	9

Table 3: Comparison with prior scientific figure datasets.

outperforms the PDFFigures approach on all metrics. Upon qualitative evaluation, we find that there are two major reasons for the performance of PDFFigures: firstly there are a lot of tables extracted along with the figures, and secondly, this algorithm randomly extracts many blank strips.

**PaperMage** (Lo et al., 2023): It is an open-source Python toolkit that allows the representation and manipulation of both textual and visual elements in a document. For its evaluation, we used a test set of 5532 PDF page images, out of which 4325 pages contained figures. In 55 out of 4325 pages, PaperMage showed some signs of figure recognition. In the remaining 4270 pages, no figures were detected, indicating false negatives across these pages. In 27% of the 55 pages, PaperMage exhibits poor extraction quality, with the bounding box placed in the middle of the figure, failing to properly define the figure’s boundaries. As a result, most of these extractions are sliced and unsuitable for evaluation. However, in 41 images, the extractions were suitable for evaluation, where we observed an average Intersection over Union (IoU) score of 0.6818.

## 2.8 Dataset Statistics

Our final dataset contains 1,33,749 scientific figure-caption pairs. We present the class-wise statistics of both the human-annotated dataset and the larger inferred dataset in Table 4. Higher num-

Model	AI-FIGURES-Human	AI-FIGURES
Pipeline	519	2,154
Model Arch.	500	12,169
Auxiliary Diagram	402	12,975
Illustrations	1,351	39,359
Real Image	296	1,910
Graph Plots	956	52,932
Model Perf. Plots	324	4,305
Spatial Charts	313	7,945
Algo./Flowchart	183	-
<b>Total</b>	<b>4,844</b>	<b>133,749</b>

Table 4: Class-wise statistics of our proposed dataset

ber of figures in a category reflect higher presence of that class of figures in research documents. Our dataset contains a total of 4,925,626 words with the average caption length being 36.83 words and the quartile lengths being [13, 27, 49].

## 3 Comparison with Related Datasets

Table 3 contains a list of related datasets and shows them in comparison with our dataset. CS-150 (Clark and Divvala, 2015) and CS-Large datasets (Clark and Divvala, 2016) are very small datasets and do not include fine-grained classification. The is also true for Paper2Fig100k (Rodriguez et al., 2023b), ArXivCap (Li et al., 2024b) and Multimodal ArXiv (Li et al., 2024b), which are all sourced from ArXiv. ACL-FIG (Karishma et al., 2023) contains 112,05 unlabeled figures from 55,760 papers in the ACL Anthology. It is ac-

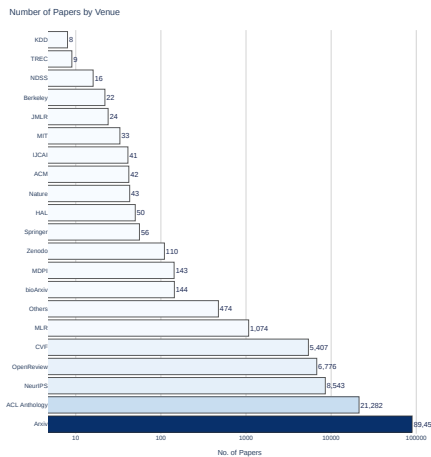


Figure 4: Domain Distribution of the figures in our dataset

348 accompanied by its labeled subset ACL-FIG-PILOT, 349 containing 1671 figures with 19 manually verified 350 labels. However, ACL-FIG figures do not con- 351 tain captions, and this limits their utility. For the 352 MMSci dataset (Li et al., 2025), the data collec- 353 tion pipeline is through web-scraping, which is not 354 adaptable to all domains and data sources. Ad- 355 ditionally, the annotation schema introduced in 356 MMSci is skewed towards the natural sciences 357 domain since it contains classes such as “Micro- 358 scopic/Macroscopic photographs” or “Geographi- 359 cal and Environmental Maps”, whereas our focus 360 is CS-related areas such as AI. DocFigures (Jobin 361 et al., 2019) contains figures from computer vi- 362 sion conferences like CVPR, ECCV, ICCV, etc., 363 extracted using a web scraper and classified with 364 Support Vector Machine (SVM). However, the fo- 365 cus is mainly on plots, the domain diversity is 366 low, and the reliability of the dataset is also not 367 high due to minimal and simplistic human anno- 368 tation. For the curation of SciFiBench (Roberts et al., 369 2024), scientific figure-caption pairs are sourced 370 from the SciCap dataset, which contains 7 cate- 371 gories. Therefore, the annotation schema here is 372 skewed towards graphical plots and flowcharts. In 373 contrast, our annotation schema provides more cov- 374 erage towards all types of images in the AI domain. 375 The same argument holds for the CharXiv (Wang 376 et al., 2024c), ChartMimic (Shi et al., 2024), and 377 VisImages (Deng et al., 2022) datasets, which are 378 constrained to the categories of graph plots and 379 charts. MedICaT (Subramanian et al., 2020) con- 380 tains 217K medical images while SCI-3000 (Dar- 381 manović et al., 2023) includes figures extracted 382 from PDFs from open-access journals across com-

puter science, medicine, physics, chemistry, and 383 technology without explicit classification. The 384 MMMU (Yue et al., 2024) and Science QA (Lu 385 et al., 2022) datasets figures are derived from text- 386 books and exam materials targeting elementary and 387 high school science curricula. They differ signifi- 388 cantly from our domain of interest and fail to cap- 389 ture the complexity of figures in scholarly papers. 390 VGBench (Zou et al., 2024) and DaTikz (Belouadi 391 et al., 2024) are specialized benchmarks for vector 392 graphics and are thus constrained to figures ren- 393 dered using TikZ, Graphviz, or SVG codes. 394

## 4 Downstream tasks 395

### 4.1 Figure Captioning 396

397 This objective of this task is to generate visual de- 398 scriptions from scientific figures. We benchmark 399 the following LVLMS on this task: MOLMO-7B 400 (Deitke et al., 2024), InternVL2\_5-8B (Chen et al., 401 2024), Qwen2-VL-7B (Wang et al., 2024a), Janus- 402 Pro-7B (Chen et al., 2025), MiniCPM-V (Yao et al., 403 2024), and advanced proprietary models such as 404 GPT-4o, GPT-4o mini, Claude 3.7 Sonnet and Gem- 405 ini 2.0 Flash. We tested the proprietary models on 406 a 500-sample subset due to budget constraints. For 407 each model, we report the BLEU-2 (Papineni et al., 408 2002), ROUGE-L (Lin, 2004), and the BERT-Score 409 (Zhang et al., 2019) in all three settings, i.e., cap- 410 tioning without context, captioning with title as 411 context, and captioning with both title and abstract 412 as context. We fine-tune the GIT-base and GIT- 413 large (Wang et al., 2022) models on the uncleaned 414 version of our dataset so that we may train on as 415 many figure caption pairs as possible, but while 416 testing, we use the cleaned set. Tables 5 and 6 417 show the results of the various models on the figure 418 captioning task. We see that models finetuned on 419 our dataset perform the best followed by the open- 420 source models (especially MiniCPM and Qwen) 421 and the proprietary models like Gemini or GPT.

