# NOT SEARCH, BUT SCAN: BENCHMARKING MLLMS ON SCAN-ORIENTED ACADEMIC PAPER REASONING

**Anonymous authors** 

000

001

002003004

010 011

012

013

014

016

017

018

019

021

023

024

025

026

027

028

029

031

033

035

037

038

040

041

042 043

044

046

047

048

051

052

Paper under double-blind review

## **ABSTRACT**

With the rapid progress of multimodal large language models (MLLMs), AI already performs well at literature retrieval and certain reasoning tasks, serving as a capable assistant to human researchers, yet it remains far from autonomous research. The fundamental reason is that current work on scholarly paper reasoning is largely confined to a search-oriented paradigm centered on pre-specified targets, with reasoning grounded in relevance retrieval, which struggles to support researcher-style full-document understanding, reasoning, and verification. To bridge this gap, we propose ScholScan, a new benchmark for scholarly paper reasoning. ScholScan introduces a scan-oriented task setting that asks models to read and cross-check entire papers like human researchers, scanning the document to identify consistency issues. The benchmark comprises 1,800 carefully annotated questions drawn from 9 error families across 13 natural-science domains and 715 papers, and provides detailed annotations for evidence localization and reasoning traces, together with a unified evaluation protocol. We assessed 15 models across 24 input configurations and conduct a fine-grained analysis of MLLM capabilities across error families. Across the board, retrieval-augmented generation (RAG) methods yield no significant improvements, revealing systematic deficiencies of current MLLMs on scan-oriented tasks and underscoring the challenge posed by ScholScan. We expect ScholScan to be the leading and representative work of the scan-oriented task paradigm.

## 1 Introduction

Scientific papers are crystallizations of human intelligence. Enabling multimodal large language models (MLLMs) (OpenAI, 2025; Anthropic, 2025; ByteDance Seed Team, 2025; Meta, 2025; xAI, 2025) to conduct comprehensive understanding and generation based on academic literature is the ultimate goal of Deep Research, and a critical milestone on the path toward artificial general intelligence (AGI) (Ge et al., 2023; Morris et al., 2024; Phan et al., 2025). With rapid advances, MLLMs are increasingly capable of supporting academic workflows through retrieval, reading, and writing. For example, PaSa (He et al., 2025) can invoke a series of tools to answer complex academic queries with high-quality results, while Google Deep Research (Comanici et al., 2025) is capable of producing human-level research reports based on specific queries.

However, most of the existing work still follows *a search-oriented paradigm*, where models retrieve a few relevant passages and reason over local evidence based on prespecified targets (Gao et al., 2023; Lou et al., 2025). Such methods are effective for tasks with clearly predefined targets, but struggle with researcher-style full-document reasoning and verification (Zhou et al., 2024). *To function as researchers, models must move beyond reactive question answering and toward proactive discovery of implicit problems.* 

To fill this gap, as shown in Figure 1, we introduce *a scan-oriented paradigm*, where models address queries with targets absent and are required to actively **construct a document-level evidence view**, **perform exhaustive scanning over the full paper**, and **conduct evidence-based reasoning**. In contrast to search-oriented tasks that assess a model's ability to identify and reason over *relevant* fragments, scan-oriented tasks emphasize *consistency*. *Instead of relying on prespecified targets or hints, models must derive all necessary concepts and inferences solely from given documents*.

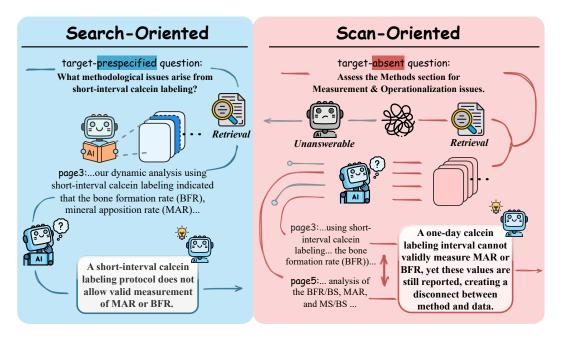


Figure 1: A comparison between search-oriented and scan-oriented task paradigms. Unlike the former, the scan-oriented paradigm provides no prespecified targets, requiring the model to actively scan the entire paper, construct a document-level evidence view.

We instantiate this setting via scientific error detection, as it naturally demands discovering nonobvious flaws without target cues, and present ScholScan, a new multimodal benchmark for scholarly reasoning. ScholScan features the following key highlights:

- Scan-Oriented Task Paradigm. ScholScan receive one or more complete academic papers together with target-absent queries, presenting a rigorous challenge to their evidence-based reasoning capabilities. The benchmark comprises 715 papers spanning 13 natural science disciplines.
- Comprehensive Error Types. ScholScan covers 9 categories of scientific errors across the entire research workflow. It also includes citation and referencing errors, providing a rigorous test of a model's cross-source reasoning ability.
- Process-Aware Rvaluation Framework. ScholScan provides fine-grained annotations for both
  evidence location and reasoning steps, enabling a comprehensive evaluation framework that assesses model performance in terms of both process and outcome.

We evaluate 15 models across 24 input configurations and 8 retrieval-augmented generation (RAG) frameworks. All models exhibit limited performance, and none of the RAG methods deliver significant improvements. These results highlight the inadequacy of search-oriented frameworks when applied to scan-oriented tasks, and underscore both the challenges and the potential of enabling MLLMs to perform reliable, document-level reasoning over full academic papers.

### 2 Related Work

## 2.1 MULTIMODAL LARGE LANGUAGE MODELS

With the rapid progress of MLLMs, models have evolved beyond perception tasks (e.g., image recognition and explanation) (Liu et al., 2024) toward deep understanding of structured, multimodal long documents. Their strengths lie in the ability to integrate cross-modal information and perform multi-hop reasoning over extended contexts. These capabilities are not only valuable for specific question answering or instruction-following tasks (Yue et al., 2024) but are particularly well suited for simulating human thought processes and generating explainable reasoning trajectories (Zheng et al., 2023). Consequently, achieving comprehensive understanding of entire documents has emerged as a core challenge that MLLMs are inherently equipped to address.

### 2.2 DOCUMENT UNDERSTANDING BENCHMARK

Document understanding tasks challenge models to identify relevant context and perform accurate reasoning grounded in that information. Progress in document understanding benchmarks has followed two main axes. Along the input dimension, it has evolved from short to long contents, from everyday to specialized domains, and from plain text to multimodal format (Chen et al., 2021; Yang et al., 2018; Tito et al., 2021; Deng et al., 2025). Along the scenario dimension, it has shifted from limited-output formats to more open-ended responses (Pramanick et al., 2024). DocMath-Eval (Zhao et al., 2024) evaluates numerical reasoning on long, specialized documents, revealing large performance gaps even for strong models in expert domains, while MMLongBench-Doc (Ma et al., 2024) builds a multimodal benchmark with layout-rich documents. However, a comprehensive benchmark that integrates all challenges above has yet to be introduced.