422 **Manual Analysis:** We also construct a manual 423 evaluation setup for the Figure Captioning task. We 424 construct two different sets with 25 captions each 425 and ask three annotators to analyze the quality of 426 the generated captions for each model. Each anno- 427 tator is provided with the figure, the gold standard 428 caption from our dataset, and the generated cap- 429 tions from the models. The annotators are doctoral 430 students who work in allied areas and have at least 431 one publication in the domain of AI/ML/CV/NLP. 432 The annotation guidelines for this task are provided

Model	Zero-shot Captioning			Context = Title			Context = Title + Abstract		
	BLEU-2	R-L	B-S	BLEU-2	R-L	B-S	BLEU-2	R-L	B-S
MOLMO-7B	1.42	8.13	81.41	1.27	7.99	81.49	1.35	7.91	81.61
InternVL2_5-8B	1.31	7.83	81.01	1.40	7.96	81.18	1.15	7.09	80.83
Qwen2-VL-7B	1.91	9.00	81.40	<b>2.35</b>	<b>9.62</b>	81.92	<b>2.28</b>	<b>9.50</b>	<b>81.75</b>
MiniCPM-V	<b>1.94</b>	<b>9.54</b>	<b>81.66</b>	1.47	8.53	<b>82.38</b>	1.64	7.56	81.07
Janus-Pro-7B	1.60	8.59	81.10	1.62	8.73	81.22	1.79	8.86	81.42

Table 5: Evaluation results of the Figure Captioning task by prompting over open-source models.

Model	BLEU-2	ROUGE-L	BERTScore
<b>Proprietary Models (500 samples)</b>			
GPT-4o	1.28	8.08	81.55
GPT-4o mini	<b>1.44</b>	8.12	81.51
Claude 3.7 Sonnet	1.13	7.13	81.06
Gemini 2.0 Flash	1.13	<b>8.47</b>	<b>81.69</b>
<b>Finetuned Models</b>			
GIT-base	1.58	<b>10.74</b>	<b>83.22</b>
GIT-large	<b>3.01</b>	10.01	81.61

Table 6: Results for Figure Captioning using Proprietary Models as well as Finetuned models

in the Appendix. In line with the quantitative results, we see in Table 7 that there are only a few acceptable responses across all models, with the Qwen model performing the best among the given models, whereas MOLMO performs the worst.

Model	Accept.	Simp.	Misrep.	Recog.
MOLMO-7B	1	3	37	43
InternVL2.5-8B	18	0	20	36
Qwen2-VL-7B	<b>28</b>	1	31	13
MiniCPM-V	5	8	39	29
Janus-Pro-7B	8	3	35	35

Table 7: Manual Analysis results for Figure Captioning

## 4.2 Text-to-Figure

The text-to-image (T2I) generation is quite successful in the open-domain, but the same cannot be said for the scientific domain (Rodriguez et al., 2023a; Zala et al., 2023). The text-to-figure task involves generating a figure from its corresponding textual mention in the paper.

To test the effectiveness of our fine-grained classification, we train a Stable Diffusion model (SDXL) both on our whole dataset and on a smaller subset of classes for which the text-to-figure generation task is more plausible, i.e., classes containing procedural or architectural figures. We select only such figures because it is inexplicable to generate plot-based figures or other auxiliary diagrams or illustrations from only the figure captions. We select the *Model Architecture* and *Pipeline* classes to create a training set comprising 21, 839 images and the test set with 5, 461 images for the targeted

finetuning. We fine-tune the Stable Diffusion-XL model using LoRA (Hu et al., 2022) on both training sets (full and targeted) for 5 epochs with a batch size of 12 and a learning rate of 1e-06. We compute the Fréchet Inception Distance (FID) (Heusel et al., 2017), Inception Score (Salimans et al., 2016) and the Kernel Inception Distance (Bińkowski et al., 2021) metrics for this task.

Train Corpus	FID ↓	IS ↑	KID ↓
AI-FIGURES (full)	118.15	6.79 ± 0.58	0.086 ± 0.002
AI-FIGURES (targeted)	<b>103.153</b>	<b>7.26 ± 0.48</b>	<b>0.075 ± 0.002</b>

Table 8: Text-to-Figure using SDXL\_base\_1.0

Table 8 shows the results for the text-to-figure task, which clearly show the effectiveness of targeted finetuning on task-specific classes. The results are also better than those obtained by (Rodriguez et al., 2023a). However, manual review shows that the generated images are hardly comprehensible. A set of images received as output from both the finetuned models are presented on our website.

## 4.3 Scholarly Question Answering

MathVista (Lu et al., 2024) is a mathematical reasoning benchmark within visual contexts. We show the performance of five subsets of MCQ questions and three subsets of Free form questions of the MathVista dataset, which are present in its testmini version. We choose the subsets such that they align with our problem setup, i.e., they are dependent on either academic papers, charts, or scientific knowledge and require numerical rationale.

We choose a certain subset of categories from our dataset and then fine-tune an LVLm on this subset. We posit that the finetuned model will work better than the naive model. We use the "Graph Plots" and the "Spatial Charts" classes to form a subset of our entire dataset since the benchmark is mainly based on reasoning over plots and charts, i.e., we choose a combined 48, 716 figures from our

Model	All	Modality		Reasoning Type			
		Table	Figure	Comp.	Data Extraction	Loc.	Visual Understanding
GPT-4o	<b>0.5</b>	<b>0.52</b>	<b>0.454</b>	<b>0.443</b>	<b>0.565</b>	<b>0.57</b>	<b>0.418</b>
Qwen2-VL-7B-Instruct	0.1334	0.112	0.1863	0.1233	0.1135	0.1667	0.1708
Qwen2-VL-7B-Instruct(fine-tuned)	<u>0.2021</u>	<u>0.1848</u>	<u>0.2446</u>	<u>0.2187</u>	<u>0.1727</u>	<u>0.2272</u>	<u>0.1966</u>

Table 9: Mean reciprocal rank (MRR) results on M3SciQA. GPT-4o results are from (Li et al., 2024a). The second-best results are underlined.

training set to create this subset. We then prompt the InternVL2\_5-8B (Chen et al., 2024) model to generate 3 unique sets of question-answer pairs for each figure. We then finetune Qwen2-VL-7B (Wang et al., 2024a) on this derived question answering set. We use QLoRA (Detmers et al., 2023) in the HuggingFace Ecosystem (TRL) to train the model for 10 epochs with a train batch size of 4, learning rate of  $2e-04$ , and maximum sequence length of 1024.

Data	Qwen2-VL-7B	QWEN2-VL-7B (Finetuned)
<b>Multiple Choice Questions (MCQs)</b>		
PaperQA	53.85	<b>61.54</b>
ScienceQA	<b>55.56</b>	<b>55.56</b>
IQTest	25	<b>50</b>
TabMWP	44.44	<b>55.56</b>
ChartQA	66.67	<b>83.33</b>
<b>Free-form questions</b>		
CLEVR-Math	67.74	<b>72.58</b>
DVQA	82.26	<b>85.48</b>
TabMWP	26.42	<b>43.4</b>

Table 10: Accuracy scores in the MathVista dataset

Table 10 shows the results on the five MCQ subsets and three free form subsets of MathVista. We observe that category-specific fine-tuning helps in achieving better results over all the subsets.

#### 4.4 Multimodal Document Retrieval

This task necessitates both multimodal and multi-document reasoning over scientific papers. M3SciQA (Li et al., 2024a) is a benchmark with expert-annotated questions from paper clusters. The questions are divided into four reasoning categories: comparisons, data extraction, locations, and visual understanding. Therefore, given a locality-specific question  $Q$ , the corresponding image  $I$  and the list of documents  $D = \{d_1, d_2, \dots, d_n\}$ , the task is to determine the ranking  $R = \{r_1, r_2, \dots, r_n\}$  of papers based on the relevance of  $D$  to  $Q$  and  $I$ .

$$R = \mathcal{M}(Q, I, D) \quad (1)$$

We only consider the locality-specific document retrieval setup here, since it tests the capability of

LVLMS. We use the same finetuned Qwen model as the one used in Section 4.3 as well as the naive Qwen 7B model for this task. Table 9 shows the results for this task. The fine-tuned Qwen model outperforms the naive Qwen model across all modalities and reasoning types although the proprietary GPT model gives the best results.