#### 2.3 ACADEMIC PAPER UNDERSTANDING BENCHMARK

Compared with general documents, academic papers are distinguished by their rich domain knowledge and logical rigor. Reasoning over papers has emerged as a major challenge in recent research. Some studies ask for local elements like charts or snippets, leveraging their internal complexity, but neglect the need for cross-source integration and domain-specific interpretation within the full document (Wang et al., 2024; Li et al., 2024). Recent studies extend inputs to the document level and adopt image-based formats to better simulate real-world reading scenarios. (Auer et al., 2023; Yan et al., 2025) However, benchmarks based on the QA paradigm face inherent limitations, as they typically presuppose answer existence and embed explicit cues in the question itself, reducing the need for comprehensive understanding and information organization. Moreover, mainstream evaluation protocols focus on the final outcome, with limited assessment of whether intermediate reasoning is evidentially grounded and logically valid. More examples and analysis are shown in Appendix C.

## 3 THE SCHOLEVAL BENCHMARK



Benchmark	Mod.	Para.	Eval.	# Dom.
Document Understand	ing			
DocMath-Eval <sub>CompLong</sub>	T+TD	Search	Α	N/A
MMLongbench-Doc	T+MD	Search	Α	N/A
LongDocURL	T+MD	Search	Α	N/A
SlideVQA	T+MD	Search	A	N/A
Academic Paper Under	rstandin	g		
CharXiv	I	Search	Α	8
ArXivQA	I	Search	Α	10
MMCR	T+MD	Search	Α	CS
AAAR-1.0	T+MD	Search	A	CS
ScholScan (ours)	T+MD	Scan	A+P	13

Figure 2: Left: Overview of ScholScan. Right: Comparison to related benchmarks. **Mod.**: Modalities; **Para.**: Task Paradigm; **Eval.**: Evaluation; **T**: Text; **I**: Image; **TD**: Text-Form Document; **MD**: Multimodal Document; **A**: Answer; **P**: Process; **Dom**: Number of academic domains in the dataset.

## 3.1 OVERVIEW OF SCHOLSCAN

We introduce ScholScan, a benchmark designed to comprehensively evaluate MLLMs' ability to detect scientific flaws in academic papers under scan-oriented task settings. As illustrated in Figure 2, ScholScan spans 13 disciplines across the natural sciences, including physics, chemistry, and computer science, and spans over 100 subfields such as immunology, total synthesis, and machine learning. The benchmark comprises 1,800 questions derived from 715 real academic papers, and covers 9 major error categories (Figure 3) that commonly observed in real-world research scenarios. These include issues in numerical and formulaic computation, experimental design, inference and conclusion, and citation misuse, among others. Figure 2 also provides a comparison ScholScan with existing benchmarks for multimodal paper understanding and long-document reasoning.

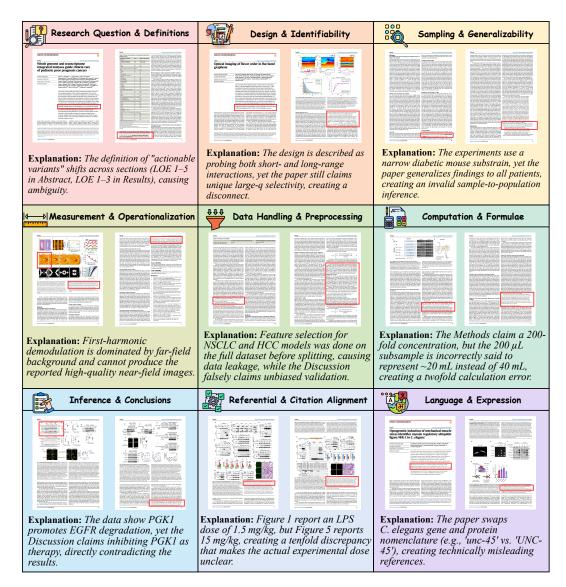


Figure 3: Sampled ScholScan examples with 9 error types, covering the whole process of scientific research, each requiring the model to perform thorough cross-source evidence-based reasoning.

## 3.2 Data Collection & Question Generation

We curated papers from ICLR 2024/2025 and Nature Communications, and collected public reviews for the former. Questions were constructed based on two dimensions, where the source is either generated or sampled, and the context is either within-paper or cross-paper.

**Generation.** On high-quality accepted papers, we prompt Gemini 2.5 Pro to perform coordinated sentence-level edits spanning multiple sections or pages. It then synthesizes composite errors and generates the corresponding question along with an explanation grounded in the edited context.

**Sampling.** From rejected ICLR submissions and their public reviews, we prompt Gemini 2.5 Pro to extract explicit, falsifiable scientific errors and convert them into questions with initial explanations. Subjective remarks about novelty or writing quality are excluded.

**Within-paper.** This setting focuses on verifiable facts and internal consistency within a single paper, and supports both Generation and Sampling.

**Cross-paper.** This setting examines citation consistency across papers. For each instance, Gemini 2.5 Pro receives an accepted paper and one of its cited sources, then edits the accepted paper to

introduce paraphrases or reasoning errors about the citation. As public reviews mainly address nonfalsifiable aspects such as appropriateness, all cross-paper instances are constructed exclusively using the generation method.

## 3.3 QUALITY CONTROL & ANNOTATION

Despite explicit instructions, initial outputs exhibited substantial hallucinations, logical inconsistencies, and low-quality questions. To ensure the quality, 10 domain experts conducted a rigorous annotation process. Each instance underwent independent dual review, and disagreements were resolved by a third expert. Among the 3,500 initially generated candidates, 1,700 were discarded, and 1,541 of the remaining were revised, including 535 question rewrites, 1,207 explanation edits, and 1,141 corrections to error categories or metadata. Further details are provided in Appendix D.

## 4 EXPERIMENTS

#### 4.1 Experiments Setting

**Models.** We benchmark a total of 24 input configurations by feeding academic papers as either images or OCR text using the Tesseract (Smith, 2007) engine, covering 15 mainstream models (Yang et al., 2025; Bai et al., 2025; DeepSeek-AI et al., 2025; Guo et al., 2025; OpenAI et al., 2025).

**Evaluation Protocol.** Inspired by MMLongBench-Doc (Ma et al., 2024), we prompt models to generate necessary reasoning chains from evidence to detected anomalies without constraining the output format, which aims to assess the ability for evidence-grounded reasoning rather than mere instruction-following. For open-ended responses, we use GPT-4.1 (OpenAI, 2025) to extract cited evidence and reasoning steps, and quantify alignment with annotated explanations. Human evaluation confirms high agreement between our pipeline and expert annotations. Further implementation details are provided in Appendix F.