## 5 Related work

**Figure extraction:** Software tools which are ideal for the off-the-shelf-processing of scientific documents include GROBID (GRO, 2008–2024), ParsCit (Isaac Councill and Kan, 2008) and CER-MINE (Tkaczyk et al., 2015). They use various Machine Learning algorithms like CRFs (Conditional Random Fields), recurrent neural networks, and even recent deep learning models. PDFFigures (Clark and Divvala, 2015), a widely used figure extractor, performs structural analysis of individual pages of a document and can identify figures, tables, and captions in the pages. PDFFigCapX (Li et al., 2018) is an algorithm which identifies regions devoid of body text that are adjacent to captions.

**Related datasets:** Datasets of figure-caption pairs in the domain of scientific literature typically focus only on scientific plots. Example datasets include FigureQA (Kahou et al., 2017), DVQA-cap (Kafle et al., 2018), FigJAM (Qian et al., 2021), SciCap (Hsu et al., 2021), DocFigure (Jobin et al., 2019), Paper2Fig100k (Rodriguez et al., 2023b), and ACL-FIG (Karishma et al., 2023).

## 6 Conclusion

We introduce the AI-FIGURES dataset, and thereby, propose a dataset construction pipeline that can be used to extract and label figure-caption pairs. Our dataset is divided into 8 fine-grained categories. We show the challenging nature of captioning, text-to-figure, scholarly QA and multimodal document retrieval. Results show the improvements achieved on targeted fine-tuning of generative models over specific classes of our dataset.

## 561 Limitations

562 We hereby state the limitations of our work. We  
563 understand that the scientific domain is extremely  
564 challenging and large, and Artificial Intelligence,  
565 i.e., the area we choose for the creating the dataset  
566 here is a very niche and evolving area. So the  
567 dataset may need to be regularly updated for high  
568 practical utility to researchers.

569 Furthermore, the space of language-vision mod-  
570 els and language models is rapidly evolving and  
571 therefore, we have not been able to exhaustively  
572 test on many of these models.

## 573 References

574 2008–2024. Grobid. [https://github.com/](https://github.com/kermitt2/grobid)  
575 [kermitt2/grobid](https://github.com/kermitt2/grobid).

576 Jonas Belouadi, Anne Lauscher, and Steffen Eger. 2024.  
577 [AutomaTikZ: Text-guided synthesis of scientific vec-](#)  
578 [tor graphics with TikZ](#). In *The Twelfth International*  
579 *Conference on Learning Representations*.

580 Mikołaj Bińkowski, Danica J. Sutherland, Michael Ar-  
581 bel, and Arthur Gretton. 2021. [Demystifying mmd](#)  
582 [gans](#). *Preprint*, arXiv:1801.01401.

583 Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan,  
584 Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan.  
585 2025. [Janus-pro: Unified multimodal understanding](#)  
586 [and generation with data and model scaling](#). *arXiv*  
587 *preprint arXiv:2501.17811*.

588 Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo  
589 Chen, Sen Xing, Muyan Zhong, Qinglong Zhang,  
590 Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu,  
591 Yu Qiao, and Jifeng Dai. 2024. [Internvl: Scaling up](#)  
592 [vision foundation models and aligning for generic](#)  
593 [visual-linguistic tasks](#). *Preprint*, arXiv:2312.14238.

594 Christopher Clark and Santosh Divvala. 2016. [Pdf-](#)  
595 [figures 2.0: Mining figures from research papers](#).  
596 In *Proceedings of the 16th ACM/IEEE-CS on Joint*  
597 *Conference on Digital Libraries, JCDL '16*, page  
598 143–152, New York, NY, USA. Association for Com-  
599 puting Machinery.

600 Christopher Clark and Santosh Kumar Divvala. 2015.  
601 [Looking beyond text: Extracting figures, tables and](#)  
602 [captions from computer science papers](#). In *AAAI*  
603 *Workshop: Scholarly Big Data*.

604 Filip Darmanović, Allan Hanbury, and Markus  
605 Zlabinger. 2023. [SCI-3000: A dataset for figure,](#)  
606 [table and caption extraction from scientific PDFs](#). In  
607 *International Conference on Document Analysis and*  
608 *Recognition*, pages 234–251. Springer.

609 Rocktim Das, Simeon Hristov, Haonan Li, Dimitar  
610 Dimitrov, Ivan Koychev, and Preslav Nakov. 2024.

[EXAMS-V: A multi-discipline multilingual multi-](#)  
611 [modal exam benchmark for evaluating vision lan-](#)  
612 [guage models](#). In *Proceedings of the 62nd Annual*  
613 *Meeting of the Association for Computational Lin-*  
614 *guistics (Volume 1: Long Papers)*, pages 7768–7791,  
615 Bangkok, Thailand. Association for Computational  
616 Linguistics. 617

Matt Deitke, Christopher Clark, Sangho Lee, Rohun  
618 Tripathi, Yue Yang, Jae Sung Park, Mohammadreza  
619 Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini,  
620 et al. 2024. [Molmo and pixmo: Open weights and](#)  
621 [open data for state-of-the-art vision-language models](#).  
622 *arXiv preprint arXiv:2409.17146v2*. Allen Institute  
623 for AI and University of Washington. 624

Dazhen Deng, Yihong Wu, Xinhuan Shu, Jiang Wu, Si-  
625 wei Fu, Weiwei Cui, and Yingcai Wu. 2022. [Visim-](#)  
626 [ages: A fine-grained expert-annotated visualization](#)  
627 [dataset](#). *IEEE Transactions on Visualization and*  
628 *Computer Graphics*, 29(7):3298–3311. 629

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and  
630 Luke Zettlemoyer. 2023. [Qlora: Efficient finetuning](#)  
631 [of quantized llms](#). In *Advances in Neural Information*  
632 *Processing Systems*, volume 36, pages 10088–10115.  
633 Curran Associates, Inc. 634

Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim  
635 Entezari, Jonas Müller, Harry Saini, Yam Levi, Do-  
636 minik Lorenz, Axel Sauer, Frederic Boesel, Dustin  
637 Podell, Tim Dockhorn, Zion English, and Robin  
638 Rombach. 2024. [Scaling rectified flow trans-](#)  
639 [formers for high-resolution image synthesis](#). In *Proceed-*  
640 *ings of the 41st International Conference on Machine*  
641 *Learning*, volume 235 of *Proceedings of Machine*  
642 *Learning Research*, pages 12606–12633. PMLR. 643

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,  
644 Abhinav Pandey, Abhishek Kadian, Ahmad Al-  
645 Dahle, Aiesha Letman, and Akhil Mathur et al.  
646 2024. [The llama 3 herd of models](#). *Preprint*,  
647 arXiv:2407.21783. 648

Martin Heusel, Hubert Ramsauer, Thomas Unterthiner,  
649 Bernhard Nessler, and Sepp Hochreiter. 2017. [Gans](#)  
650 [trained by a two time-scale update rule converge to](#)  
651 [a local nash equilibrium](#). In *Advances in Neural*  
652 *Information Processing Systems*, volume 30. Curran  
653 Associates, Inc. 654

Ting-Yao Hsu, C Lee Giles, and Ting-Hao Huang. 2021.  
655 [SciCap: Generating captions for scientific figures](#).  
656 In *Findings of the Association for Computational*  
657 *Linguistics: EMNLP 2021*, pages 3258–3264, Punta  
658 Cana, Dominican Republic. Association for Compu-  
659 tational Linguistics. 660

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan  
661 Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and  
662 Weizhu Chen. 2022. [LoRA: Low-rank adaptation of](#)  
663 [large language models](#). In *International Conference*  
664 *on Learning Representations*. 665



780	2021, WWW '21, page 2792–2804, New York, NY, USA. Association for Computing Machinery.		
781			
782	Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. 2022. Hierarchical text-conditional image generation with clip latents. <i>arXiv preprint arXiv:2204.06125</i> , 1(2):3.		
783			
784			
785			
786	Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. In <i>International conference on machine learning</i> , pages 8821–8831. Pmlr.		
787			
788			
789			
790			
791	Joseph Redmon, Santosh Kumar Divvala, Ross B. Girshick, and Ali Farhadi. 2015. <b>You only look once: Unified, real-time object detection</b> . <i>CoRR</i> , abs/1506.02640.		
792			
793			
794			
795	Jonathan Roberts, Kai Han, Neil Houlsby, and Samuel Albanie. 2024. Scifibench: Benchmarking large multimodal models for scientific figure interpretation. <i>arXiv preprint arXiv:2405.08807</i> .		
796			
797			
798			
799	Juan A Rodriguez, David Vazquez, Issam Laradji, Marco Pedersoli, and Pau Rodriguez. 2023a. FigGen: Text to scientific figure generation. <i>arXiv preprint arXiv:2306.00800</i> .		
800			
801			
802			
803	Juan A Rodriguez, David Vazquez, Issam Laradji, Marco Pedersoli, and Pau Rodriguez. 2023b. OCR-VQGAN: Taming text-within-image generation. In <i>Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision</i> , pages 3689–3698.		
804			
805			
806			
807			
808	Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pages 10684–10695.		
809			
810			
811			
812			
813			
814	Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. 2022. Photorealistic text-to-image diffusion models with deep language understanding. <i>Advances in neural information processing systems</i> , 35:36479–36494.		
815			
816			
817			
818			
819			
820			
821	Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen, and Xi Chen. 2016. <b>Improved techniques for training gans</b> . In <i>Advances in Neural Information Processing Systems</i> , volume 29. Curran Associates, Inc.		
822			
823			
824			
825			
826	Chufan Shi, Cheng Yang, Yaxin Liu, Bo Shui, Junjie Wang, Mohan Jing, Linran Xu, Xinyu Zhu, Siheng Li, Yuxiang Zhang, Gongye Liu, Xiaomei Nie, Deng Cai, and Yujiu Yang. 2024. Chartmimic: Evaluating Imm’s cross-modal reasoning capability via chart-to-code generation. <i>arXiv preprint arXiv:2406.09961</i> .		
827			
828			
829			
830			
831			
832	Sanjay Subramanian, Lucy Lu Wang, Ben Bogin, Sachin Mehta, Madeleine Van Zuylen, Sravanthi Parasa, Sameer Singh, Matt Gardner, and Hannaneh		
833			
834			
		Hajishirzi. 2020. MedICAT: A dataset of medical images, captions, and textual references. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2112–2120.	835
			836
			837
			838
		Qwen Team. 2024. <b>Qwen2.5: A party of foundation models</b> .	839
			840
		D. Tkaczyk, P. Szostek, M. Fedoryszak, P.J. Dendek, and Ł. Bolikowski. 2015. <b>Cerminet: automatic extraction of structured metadata from scientific literature</b> . In <i>International Journal on Document Analysis and Recognition (IJDAR)</i> , pages 1433–2825.	841
			842
			843
			844
			845
		Jianfeng Wang, Zhengyuan Yang, Xiaowei Hu, Linjie Li, Kevin Lin, Zhe Gan, Zicheng Liu, Ce Liu, and Lijuan Wang. 2022. <b>Git: A generative image-to-text transformer for vision and language</b> . <i>Preprint</i> , arXiv:2205.14100.	846
			847
			848
			849
			850
		Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024a. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. <i>arXiv preprint arXiv:2409.12191</i> .	851
			852
			853
			854
			855
			856
			857
			858
		Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R. Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. 2024b. Scibench: evaluating college-level scientific problem-solving abilities of large language models. In <i>Proceedings of the 41st International Conference on Machine Learning, ICML’24</i> . JMLR.org.	859
			860
			861
			862
			863
			864
			865
		Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Hao-tian Liu, Sadhika Malladi, Alexis Chevalier, Sanjeev Arora, and Danqi Chen. 2024c. Chartxiv: Charting gaps in realistic chart understanding in multimodal llms. <i>arXiv preprint arXiv:2406.18521</i> .	866
			867
			868
			869
			870
			871
		Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. 2018. AttnGAN: Fine-grained text to image generation with attentional generative adversarial networks. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> , pages 1316–1324.	872
			873
			874
			875
			876
			877
		Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin Zhao, Zhihui He, Qianyu Chen, Huarong Zhou, Zhensheng Zou, Haoye Zhang, Shengding Hu, Zhi Zheng, Jie Zhou, Jie Cai, Xu Han, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. 2024. <b>Minicpm-v: A gpt-4v level mllm on your phone</b> . <i>Preprint</i> , arXiv:2408.01800.	878
			879
			880
			881
			882
			883
			884
			885
		Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, Huan Sun, Yu Su, and Wenhao Chen. 2024. Mmmu:	886
			887
			888
			889
			890
			891

892	A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In <i>Proceedings of CVPR</i> .	– Carefully examine the content and concept behind each figure.	936
893		– Assign the most appropriate class from the predefined categories listed below.	937
894		– If a figure could belong to multiple categories, choose the most dominant or relevant one.	938
895	Abhay Zala, Han Lin, Jaemin Cho, and Mohit Bansal. 2023. DiagrammerGPT: Generating open-domain, open-platform diagrams via llm planning. <i>arXiv preprint arXiv:2310.12128</i> .		939
896			940
897			941
898			942
899	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2019. <a href="#">Bertscore: Evaluating text generation with BERT</a> . <i>CoRR</i> , abs/1904.09675.	• Subfigures and Complex Figures	943
900		– If a figure consists of multiple subfigures labeled as (a), (b), (c), etc., annotate the entire figure as one bounding box.	944
901		– If subfigures have separate captions, annotate them individually with their respective captions.	945
902			946
903	Bocheng Zou, Mu Cai, Jianrui Zhang, and Yong Jae Lee. 2024. <a href="#">VGBench: Evaluating large language models on vector graphics understanding and generation</a> . In <i>Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing</i> , pages 3647–3659, Miami, Florida, USA. Association for Computational Linguistics.	• For the Algorithms Code or Flowchart class, ensure that only the code or flowchart is included in the bounding box, excluding body text explanations.	947
904			948
905			949
906			950
907			951
908			952
909			953
910	<b>A Dataset Statistics</b>		
911	Our dataset contains the following fields:		
912	• Figure Filename	• Do not include surrounding text or unrelated parts of the page in the bounding box.	954
913	• Figure Caption	• Do not annotate tables or equations; this task is only for figures and captions.	955
914	• Paper Abstract	• Do not annotate tables or equations; this task is only for figures and captions.	956
915	• Paper Title	• There should be no overlapping or duplicate annotations.	957
916	• PDF URL	• Class Definitions: Each figure must be assigned exactly one of the following classes:	958
917		– <b>Caption:</b> Text that describes a figure.	959
918	<b>B Dataset Construction - Annotation Guidelines</b>	Example: "Figure 3: Architecture of the proposed model."	960
919	• Each figure and its corresponding caption must have a separate bounding box.	– <b>Diagrams:</b> Schematic representations, flowcharts, or conceptual illustrations.	961
920		Examples: System design diagrams, logic flow representations.	962
921	– Figures should be assigned to exactly 1 of the 10 predefined classes .	– <b>Graphs/Plots:</b> Graphs, charts and mathematical plots.	963
922	– Captions are always assigned the class "Caption".	Examples: Line graphs, bar charts, scatter plots, histograms.	964
923	– Ensure no overlap between figure and caption annotations.	– <b>Illustrations and Examples:</b> Figures providing explanatory visual aids for a concept or process.	965
924		Examples: Illustrative sketches, educational examples, artistic depictions.	966
925		– <b>Model Architecture:</b> Figures depicting the structural design of machine learning or deep learning models.	967
926		Examples: Transformer model, LSTM or YOLO architectural diagram	968
927	• Bounding Box Rules		969
928	– Draw tight bounding boxes around each figure and its caption.		970
929	– The caption box should cover only the text of the caption, not surrounding text.		971
930	– The figure box should include only the visual content of the figure, avoiding page borders or surrounding text.		972
931			973
932			974
933			975
934			976
935	• Classifying Figures		977