**Metrics.** We define a structured evaluation framework by parsing the model response a into a tuple:

$$\Psi(a) \Rightarrow (\mathbf{1}_{\text{exist}}, \mathbf{1}_{\text{contain}}, \widehat{\mathcal{E}}, \widehat{\mathcal{R}}, n).$$
 (1)

Here,  $\mathbf{1}_{\text{exist}}$  and  $\mathbf{1}_{\text{contain}}$  are binary indicators for whether output contains any error and includes the annotated target error;  $\widehat{\mathcal{E}}, \widehat{\mathcal{R}}$  and  $\mathcal{E}^*, \mathcal{R}^*$  are the predicted and gold evidence sets and reasoning chains;  $\widehat{g} = \text{prefix\_match}(\widehat{\mathcal{R}}, \mathcal{R}^*)$  counts matched reasoning steps;  $n \in \mathbb{N}$  is the number of unrelated errors. HasError(a) is 1 if the output contains any predicted error, and 0 otherwise. Based on  $\Psi(a)$ , we define an end-to-end score  $S(m) \in [0,1]$  that combines all aspects of prediction quality:

(i) Existence.  $S_{\text{exist}}(a) = 1$  if and only if the response includes the annotated target error.

$$S_{\text{exist}}(a) = \mathbf{1}\{\text{HasError}(a)\} \cdot \mathbf{1}\{\hat{\mathcal{E}} \cap \mathcal{E}^* \neq \emptyset\}$$
 (2)

(ii) Evidence location score. Even when the target error is identified, the cited evidence may be incomplete or noisy. We compute a Dice score with a squared penalty for over-reporting:

$$S_{\text{location}} = \max \left\{ 0, \ \frac{2 \left| \widehat{\mathcal{E}} \cap \mathcal{E}^* \right| + 1 \left\{ \left| \widehat{\mathcal{E}} \right| + \left| \mathcal{E}^* \right| = 0 \right\}}{\max \left( \left| \widehat{\mathcal{E}} \right| + \left| \mathcal{E}^* \right|, 1 \right)} - 0.8 \left( \frac{\left| \widehat{\mathcal{E}} \setminus \mathcal{E}^* \right|}{\max \left( \left| \widehat{\mathcal{E}} \right|, 1 \right)} \right)^2 \right\}.$$
(3)

(iii) Reasoning process score. Even if the target error is detected, the reasoning may diverge from the gold chain. We use prefix match to assess reasoning completeness:

$$S_{\text{reasoning}} = \mathbf{1} \{ g_r = 0 \} + \mathbf{1} \{ g_r > 0 \} \left( \frac{\hat{g}}{g_r} \right)^2.$$
 (4)

(iv) Unrelated-error penalty. Models may list unrelated items to inflate recall at the cost of precision. We penalize this with a rapidly increasing function of unrelated error count:

$$P_{\text{unrelated\_err}}(n) = 0.9^{\min(n,2)} \exp(-0.6 \left[\max(n-2,0)\right]^{1.5}).$$
 (5)

(v) Overall outcome score. The final score for a is defined as:

$$S(m) = S_{\text{exist}}(a) \sqrt{S_{\text{location}} \cdot S_{\text{reasoning}}} \cdot P_{\text{unrelated\_err}}(n). \tag{6}$$

Table 1: Model performance (scaled by 100) across input configurations. **RQD**: Research Question & Definitions; **DI**: Design & Identifiability; **SG**: Sampling & Generalizability; **MO**: Measurement & Operationalization; **DHP**: Data Handling & Preprocessing; **CF**: Computation & Formulae; **IC**: Inference & Conclusions; **RCA**: Referential and Citation Alignment; **LE**: Language & Expression.

Models	Avg.	RQD	DI	SG	МО	DHP	CF	IC	RCA	LE
MLLM (Image Input)										
Proprietary MLLMs										
Gemini 2.5 Pro	<u>15.6</u>	11.9	12.6	35.7	12.3	27.0	4.6	14.7	<u>15.2</u>	7.4
GPT-5	19.2	10.1	9.7	28.2	14.6	26.6	13.8	25.3	25.3	6.9
Grok 4	4.0	0.0	1.9	16.7	3.2	7.4	0.7	1.9	3.6	0.0
Doubao-Seed-1.6-thinking	10.2	3.4	3.5	22.3	7.5	15.1	10.2	12.2	10.9	3.3
Doubao-Seed-1.6	9.9	3.0	4.4	29.2	4.9	15.0	6.3	17.9	8.0	3.9
Open-source LLMs										
Llama 4 Maverick	7.0	7.0	7.3	9.4	4.5	4.0	6.5	6.7	8.8	3.0
Gemma 3 27B	1.7	0.5	2.7	2.3	1.7	1.0	1.0	1.3	2.6	0.0
Mistral Small 3.1	3.3	0.1	2.0	2.0	1.5	0.1	1.0	2.2	8.6	1.0
Qwen2.5 VL 72B	0.1	0.0	0.7	0.0	0.0	0.0	0.0	0.0	0.2	0.0
		OCI	R + LLN	I (Text ]	Input)					
Proprietary LLMs										
Gemini 2.5 Pro	30.3	21.5	34.2	44.3	27.6	56.6	10.3	28.8	35.6	8.1
GPT-5	22.5	16.1	21.4	26.0	20.3	36.7	4.7	29.8	30.0	2.6
Claude Sonnet 4	5.7	3.7	2.5	10.8	4.3	10.3	1.4	8.4	6.6	3.5
Grok 4	20.8	9.3	7.7	37.4	12.3	34.4	9.0	20.0	31.2	7.2
Doubao-Seed-1.6-thinking	15.3	8.2	10.1	24.3	10.1	24.2	6.4	19.2	21.0	4.2
Doubao-Seed-1.6	13.9	5.4	6.9	26.4	10.3	23.6	6.3	20.1	17.5	2.3
Open-source LLMs										
Qwen3 A22B (Thinking)	17.4	8.9	16.2	31.9	15.1	23.7	5.6	22.3	21.1	2.3
Qwen3 A22B	1.7	1.2	0.0	2.7	0.4	1.0	0.1	4.3	2.5	1.1
gpt-oss-120b	7.3	6.3	5.7	18.3	4.9	14.5	1.6	12.5	5.5	0.0
DeepSeek-R1	11.4	5.1	11.9	25.4	8.7	22.5	4.7	16.3	9.8	3.5
DeepSeek-V3.1	1.7	1.2	2.0	1.7	1.0	5.8	0.5	2.2	2.1	0.0
Llama 4 Maverick	2.3	1.5	2.0	4.8	3.0	3.6	0.0	5.8	1.6	0.2
Gemma 3 27B	2.0	2.1	1.6	3.0	2.7	0.2	0.7	7.7	1.0	0.0
Mistral Small 3.1	6.9	3.0	2.7	5.5	7.0	2.0	8.5	4.0	12.2	3.0
Owen2.5 VL 72B	0.2	0.0	0.7	0.0	0.0	0.0	0.0	0.0	0.6	0.0

### 4.2 Main Result

Table 1 presents our evaluation results. Our main findings are summarized as follows:

Overall performance remains unsatisfactory. GPT-5 achieves the highest average score in the image input group (19.2), while Gemini 2.5 Pro, the best-performing model in the text input setting, still fails to surpass the 60-point threshold on any subtask. Even in the SG category, which yields the best performance overall, nearly half of the models receive single-digit scores. Most models perform poorly under the scan-oriented task formulation and fail to detect any issues in many papers. This challenge is particularly pronounced for open-source models.

Reasoning-enhanced models demonstrate clear advantages. Across both input configurations, reasoning-enhanced variants consistently achieve higher scores. Almost all top-performing models, measured by both subtask-specific and overall metrics, fall into this category. Notably, Qwen3-Thinking and Deepseek-R1 outperform their base versions by more than 10% in average scores, with substantial gains observed across all error types. These results indicate that reasoning-enhanced models are better able to simulate the iterative process of extraction followed by reasoning, which is essential for effectively handling scan-oriented tasks and producing higher-quality responses.

MLLMs face significant bottlenecks in handling long multimodal inputs. Across most evaluation metrics, text inputs outperform image inputs. Among the nine MLLMs tested, the average performance gap between text and image inputs reaches 4.81 points, highlighting visual processing as a key limitation in current MLLM capabilities.

In most evaluation metrics, text inputs consistently outperform image inputs. Among the nine MLLMs evaluated, the average performance gap between text and image inputs is 4.81 points, underscoring visual processing as a key limitation in current MLLM capabilities.

Although overall performance is generally weaker, multimodal input remains indispensable. In certain categories such as CF, where OCR-based text extraction leads to substantial loss of formulaic or tabular content, image inputs outperform their text counterparts. This highlights the essential role of multimodal reasoning and the irreplaceable value of visual information in addressing specific types of errors.

#### 4.3 FINE-GRAINED ANALYSIS

**Capability Dimensions.** We compute pairwise Spearman correlations between error types across two input configurations (text and image) for the eight evaluated MLLMs excluding Qwen2.5-VL-72B, as shown in Figure 4. We derive the following insights:

(i) With image input, CF exhibits consistently low correlations with other error categories, suggesting that the skills required for mathematical reasoning are relatively distinct. In contrast, with text input, CF shows moderate correlation with LE, indicating that OCR-flattened formulas lose their structural specificity and are interpreted by models in a manner more akin to natural language. Combined with the overall poor performance on CF tasks, this underscores the unique challenges of this category and the need for targeted improvements.

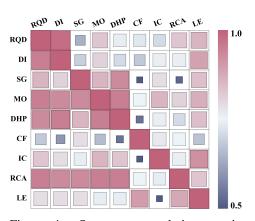


Figure 4: Spearman correlation matrix among the 9 error types.

(ii) Although DI is also related to experimental settings, it does not exhibit strong correlations with SG,

*MO*, *or DHP*. This indicates that DI primarily emphasizes causal framing and variable identifiability, rather than the procedural understanding of experimental operations.

(iii) OCR severely degrades structured content such as figures and formulas, making questions that depend on multimodal information unanswerable. This diminishes the expression of multimodal reasoning capabilities and artificially inflates inter-category correlations under text input.

Based on the above analysis, we consolidate the original 9 error categories, each defined by its objective target, into 5 core latent skill dimensions evaluated by ScholScan under the image input setting. While each dimension highlights the primary competence emphasized by its corresponding error types, they are not mutually exclusive, as many questions involve overlapping reasoning abilities.

RQD and DI correspond to research concept comprehension, which requires models to *identify the scope and definition* of research objectives by integrating contextual cues and prior knowledge. SG, MO, and DHP fall under *experimental process modeling*, which tests a model's ability to reconstruct procedural workflows such as sampling, measurement, and data handling. CF captures *formal reasoning and symbolic computation*, focusing on syntactic parsing and numerical logic. IC evaluates causal inference, where models must *synthesize dispersed causal evidence* to reach sound conclusions. RCA and LE reflect referential alignment and linguistic consistency, which assess the ability to *verify citations and maintain coherent expression* throughout the document.

**Hidden Complexity in Scan-Oriented Tasks.** We analyze the reasoning traces of GPT-5 and Gemini 2.5 Pro under both input configurations, focusing on the number of evidence pieces scanned and the reasoning steps performed. As illustrated in Figure 5, even the most advanced models often scan up to 8 times more evidence and execute 3.5 times more reasoning steps than the reference answers, merely to approximate a correct response, yet they still frequently fail. This highlights the substantial hidden complexity inherent in scan-oriented tasks, which significantly amplifies the challenge of successful task completion.

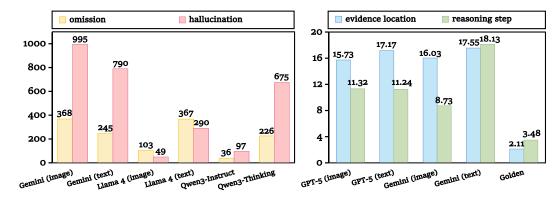


Figure 5: Left: Distribution of omission and hallucination errors. Right: Average reasoning steps and evidence locations involved in the answer generation, compared against the golden reference.

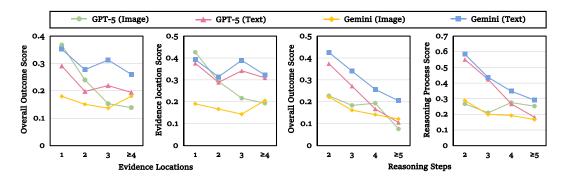


Figure 6: Performance trends across varying reasoning depths and evidence counts.

### 4.4 ERROR ANALYSIS

**Omission and Hallucination.** Most zero-score cases fall into two categories: either the model fails to detect any errors in the paper, or it becomes overwhelmed by hallucinations and entirely overlooks the actual errors present in the reference answer. We analyze the number of zero-score questions and the proportion of these two failure modes across models, as shown in Figure 5. Stronger models tend to have fewer zero-score cases overall, but are more prone to overconfident hallucinations.

**Fragile Reasoning under Complex Evidence.** Figure 6 shows how top-performing models behave under different numbers of reasoning steps and evidence locations. As reasoning steps increase, both reasoning and overall scores steadily decline, revealing a clear bottleneck in MLLMs' ability to construct long causal chains. In contrast, variation in evidence count has a weaker and less consistent impact. However, this does not imply that multi-evidence questions pose only marginal difficulty. Since the evaluation metric allows partial evidence omissions, more evidence items do not necessarily incur large score penalties. Still, heavier evidence loads often require longer reasoning chains, which substantially affect the coherence and completeness of inferred logic. These results highlight the persistent challenge for MLLMs in integrating evidence and maintaining logical structure as task complexity grows.

## 4.5 RAG ANALYSIS

We evaluated 8 RAG methods under both input configurations (Robertson et al., 1994; Chen et al., 2024; Lee et al., 2025; Faysse et al., 2025; Yu et al., 2025; Wang et al., 2025; Izacard et al., 2022). Key findings are presented below, with detailed results shown in Tables 2 and 3.

**Oracle Condition Yields Significant Accuracy Gains.** Providing gold-standard images alleviates the scanning burden in long-context inputs, increasing the chances of generating correct answers. While overall performance improves, gains are limited for CF errors and minimal for LE errors. For

Table 2: Scores of RAG methods across the 9 error types (scaled by 100).