Model	AI-FIGURES-HUMAN	AI-FIGURES-Before Clean	AI-FIGURES
Algo./Flowchart	183	10,014	-
Diagram	402	14,704	12,975
Graph Plots	956	55,676	52,932
Illustrations	1,351	42,066	39,359
Model Arch.	500	15,839	12,169
Metrics	324	4,548	4,305
Overview	340	2,201	2,095
Pipeline	179	59	59
Real Image	296	2,321	1,910
Stat./Analysis	313	8,756	7,945
<b>Total</b>	<b>4,844</b>	<b>1,56,184</b>	

Table 11: AI-FIGURES Dataset Across Categories and Cleaning Stages

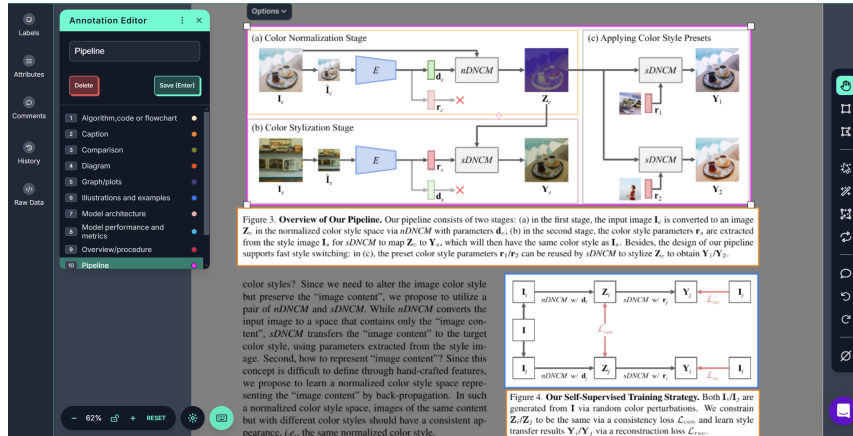


Figure 5: Roboflow Annotate Platform

- **Statistics and Analysis:** Figures containing statistical results, experimental comparisons, or analytical visualizations. Examples: Performance comparison graphs, confusion matrices, regression analysis plots.
- **Overview/Procedure:** Figures illustrating a high level overview of multi-step processes, workflows, or methodologies without detailed representations of each component. Examples: Overview of the object detection process using YOLO
- **Pipeline:** Figures representing entire processing workflows, often spanning multiple steps and modules. Examples: End-to-end ML pipeline diagrams
- **Model Performance and Metrics:** Figures showing model evaluation results, benchmarking, and performance graphs. Examples: Precision-recall curves, accuracy vs. epochs graphs, performance tables.
- **Real Images:** Photographic images or re-

- alistic visual content extracted from real-world sources. Examples: Images from datasets, captured photographs, images of people, animals, places or objects.
- **Algorithms/Code/Flowchart:** Figures containing algorithmic representations, such as code snippets, pseudo code, or flowcharts. Examples: Code blocks (e.g., Python, C++, pseudo code), Flowcharts detailing algorithmic steps, Structured representations of an algorithm's execution flow.

## C Dataset Construction - Manual Evaluation

### C.1 Annotation Guidelines

For each figure, its corresponding caption, class category and a link to the original pdf from which the figure was extracted is provided. The evaluator must :

- Choose "Yes" under "Acceptable" if the figure is segmented and classified correct and paired

with the correct caption as per the parent paper pdf supplied.

- If "No" is selected, then the issue must be narrowed down to one of the following cases:
  - Choose "Figure Segmentation Error" if the figure is cropped in a wrong fashion.
  - Choose "Figure Classification Error" if the figure is classified into the wrong category.
  - Choose "Figure-Caption Pairing Error" if the figure is paired with the wrong caption.

## D Captioning - Manual Evaluation

### D.1 Annotation Guidelines

For each model, evaluate whether the caption generated provides a comprehensive description of its figure. An exact match is not expected with the ground truth caption, but there must be some degree of alignment in the content.

- If the caption generated by a particular model is acceptable, select "Yes".
- If you have selected "No" then narrow down the issue to one of the following:
  - Oversimplification: The oversimplified caption is too short compared with the original ground truth caption.
  - Contextual misunderstanding : Contextual Misinterpretation refers to captions with unmentioned content in the figure.
  - Recognition Error : Recognition Error denotes the model wrongly identified the number or text in the figure.

## E Model Description

### E.1 Models for Figure Captioning

**Molmo-7B:** Molmo-7B-D-0924 (Deitke et al., 2024) is a multimodal AI model developed by the Allen Institute for AI, designed to integrate vision and language understanding. Built upon the Qwen2-7B architecture and utilizing OpenAI’s CLIP as its vision backbone, this model has been trained on PixMo, a curated dataset of 1 million image-text pairs. Molmo-7B-D-0924 achieves an average score of 77.3% across 11 academic benchmarks and holds a human preference Elo rating of 1056, positioning its performance between GPT-4V and GPT-4o. The model is fully open-source,

with all associated artifacts, including the PixMo dataset, training code, evaluations, and intermediate checkpoints, available to the public.

**InternVL2\_5-8B:** This is a multimodal large language model developed as part of their InternVL 2.5 series. This model integrates a vision component, InternViT-300M-448px-V2\_5, with a language component, InternLM2\_5-7B-Chat, connected through a randomly initialized MLP projector. The architecture follows the "ViT-MLP-LLM" paradigm, employing a pixel unshuffle operation to reduce the number of visual tokens and a dynamic high-resolution strategy to handle various data types, including single images, multiple images, and videos. The training process is structured across three stages: MLP warmup, contrastive learning, and generative learning, aiming to enhance the model’s visual perception and multimodal capabilities. InternVL2\_5-8B (Chen et al., 2024) has demonstrated proficiency in tasks such as multimodal reasoning, OCR, chart and document understanding, and video comprehension.

**Qwen2-VL-7B-Instruct:** This is an advanced vision-language model developed by Qwen, designed to handle a variety of visual and textual tasks. This model supports arbitrary image resolutions, dynamically converting them into visual tokens for more human-like visual processing. Qwen2-VL (Wang et al., 2024a) achieves state-of-the-art performance on visual understanding benchmarks, including MathVista, DocVQA, RealWorldQA, MTVQA, etc. Additionally, it offers multilingual support, understanding texts in languages such as English, Chinese, most European languages, Japanese, Korean, Arabic, and Vietnamese.

**MiniCPM-V:** MiniCPM-V (Yao et al., 2024) is a multimodal large language model designed for deployment on devices ranging from GPU cards to mobile phones. By compressing image representations into 64 tokens via a perceiver resampler, it achieves high efficiency with reduced memory usage and faster inference speeds. Despite its compact size of 3 billion parameters, MiniCPM-V demonstrates state-of-the-art performance on multiple benchmarks, surpassing existing models of comparable size and even rivaling larger models like Qwen-VL-Chat. Notably, it supports bilingual multimodal interactions in English and Chinese, making it versatile for diverse applications.