Models	Avg	RQD	DI	SG	MO	DHP	CF	IC	RCA	LE
Text Input (Base Mod	lel: Qwei	n3 Thinkir	ıg)							
Baseline	17.4	8.9	16.2	31.9	15.1	23.7	5.6	22.3	21.1	2.3
Oracle	24.5	20.6	27.9	43.6	21.3	40.8	7.4	26.9	26.0	1.9
bm25	16.7	9.7	13.7	33.0	17.3	23.8	6.8	25.4	16.5	3.0
BGE-M3	11.3	8.6	7.5	24.8	9.1	15.4	5.3	15.6	11.4	1.0
Contriever-msmacro	16.6	9.7	18.2	33.7	10.7	20.8	6.4	18.5	19.8	1.8
nv-embed-v2	6.8	4.0	4.0	9.4	6.1	4.9	5.5	5.7	10.0	2.0
Image Input (Base M	odel: Lla	ama4 Mav	erick)							
Baseline	7.0	7.0	7.3	9.4	4.5	4.0	6.5	6.7	8.8	3.0
Oracle	6.5	3.0	4.5	15.6	8.2	9.4	4.9	10.0	4.4	1.4
ColPali-v1.3	0.8	1.5	0.0	0.5	0.0	0.9	0.5	1.3	1.4	0.0
ColQwen2.5	1.2	2.1	0.7	0.5	0.0	1.2	0.2	2.7	2.0	0.0
VisRAG	1.0	2.0	0.0	1.0	0.0	1.0	1.6	1.3	1.2	0.0
VRAG-RL	10.9	9.8	11.6	17.8	8.2	11.0	6.8	13.1	10.8	8.1

CF, sparse formulaic content means gold images offer slight help. For LE, dense text distribution makes even direct access to target regions insufficient to reduce complexity for current models.

In consistency-centric scan-oriented tasks, most retrieval-based enhancement methods show minimal effectiveness. All embedding models exhibit poor retrieval accuracy. None achieves recall of 50% within the top-5 retrieved items. More critically, performance deteriorates after retrieval, especially for multimodal embedding models, where post-retrieval responses are almost entirely incorrect and scores approach 0.

Complex embedding model architectures do not yield better performance. Providing gold-standard images alleviates the scanning burden in long-context inputs, increasing the chances of retrieving correct answers. While overall performance improves, gains are limited for CF and minimal for LE errors. For CF, sparse formulaic content means gold images offer only

Table 3: Summary of retrieval performance for RAG methods.

Models	MRR@5	Recall@5
Text Input (Base Mo	del: Qwen3 '	Thinking)
bm25	0.41	0.48
BGE-M3	0.16	0.21
Contriever-msmacro	0.31	0.39
nv-embed-v2	0.30	0.38
Image Input (Base M	Iodel: Llama	14 Maverick)
ColPali-v1.3	0.26	0.31
ColQwen2.5	0.30	0.35
VisRAG	0.41	0.46

slight localization help. For LE, dense error distribution makes even direct access to target regions insufficient to reduce task complexity for current models.

Reinforcement learning frameworks with a visual-centric focus have distinguished themselves as leading approaches. Despite being built on a compact 7B model, VRAG-RL consistently delivers improved performance and is the only method that achieves gains in the image-input setting following RL optimization. Its enhanced retrieval sharpens evidence selection, while strong reasoning provides effective guidance during document scanning. The retrieval and reasoning components are interleaved in design, with each stage informing the other in an iterative loop. This tightly coupled interaction contributes to the method's superior performance potential.

## 5 CONCLUSION

In this paper, we introduce ScholScan, a benchmark designed to evaluate the performance of MLLMs on scan-oriented tasks that require detecting scientific errors across entire academic papers. We conduct a comprehensive evaluation and in-depth analysis of mainstream MLLMs and RAG methods. The results demonstrate that current MLLMs remain far from capable of reliably addressing such tasks, and that existing RAG approaches provide little to no improvement. This highlights the complexity, integrative demands, and originality of the ScholScan benchmark. Looking ahead, we aim to develop scan-oriented task paradigms suited to diverse academic scenarios and explore new techniques for enhancing model performance on target-suppressed inputs. These directions support the broader goal of advancing MLLMs from passive assistants to active participants in scientific research.

## 6 ETHICS STATEMENT

All data used in this paper were constructed by the authors and do not include any external public or proprietary datasets. The included academic papers and author names are publicly available through arXiv and OpenReview and can be freely accessed.

A team of 10 domain experts was assembled to comprehensively review all task instances initially generated by Gemini 2.5 Pro. All annotators gave informed consent to participate. To ensure the accuracy and neutrality of both model-generated and human-verified content, we employed a rigorous multi-stage validation process involving cross-review and third-party adjudication.

Evaluation across 15 mainstream models and 24 input configurations was conducted via legally authorized API access through the VolcEngine, Alibaba Cloud's LLM services, and OpenRouter.

ScholScan is fully open-sourced and freely available for academic and non-commercial research purposes. We provide the complete download link and documentation through an anonymous GitHub repository. All personally identifiable information has been removed from the dataset, and its collection and release comply with the ethical and legal requirements in place at the time of data acquisition.

## 7 REPRODUCIBILITY STATEMENT

All results presented in this paper are fully reproducible. To facilitate verification and extension, we provide an anonymous repository (https://anonymous.4open.science/r/ScholScan-6657/) that contains the complete dataset, source code, and detailed documentation. The repository also includes step-by-step instructions and the exact hyperparameter configurations used in our experiments, ensuring that other researchers can replicate our findings with minimal effort.

The retrieval components in all retrieval-augmented generation (RAG) experiments were executed on a server equipped with 8 NVIDIA A40 GPUs.

### REFERENCES

- Anthropic. System card: Claude opus 4 & claude sonnet 4. https://www.anthropic.com/claude-4-system-card, May 2025. Updated Sep 2, 2025.
- S. Auer, Dante Augusto Couto Barone, Cassiano Bartz, E. Cortes, Mohamad Yaser Jaradeh, Oliver Karras, Manolis Koubarakis, Dmitry I. Mouromtsev, Dmitrii Pliukhin, Daniil Radyush, Ivan Shilin, Markus Stocker, and Eleni Tsalapati. The sciqa scientific question answering benchmark for scholarly knowledge. *Scientific Reports*, 13, 2023. URL https://api.semanticscholar.org/CorpusID:258507546.
- Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report, 2025. URL https://arxiv.org/abs/2502.13923.
- ByteDance Seed Team. Introduction to techniques used in seed1.6. https://seed.bytedance.com/en/blog/introduction-to-techniques-used-in-seed1-6, June 2025. Official blog post describing Seed1.6 techniques.
- Jianly Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge m3-embedding: Multi-lingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation, 2024. URL https://arxiv.org/abs/2402.03216.
- Zhiyu Chen, Wenhu Chen, Charese Smiley, Sameena Shah, Iana Borova, Dylan Langdon, Reema Moussa, Matt Beane, Ting-Hao Huang, Bryan Routledge, and William Yang Wang. FinQA: A dataset of numerical reasoning over financial data. In Marie-Francine Moens, Xuanjing Huang,

541

542

543

544

546

547

548

549

550

551

552

553

554

558

559

561

564

565

566

567

568

569

571

572

575

576

577

578

579

581

582

583

584

585

588

592

Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 3697–3711, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021. emnlp-main.300. URL https://aclanthology.org/2021.emnlp-main.300/.

Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni, Nathan Lintz, Tiago Cardal Pais, Henrik Jacobsson, Idan Szpektor, Nan-Jiang Jiang, Krishna Haridasan, Ahmed Omran, Nikunj Saunshi, Dara Bahri, Gaurav Mishra, Eric Chu, Toby Boyd, Brad Hekman, Aaron Parisi, Chaoyi Zhang, Kornraphop Kawintiranon, Tania Bedrax-Weiss, Oliver Wang, Ya Xu, Ollie Purkiss, Uri Mendlovic, Ilaï Deutel, Nam Nguyen, Adam Langley, Flip Korn, Lucia Rossazza, Alexandre Ramé, Sagar Waghmare, Helen Miller, Nathan Byrd, Ashrith Sheshan, Raia Hadsell Sangnie Bhardwaj, Pawel Janus, Tero Rissa, Dan Horgan, Sharon Silver, Ayzaan Wahid, Sergey Brin, Yves Raimond, Klemen Kloboves, Cindy Wang, Nitesh Bharadwaj Gundavarapu, Ilia Shumailov, Bo Wang, Mantas Pajarskas, Joe Heyward, Martin Nikoltchev, Maciej Kula, Hao Zhou, Zachary Garrett, Sushant Kafle, Sercan Arik, Ankita Goel, Mingyao Yang, Jiho Park, Koji Kojima, Parsa Mahmoudieh, Koray Kayukcuoglu, Grace Chen, Doug Fritz, Anton Bulyenov, Sudeshna Roy, Dimitris Paparas, Hadar Shemtov, Bo-Juen Chen, Robin Strudel, David Reitter, Aurko Roy, Andrey Vlasov, Changwan Ryu, Chas Leichner, Haichuan Yang, Zelda Mariet, Denis Vnukov, Tim Sohn, Amy Stuart, Wei Liang, Minmin Chen, Praynaa Rawlani, Christy Koh, JD Co-Reyes, Guangda Lai, Praseem Banzal, Dimitrios Vytiniotis, Jieru Mei, Mu Cai, Mohammed Badawi, Corey Fry, Ale Hartman, Daniel Zheng, Eric Jia, James Keeling, Annie Louis, Ying Chen, Efren Robles, Wei-Chih Hung, Howard Zhou, Nikita Saxena, Sonam Goenka, Olivia Ma, Zach Fisher, Mor Hazan Taege, Emily Graves, David Steiner, Yujia Li, Sarah Nguyen, Rahul Sukthankar, Joe Stanton, Ali Eslami, Gloria Shen, Berkin Akin, Alexey Guseynov, Yiqian Zhou, Jean-Baptiste Alayrac, Armand Joulin, Efrat Farkash, Ashish Thapliyal, Stephen Roller, Noam Shazeer, Todor Davchev, Terry Koo, Hannah Forbes-Pollard, Kartik Audhkhasi, Greg Farquhar, Adi Mayrav Gilady, Maggie Song, John Aslanides, Piermaria Mendolicchio, Alicia Parrish, John Blitzer, Pramod Gupta, Xiaoen Ju, Xiaochen Yang, Puranjay Datta, Andrea Tacchetti, Sanket Vaibhay Mehta, Gregory Dibb, Shubham Gupta, Federico Piccinini, Raia Hadsell, Sujee Rajayogam, Jiepu Jiang, Patrick Griffin, Patrik Sundberg, Jamie Hayes, Alexey Frolov, Tian Xie, Adam Zhang, Kingshuk Dasgupta, Uday Kalra, Lior Shani, Klaus Macherey, Tzu-Kuo Huang, Liam MacDermed, Karthik Duddu, Paulo Zacchello, Zi Yang, Jessica Lo, Kai Hui, Matej Kastelic, Derek Gasaway, Qijun Tan, Summer Yue, Pablo Barrio, John Wieting, Weel Yang, Andrew Nystrom, Solomon Demmessie, Anselm Levskaya, Fabio Viola, Chetan Tekur, Greg Billock, George Necula, Mandar Joshi, Rylan Schaeffer, Swachhand Lokhande, Christina Sorokin, Pradeep Shenoy, Mia Chen, Mark Collier, Hongji Li, Taylor Bos, Nevan Wichers, Sun Jae Lee, Angéline Pouget, Santhosh Thangaraj, Kyriakos Axiotis, Phil Crone, Rachel Sterneck, Nikolai Chinaev, Victoria Krakovna, Oleksandr Ferludin, Ian Gemp, Stephanie Winkler, Dan Goldberg, Ivan Korotkov, Kefan Xiao, Malika Mehrotra, Sandeep Mariserla, Vihari Piratla, Terry Thurk, Khiem Pham, Hongxu Ma, Alexandre Senges, Ravi Kumar, Clemens Meyer, Ellie Talius, Nuo Wang Pierse, Ballie Sandhu, Horia Toma, Kuo Lin, Swaroop Nath, Tom Stone, Dorsa Sadigh, Nikita Gupta, Arthur Guez, Avi Singh, Matt Thomas, Tom Duerig, Yuan Gong, Richard Tanburn, Lydia Lihui Zhang, Phuong Dao, Mohamed Hammad, Sirui Xie, Shruti Rijhwani, Ben Murdoch, Duhyeon Kim, Will Thompson, Heng-Tze Cheng, Daniel Sohn, Pablo Sprechmann, Qiantong Xu, Srinivas Tadepalli, Peter Young, Ye Zhang, Hansa Srinivasan, Miranda Aperghis, Aditya Ayyar, Hen Fitoussi, Ryan Burnell, David Madras, Mike Dusenberry, Xi Xiong, Tayo Oguntebi, Ben Albrecht, Jörg Bornschein, Jovana Mitrović, Mason Dimarco, Bhargav Kanagal Shamanna, Premal Shah, Eren Sezener, Shyam Upadhyay, Dave Lacey, Craig Schiff, Sebastien Baur, Sanjay Ganapathy, Eva Schnider, Mateo Wirth, Connor Schenck, Andrey Simanovsky, Yi-Xuan Tan, Philipp Fränken, Dennis Duan, Bharath Mankalale, Nikhil Dhawan, Kevin Sequeira, Zichuan Wei, Shivanker Goel, Caglar Unlu, Yukun Zhu, Haitian Sun, Ananth Balashankar, Kurt Shuster, Megh Umekar, Mahmoud Alnahlawi, Aäron van den Oord, Kelly Chen, Yuexiang Zhai, Zihang Dai, Kuang-Huei Lee, Eric Doi, Lukas Zilka, Rohith Vallu, Disha Shrivastava, Jason Lee, Hisham Husain, Honglei Zhuang, Vincent Cohen-Addad, Jarred Barber, James Atwood, Adam Sadovsky, Quentin Wellens, Steven Hand, Arunkumar Rajendran, Aybuke Turker, CJ Carey, Yuanzhong Xu, Hagen Soltau, Zefei Li, Xinying Song, Conglong Li, Iurii Kemaev, Sasha Brown, Andrea Burns, Viorica Patraucean, Piotr Stanczyk, Renga Aravamudhan, Mathieu Blondel, Hila Noga, Lorenzo Blanco, Will Song, Michael Isard, Mandar