**Janus-Pro-7B:** Janus-Pro-7B (Chen et al., 2025) is an advanced multimodal AI model developed by

DeepSeek, designed to unify text and image processing capabilities within a single framework. In text-to-image tasks, Janus-Pro-7B excels in generating high-quality images from textual descriptions, outperforming models like OpenAI’s DALL-E 3 and Stability AI’s Stable Diffusion in various benchmarks. For image-to-text tasks, Janus-Pro-7B employs a decoupled visual encoding approach, utilizing the SigLIP-L vision encoder to process images at resolutions up to 384x384 pixels. This design allows the model to effectively understand and generate textual descriptions of visual content, making it versatile for applications requiring both image generation and comprehension.

**GIT (Base and Large):** GIT (Generative Image-to-Text) (Wang et al., 2022) is a Transformer-based model developed by Microsoft for vision-language tasks such as image and video captioning, visual question answering (VQA), and image classification. The model is conditioned on both CLIP image tokens and text tokens, enabling it to generate textual descriptions based on visual inputs. GIT is available in two primary configurations:

- **GIT-Base:** This version comprises approximately 177 million parameters and is trained on 10 million image-text pairs.
- **GIT-Large:** This larger variant contains around 395 million parameters and is trained on 20 million image-text pairs. The expanded parameter count enhances its capacity to generate more detailed and accurate textual descriptions from images, making it well-suited for complex vision-language tasks.

Both versions utilize a Transformer decoder architecture, where the model has full bidirectional attention over image patch tokens and causal attention over text tokens. This design enables the models to predict the next text token by considering both the visual input and the preceding text, facilitating coherent and contextually relevant text generation based on images.

## E.2 Tag Classification

**Llama 3.2-1B Instruct:** The Llama 3.2 collection of multilingual large language models (LLMs) (Grattafiori et al., 2024) is a collection of pre-trained and instruction-tuned generative models in 1B and 3B sizes (text in/text out). The Llama 3.2 instruction-tuned text only models are optimized

for multilingual dialogue use cases, including agentic retrieval and summarization tasks. They outperform many of the available open source and closed chat models on common industry benchmarks.

**Llama-3.1-8B-Instruct:** Llama-3.1-8B-Instruct (Grattafiori et al., 2024) is an 8-billion-parameter language model developed by Meta as part of the Llama 3.1 series, released in July 2024. This model is fine-tuned for instruction-based tasks, enhancing its performance in understanding and generating human-like text responses. It supports eight languages: English, German, French, Italian, Portuguese, Hindi, Spanish, and Thai. Notably, Llama-3.1-8B-Instruct features an expanded context window of up to 128,000 tokens, allowing it to process and generate longer sequences of text effectively.

**Mistral 7B Instruct v0.2:** Mistral-7B-Instruct-v0.2 (Jiang et al., 2023) is an instruction fine-tuned version of the Mistral-7B-v0.2 language model, developed by Mistral AI. This iteration introduces key improvements over its predecessor, including an expanded context window of 32,000 tokens (up from 8,000), a RoPE-theta value of  $1e6$ , and the removal of Sliding-Window Attention. These enhancements enable the model to generate coherent and contextually rich responses, making it suitable for a wide range of natural language processing tasks.

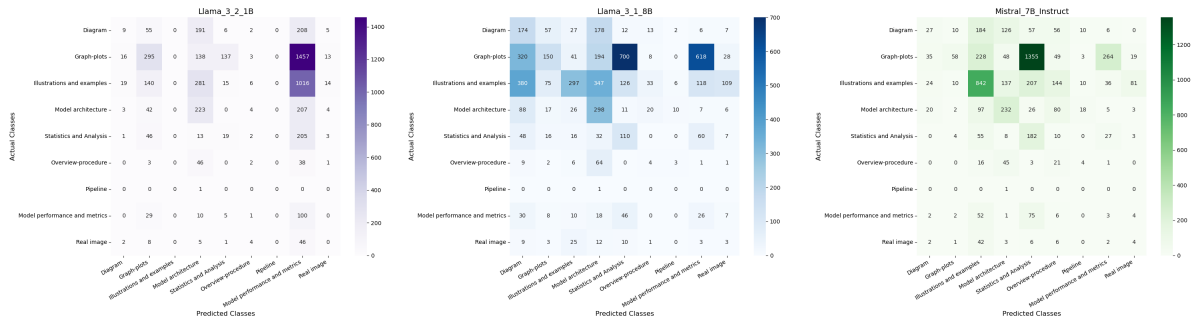
**Qwen2.5-0.5B-Instruct and Qwen2.5-7B-Instruct:** Qwen2.5-0.5B-Instruct (Team, 2024) is a 0.5 billion parameter instruction-tuned language model developed by the Qwen team at Alibaba Cloud. As part of the Qwen2.5 series, this model offers significant improvements in instruction following, coding, mathematics, and multilingual support across over 29 languages, including Chinese, English, French, and Spanish. It features a context length of up to 32,768 tokens and can generate sequences up to 8,192 tokens.

## F Task Prompts

In this section we provide the prompts that we have used for the various tasks in our study. Each prompt is designed to guide the model in performing specific operations, ensuring clarity, coherence, and consistency in the generated outputs. In a task, the same prompt is used for all models under comparative evaluation. Below, we list the prompts used under each task.

1224	<b>F.1 Figure Captioning</b>		1274
1225	<b>Zeroshot Captioning:</b>		1275
1226	"Generate a concise and articulate caption for a		1276
1227	diagram retrieved from a research paper. Focus		1277
1228	on explaining the key idea or concept represented		1278
1229	by the diagram in no more than 100 words. Avoid		1279
1230	describing the structural elements or layout of the		1280
1231	diagram, and ensure the caption is self-contained		1281
1232	and conceptually meaningful without external ref-		1282
1233	erences."		1283
1234	<b>Captioning using paper title as context:</b>		1284
1235	"Using the context provided in the following title:		1285
1236	<title>, generate a concise and meaningful caption		1286
1237	for the image that explains the key concept or core		1287
1238	idea represented by the figure in no more than 100		1288
1239	words."		1289
1240	<b>Captioning using paper title and abstract as con-</b>		1290
1241	<b>text:</b>		1291
1242	context = " Title: <title>		1292
1243	Abstract: <abstract>"		1293
1244	"Using the context provided below, generate a con-		1294
1245	cise and meaningful caption for the image that		1295
1246	explains the key concept or core idea represented		1296
1247	by the figure in no more than 100 words:		1297
1248	<context>"		
1249	<b>F.2 Tag Classification</b>		1298
1250	"You will receive a figure caption and a list of		1299
1251	predefined figure categories. Your task is to classify		1300
1252	the caption into exactly one category based on the		1301
1253	concept it represents. If the caption aligns with		1302
1254	multiple categories, choose the most appropriate		1303
1255	one that best describes the figure type.		1304
1256	<b>Categories:</b>		1305
1257	- Diagram: schematic figures or sketches.		1306
1258	- Graphs-plots: charts and plots.		1307
1259	- Illustrations and examples: figures providing		1308
1260	examples or visual aids.		1309
1261	- Model architecture: figures depicting the architec-		1310
1262	ture of models.		1311
1263	- Statistics and Analysis: figures or graphs		1312
1264	involving statistical results and analysis.		1313
1265	- Overview-procedure: figures that illustrate a high		1314
1266	level overview of methods or procedures.		1315
1267	- Pipeline: figures showing complete workflows.		1316
1268	- Model performance and metrics: figures or		1317
1269	graphs showing performance evaluation of models.		1318
1270	- Real image: photographs or realistic images.		
1271	<b>Examples:</b>		1319
1272	<i>Example 1:</i>		1320
1273	<i>Caption: A scatter plot showing the relationship</i>		1321
	<i>between training time and model accuracy, with a</i>		1322
	<i>trend line fitted to the data.</i>		1323
	<i>Your response: <b>Graph-plots</b></i>		
	<i>Example 2:</i>		
	<i>Caption: A step-by-step workflow illustrating the</i>		
	<i>data preprocessing, model training, and evaluation</i>		
	<i>stages in a deep learning pipeline.</i>		
	<i>Your response: <b>Pipeline</b></i>		
	<b>Instructions:</b>		
	- Identify the most relevant category for the caption.		
	- The classification must reflect <b>only one category</b> ,		
	avoiding overlaps. If multiple categories seem		
	relevant, choose the broadest and most appropriate		
	one		
	- Return only the category name. Do not add extra		
	explanation, reasoning, or special characters to		
	your response.		
	- Return the exact category name as it appears in		
	the list without any variations		
	<b>Figure Caption :</b>		
	<caption>		
	<b>Your Response:</b>		
	<b>F.3 Generating QA pairs using InternVL</b>		
	<image>		
	<caption>		
	Using the visual content of this image and the		
	context provided by the caption, generate 3 simple		
	and self-contained question-answer pairs.		
	Ensure that:		
	1. The questions are directly answerable using the		
	content of the image and/or the caption.		
	2. The questions are straightforward and do not		
	require multi-step reasoning.		
	3. The answers are contained entirely within the		
	image and caption.		
	4. The questions do not point to any external		
	references.		
	Provide the output in the format:		
	Q1: ...		
	A1: ...		
	Q2: ...		
	A2: ...		
	<b>F.4 Finetuning Qwen for Question-Answering</b>		
	You are a Vision Language Model specialized in		
	interpreting visual data from graphs, charts and		
	figures depicting statistical analysis. Your task is to		
	analyze the provided figure and respond to queries		