595

596

597

598

600

601

602

603

604

605

607

608

609

610

612

613

614

615

616

617

618

619

620

621

622

623

625

626

627

629

630

631

632

633

634

635

636

637

638

639

640

641

642

644

645

646

Sharma, Reid Hayes, Dalia El Badawy, Avery Lamp, Itay Laish, Olga Kozlova, Kelvin Chan, Sahil Singla, Srinivas Sunkara, Mayank Upadhyay, Chang Liu, Aijun Bai, Jarek Wilkiewicz, Martin Zlocha, Jeremiah Liu, Zhuowan Li, Haiguang Li, Omer Barak, Ganna Raboshchuk, Jiho Choi, Fangyu Liu, Erik Jue, Mohit Sharma, Andreea Marzoca, Robert Busa-Fekete, Anna Korsun, Andre Elisseeff, Zhe Shen, Sara Mc Carthy, Kay Lamerigts, Anahita Hosseini, Hanzhao Lin, Charlie Chen, Fan Yang, Kushal Chauhan, Mark Omernick, Dawei Jia, Karina Zainullina, Demis Hassabis, Danny Vainstein, Ehsan Amid, Xiang Zhou, Ronny Votel, Eszter Vértes, Xinjian Li, Zongwei Zhou, Angeliki Lazaridou, Brendan McMahan, Arjun Narayanan, Hubert Soyer, Sujoy Basu, Kayi Lee, Bryan Perozzi, Qin Cao, Leonard Berrada, Rahul Arya, Ke Chen, Katrina, Xu, Matthias Lochbrunner, Alex Hofer, Sahand Sharifzadeh, Renjie Wu, Sally Goldman, Pranjal Awasthi, Xuezhi Wang, Yan Wu, Claire Sha, Biao Zhang, Maciej Mikuła, Filippo Graziano, Siobhan McIoughlin, Irene Giannoumis, Youhei Namiki, Chase Malik, Carey Radebaugh, Jamie Hall, Ramiro Leal-Cavazos, Jianmin Chen, Vikas Sindhwani, David Kao, David Greene, Jordan Griffith, Chris Welty, Ceslee Montgomery, Toshihiro Yoshino, Liangzhe Yuan, Noah Goodman, Assaf Hurwitz Michaely, Kevin Lee, KP Sawhney, Wei Chen, Zheng Zheng, Megan Shum, Nikolay Savinov, Etienne Pot, Alex Pak, Morteza Zadimoghaddam, Sijal Bhatnagar, Yoad Lewenberg, Blair Kutzman, Ji Liu, Lesley Katzen, Jeremy Selier, Josip Djolonga, Dmitry Lepikhin, Kelvin Xu, Jacky Liang, Jiewen Tan, Benoit Schillings, Muge Ersoy, Pete Blois, Bernd Bandemer, Abhimanyu Singh, Sergei Lebedev, Pankaj Joshi, Adam R. Brown, Evan Palmer, Shreya Pathak, Komal Jalan, Fedir Zubach, Shuba Lall, Randall Parker, Alok Gunjan, Sergey Rogulenko, Sumit Sanghai, Zhaoqi Leng, Zoltan Egyed, Shixin Li, Maria Ivanova, Kostas Andriopoulos, Jin Xie, Elan Rosenfeld, Auriel Wright, Ankur Sharma, Xinyang Geng, Yicheng Wang, Sam Kwei, Renke Pan, Yujing Zhang, Gabby Wang, Xi Liu, Chak Yeung, Elizabeth Cole, Aviv Rosenberg, Zhen Yang, Phil Chen, George Polovets, Pranav Nair, Rohun Saxena, Josh Smith, Shuo yiin Chang, Aroma Mahendru, Svetlana Grant, Anand Iyer, Irene Cai, Jed McGiffin, Jiaming Shen, Alanna Walton, Antonious Girgis, Oliver Woodman, Rosemary Ke, Mike Kwong, Louis Rouillard, Jinmeng Rao, Zhihao Li, Yuntao Xu, Flavien Prost, Chi Zou, Ziwei Ji, Alberto Magni, Tyler Liechty, Dan A. Calian, Deepak Ramachandran, Igor Krivokon, Hui Huang, Terry Chen, Anja Hauth, Anastasija Ilić, Weijuan Xi, Hyeontaek Lim, Vlad-Doru Ion, Pooya Moradi, Metin Toksoz-Exley, Kalesha Bullard, Miltos Allamanis, Xiaomeng Yang, Sophie Wang, Zhi Hong, Anita Gergely, Cheng Li, Bhavishya Mittal, Vitaly Kovalev, Victor Ungureanu, Jane Labanowski, Jan Wassenberg, Nicolas Lacasse, Geoffrey Cideron, Petar Dević, Annie Marsden, Lynn Nguyen, Michael Fink, Yin Zhong, Tatsuya Kiyono, Desi Ivanov, Sally Ma, Max Bain, Kiran Yalasangi, Jennifer She, Anastasia Petrushkina, Mayank Lunayach, Carla Bromberg, Sarah Hodkinson, Vilobh Meshram, Daniel Vlasic, Austin Kyker, Steve Xu, Jeff Stanway, Zuguang Yang, Kai Zhao, Matthew Tung, Seth Odoom, Yasuhisa Fujii, Justin Gilmer, Eunyoung Kim, Felix Halim, Quoc Le, Bernd Bohnet, Seliem El-Sayed, Behnam Neyshabur, Malcolm Reynolds, Dean Reich, Yang Xu, Erica Moreira, Anuj Sharma, Zeyu Liu, Mohammad Javad Hosseini, Naina Raisinghani, Yi Su, Ni Lao, Daniel Formoso, Marco Gelmi, Almog Gueta, Tapomay Dey, Elena Gribovskaya, Domagoj Ćevid, Sidharth Mudgal, Garrett Bingham, Jianling Wang, Anurag Kumar, Alex Cullum, Feng Han, Konstantinos Bousmalis, Diego Cedillo, Grace Chu, Vladimir Magay, Paul Michel, Ester Hlavnova, Daniele Calandriello, Setareh Ariafar, Kaisheng Yao, Vikash Sehwag, Arpi Vezer, Agustin Dal Lago, Zhenkai Zhu, Paul Kishan Rubenstein, Allen Porter, Anirudh Baddepudi, Oriana Riva, Mihai Dorin Istin, Chih-Kuan Yeh, Zhi Li, Andrew Howard, Nilpa Jha, Jeremy Chen, Raoul de Liedekerke, Zafarali Ahmed, Mikel Rodriguez, Tanuj Bhatia, Bangju Wang, Ali Elqursh, David Klinghoffer, Peter Chen, Pushmeet Kohli, Te I, Weiyang Zhang, Zack Nado, Jilin Chen, Maxwell Chen, George Zhang, Aayush Singh, Adam Hillier, Federico Lebron, Yiqing Tao, Ting Liu, Gabriel Dulac-Arnold, Jingwei Zhang, Shashi Narayan, Buhuang Liu, Orhan Firat, Abhishek Bhowmick, Bingyuan Liu, Hao Zhang, Zizhao Zhang, Georges Rotival, Nathan Howard, Anu Sinha, Alexander Grushetsky, Benjamin Beyret, Keerthana Gopalakrishnan, James Zhao, Kyle He, Szabolcs Payrits, Zaid Nabulsi, Zhaoyi Zhang, Weijie Chen, Edward Lee, Nova Fallen, Sreenivas Gollapudi, Aurick Zhou, Filip Pavetić, Thomas Köppe, Shiyu Huang, Rama Pasumarthi, Nick Fernando, Felix Fischer, Daria Curko, Yang Gao, James Svensson, Austin Stone, Haroon Qureshi, Abhishek Sinha, Apoorv Kulshreshtha, Martin Matysiak, Jieming Mao, Carl Saroufim, Aleksandra Faust, Qingnan Duan, Gil Fidel, Kaan Katircioglu, Raphaël Lopez Kaufman, Dhruv Shah, Weize Kong, Abhishek Bapna, Gellért Weisz, Emma Dunleavy, Praneet Dutta, Tianqi Liu, Rahma Chaabouni, Carolina Parada, Marcus Wu, Alexandra Belias, Alessandro Bissacco, Stanislav Fort, Li Xiao, Fantine