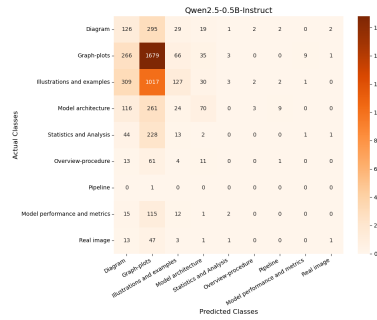
1324	<i>with concise and informative answers, usually in one or two sentences. Focus on delivering accurate, succinct answers based on the visual information. Avoid additional explanation unless absolutely necessary.</i>	model distinguishes between different tag categories.	1372
1325			1373
1326			
1327			
1328			
1329	<b>F.5 Using finetuned Qwen for question-answering on MathVista</b>	<b>H Human Evaluation Platforms</b>	1374
1330			
1331	<i>Provide a clear, short, and succinct numerical answer to the question based entirely on the given figure, without any external references or extra words. Follow the hint given below closely:</i>	<b>H.1 Evaluation of Figure Extraction and Classification in AI Figures</b>	1375
1332			1376
1333	<i>&lt;hint&gt;</i>	We present a view of the interface that was used by our evaluators for assessing the quality of figures and captions present in our dataset.	1377
1334	<i>&lt;question&gt;</i>		1378
1335	<i>Question:</i>		1379
1336	<i>&lt;question&gt;</i>	<b>H.2 Evaluation of Figure Captioning</b>	1380
1337	<i>Answer:</i>	We present a view of the interface that was used by our evaluators for assessing the quality of captions generated by each model under evaluation.	1381
1338			1382
1339	<b>F.6 Locality-Specific Question Response Generation on M3SciQA</b>		1383
1340			1384
1341	<i>You are given a figure, a question, and a list of paper candidates of titles and abstracts. Your task is to answer the question based on the figure information, then order the paper candidates that I provide to you so that the paper that is more relevant to the question comes first in the list. Return a minimum of 1 and a maximum of 5 paper candidates in the rank list. Ideally there should be 3 paper candidates.</i>	<b>I Finetuned SDXL Figure Generation results</b>	1385
1342	<i>Provide your answer at the end in a json file of this format using S2_id only:{{"ranking":[rank_1_s2_id, rank_2_s2_id] }}.</i>		1386
1343	<i>Make sure the responded list is in a valid format and that it only contains the S2_id. Do not include the title or abstract in the answer list. Also report the s2 ids in a comma separated manner.</i>		
1344	<i>&lt;question&gt;</i>		
1345	<i>{question}</i>		
1346	<i>&lt;/question&gt;</i>		
1347	<i>&lt;paper candidates&gt;</i>		
1348	<i>{reference_title_abstract_list}</i>		
1349	<i>&lt;/paper candidates&gt;</i>		
1350			
1351			
1352			
1353			
1354			
1355			
1356			
1357			
1358			
1359			
1360			
1361			
1362			
1363	<b>G Tag Classification: Confusion Matrix Evaluation</b>		
1364			
1365	We present confusion matrices for each model used in Tag classification task. These confusion matrices illustrate the distribution of predicted tags against the actual tags, highlighting patterns of correct and incorrect classifications.		
1366			
1367			
1368			
1369			
1370	By analyzing these matrices, we can identify common misclassifications and assess how well each		
1371			



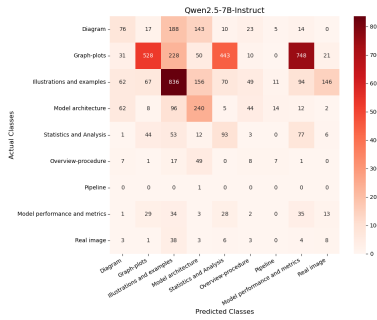
(a) Llama-3.2-1B

(b) Llama-3.1-8B

(c) Mistral-7B-v0.2



(d) Qwen2.5-0.5B



(e) Qwen2.5-7B

Figure 6: Confusion Matrix for Tag Classification

## Human Evaluation of Figures

Figure: Z26001\_page\_2\_illustrations and examples\_1\_figure.png

**Figure Category:** Illustrations And Examples

**Caption:** Figure J. Evolution of FDI among Turkey, Latin America/Caribbean, and the World

[View Full Paper](#)

Is the figure acceptable?

Yes

No

Select the issues in the figure:

Figure Classificat... x Figure Caption P... x

Figure 7: Human Evaluation of AI Figures

# Human Evaluation of Captions

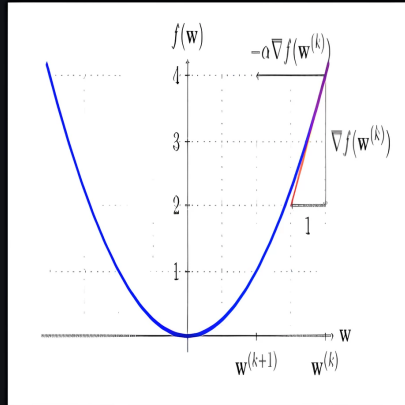


Figure: 3267\_page\_51\_Diagram\_1\_figure.png

Ground Truth: The GD step (55) amounts to a shift by  $\alpha \nabla f(w^k)$ .

## Qwen

The diagram illustrates the gradient descent method for minimizing a function  $f(w)$ . Starting at  $w^{(k)}$ , the gradient  $\nabla f(w^{(k)})$  is computed, and a step of size  $\alpha$  is taken along the negative gradient direction to reach  $w^{(k+1)}$ . The step size  $\alpha$  is chosen to ensure progress towards the minimum. The figure emphasizes the role of the gradient in guiding the descent towards the optimal solution.

Is the caption acceptable for qwen?

- Yes
- No

Select the issues in the caption:

Oversimplification x Contextual misin... x

Recognition error

minicpm

The diagram represents the relationship between the function  $f(x)$  and its derivative  $f'(x)$ . The slope of the curve indicates the rate of change of the function at any given point, while the x-axis represents the independent variable  $x$ . The y-axis represents the dependent variable  $f(x)$ . The curve shows the behavior of the function as  $x$  increases or decreases, with the slope of the curve indicating the rate of change of the function at any given point.

Is the caption acceptable for minicpm?

- Yes
- No

## Janus

The diagram illustrates the concept of the derivative of a function  $f(w)$  with respect to  $w$ , denoted as  $f'(w)$ . It shows the behavior of  $f'(w)$  as  $w$  approaches  $w^{(k)}$  and  $w^{(k+1)}$ . The graph demonstrates that  $f'(w)$  is a linear function, represented by the equation  $-\alpha \nabla f(w^{(k)})$ , where  $\alpha$

Is the caption acceptable for Janus?

- Yes
- No

## Molmo

This diagram illustrates the concept of diminishing returns in economic or ecological contexts. The blue curve represents a quantity that initially increases rapidly but then gradually slows down as it approaches a saturation point. The red line indicates the rate of change, which decreases as the quantity increases. The graph demonstrates how inputs or efforts, represented by the x-axis, have diminishing effects on outputs or benefits, represented by the y-axis, once a certain threshold is reached.

Is the caption acceptable for molmo?

- Yes
- No

## Internvl

The diagram illustrates an optimization process for a function  $f(w)$ . It shows a function  $f(w)$  with a decreasing rate, where the slope represents the gradient of the function at point  $w^{(k)}$  and  $w^{(k+1)}$ . The gradient vector  $\nabla f(w^{(k)})$  indicates the direction of the steepest ascent on the curve. The arrow labeled  $-\alpha \nabla f(w^{(k)})$  shows the descent direction, where  $\alpha$  is a step size parameter. As iterations progress, each step  $w^{(k+1)}$  results in a slight downward movement along the curve, suggesting a gradient descent approach towards minimizing  $f(w)$ .

Is the caption acceptable for Internvl?

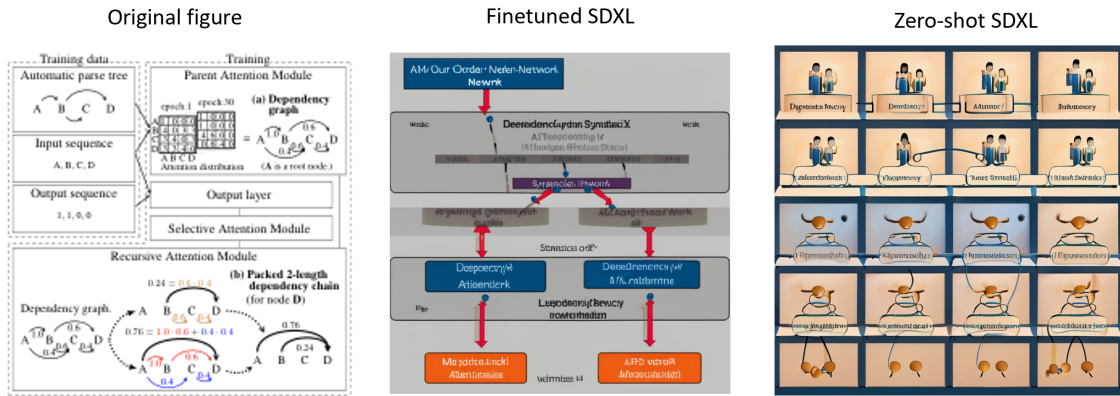
- Yes
- No

Submit Evaluation

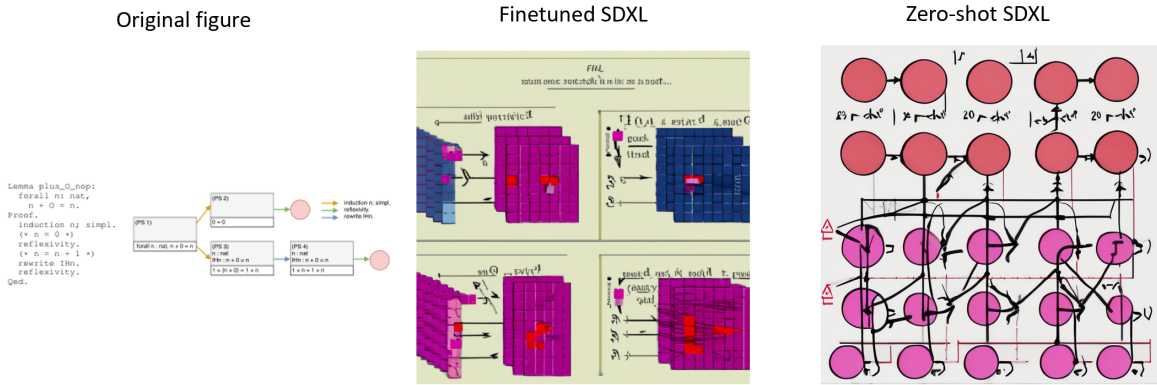
Previous

Next

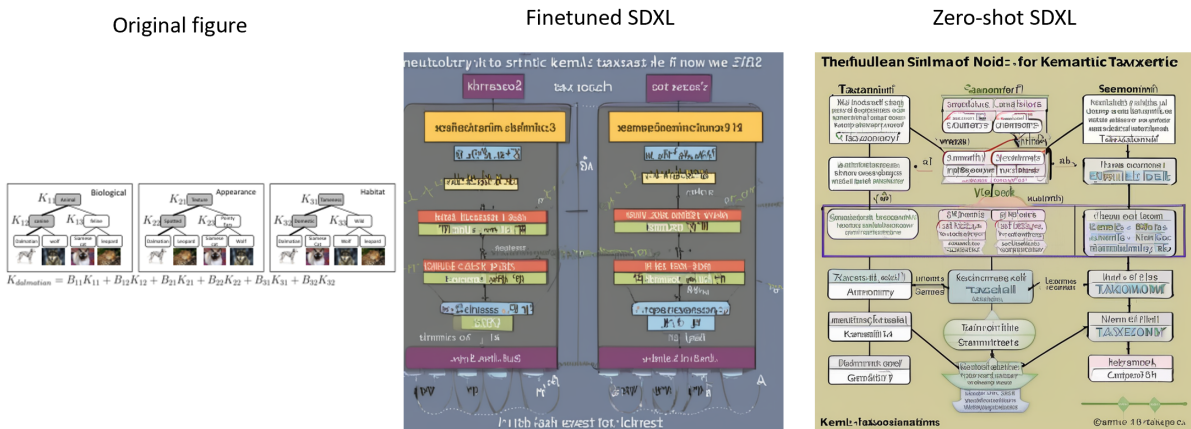
Figure 8: Human Evaluation of Captioning Results



Original Caption:  
 Dependency structures used in our higher-order syntactic attention network.



Original Caption:  
 A proof script in Cog (left) and the resulting proof states, proof steps, and the complete proof tree (right).



Original Caption:  
 Main idea: For a given set of classes, we assume multiple semantic taxonomies exist, each one representing a different view of the inter-class semantic relationships. Rather than commit to a single taxonomy which may or may not align well with discriminative visual features we learn a tree of kernels for each taxonomy that captures the granularity-specific similarity at each node. Then we show how to exploit the inter-taxonomic structure when learning a combination of these kernels from multiple taxonomies (i.e., a kernel forest) to best serve the object recognition tasks.

Figure 9: A comparative analysis of figure generations by the fine-tuned SDXL model, with the original figure and zero-shot generations from the base SDXL model.