650

651

652

653

654

655

656

657

658

659

660

661

662

665

666

667

668

669

670

671

672

673

674

675

676

677

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

696

699

700

Huot, Chris Knutsen, Yochai Blau, Gang Li, Jennifer Prendki, Juliette Love, Yinlam Chow, Pichi Charoenpanit, Hidetoshi Shimokawa, Vincent Coriou, Karol Gregor, Tomas Izo, Arjun Akula, Mario Pinto, Chris Hahn, Dominik Paulus, Jiaxian Guo, Neha Sharma, Cho-Jui Hsieh, Adaeze Chukwuka, Kazuma Hashimoto, Nathalie Rauschmayr, Ling Wu, Christof Angermueller, Yulong Wang, Sebastian Gerlach, Michael Pliskin, Daniil Mirylenka, Min Ma, Lexi Baugher, Bryan Gale, Shaan Bijwadia, Nemanja Rakićević, David Wood, Jane Park, Chung-Ching Chang, Babi Seal, Chris Tar, Kacper Krasowiak, Yiwen Song, Georgi Stephanov, Gary Wang, Marcello Maggioni, Stein Xudong Lin, Felix Wu, Shachi Paul, Zixuan Jiang, Shubham Agrawal, Bilal Piot, Alex Feng, Cheolmin Kim, Tulsee Doshi, Jonathan Lai, Chuqiao, Xu, Sharad Vikram, Ciprian Chelba, Sebastian Krause, Vincent Zhuang, Jack Rae, Timo Denk, Adrian Collister, Lotte Weerts, Xianghong Luo, Yifeng Lu, Håvard Garnes, Nitish Gupta, Terry Spitz, Avinatan Hassidim, Lihao Liang, Izhak Shafran, Peter Humphreys, Kenny Vassigh, Phil Wallis, Virat Shejwalkar, Nicolas Perez-Nieves, Rachel Hornung, Melissa Tan, Beka Westberg, Andy Ly, Richard Zhang, Brian Farris, Jongbin Park, Alec Kosik, Zeynep Cankara, Andrii Maksai, Yunhan Xu, Albin Cassirer, Sergi Caelles, Abbas Abdolmaleki, Mencher Chiang, Alex Fabrikant, Shravya Shetty, Luheng He, Mai Giménez, Hadi Hashemi, Sheena Panthaplackel, Yana Kulizhskaya, Salil Deshmukh, Daniele Pighin, Robin Alazard, Disha Jindal, Seb Noury, Pradeep Kumar S, Siyang Qin, Xerxes Dotiwalla, Stephen Spencer, Mohammad Babaeizadeh, Blake JianHang Chen, Vaibhav Mehta, Jennie Lees, Andrew Leach, Penporn Koanantakool, Ilia Akolzin, Ramona Comanescu, Junwhan Ahn, Alexey Svyatkovskiy, Basil Mustafa, David D'Ambrosio, Shiva Mohan Reddy Garlapati, Pascal Lamblin, Alekh Agarwal, Shuang Song, Pier Giuseppe Sessa, Pauline Coquinot, John Maggs, Hussain Masoom, Divya Pitta, Yaqing Wang, Patrick Morris-Suzuki, Billy Porter, Johnson Jia, Jeffrey Dudek, Raghavender R, Cosmin Paduraru, Alan Ansell, Tolga Bolukbasi, Tony Lu, Ramya Ganeshan, Zi Wang, Henry Griffiths, Rodrigo Benenson, Yifan He, James Swirhun, George Papamakarios, Aditya Chawla, Kuntal Sengupta, Yan Wang, Vedrana Milutinovic, Igor Mordatch, Zhipeng Jia, Jamie Smith, Will Ng, Shitij Nigam, Matt Young, Eugen Vušak, Blake Hechtman, Sheela Goenka, Avital Zipori, Kareem Ayoub, Ashok Popat, Trilok Acharya, Luo Yu, Dawn Bloxwich, Hugo Song, Paul Roit, Haiqiong Li, Aviel Boag, Nigamaa Nayakanti, Bilva Chandra, Tianli Ding, Aahil Mehta, Cath Hope, Jiageng Zhang, Idan Heimlich Shtacher, Kartikeya Badola, Ryo Nakashima, Andrei Sozanschi, Iulia Comşa, Ante Žužul, Emily Caveness, Julian Odell, Matthew Watson, Dario de Cesare, Phillip Lippe, Derek Lockhart, Siddharth Verma, Huizhong Chen, Sean Sun, Lin Zhuo, Aditya Shah, Prakhar Gupta, Alex Muzio, Ning Niu, Amir Zait, Abhinav Singh, Meenu Gaba, Fan Ye, Prajit Ramachandran, Mohammad Saleh, Raluca Ada Popa, Ayush Dubey, Frederick Liu, Sara Javanmardi, Mark Epstein, Ross Hemsley, Richard Green, Nishant Ranka, Eden Cohen, Chuyuan Kelly Fu, Sanjay Ghemawat, Jed Borovik, James Martens, Anthony Chen, Pranav Shyam, André Susano Pinto, Ming-Hsuan Yang, Alexandru Tifrea, David Du, Boqing Gong, Ayushi Agarwal, Seungyeon Kim, Christian Frank, Saloni Shah, Xiaodan Song, Zhiwei Deng, Ales Mikhalap, Kleopatra Chatziprimou, Timothy Chung, Toni Creswell, Susan Zhang, Yennie Jun, Carl Lebsack, Will Truong, Slavica Andačić, Itay Yona, Marco Fornoni, Rong Rong, Serge Toropov, Afzal Shama Soudagar, Andrew Audibert, Salah Zaiem, Zaheer Abbas, Andrei Rusu, Sahitya Potluri, Shitao Weng, Anastasios Kementsietsidis, Anton Tsitsulin, Daiyi Peng, Natalie Ha, Sanil Jain, Tejasi Latkar, Simeon Ivanov, Cory McLean, Anirudh GP, Rajesh Venkataraman, Canoee Liu, Dilip Krishnan, Joel D'sa, Roey Yogev, Paul Collins, Benjamin Lee, Lewis Ho, Carl Doersch, Gal Yona, Shawn Gao, Felipe Tiengo Ferreira, Adnan Ozturel, Hannah Muckenhirn, Ce Zheng, Gargi Balasubramaniam, Mudit Bansal, George van den Driessche, Sivan Eiger, Salem Haykal, Vedant Misra, Abhimanyu Goyal, Danilo Martins, Gary Leung, Jonas Valfridsson, Four Flynn, Will Bishop, Chenxi Pang, Yoni Halpern, Honglin Yu, Lawrence Moore, Yuvein, Zhu, Sridhar Thiagarajan, Yoel Drori, Zhisheng Xiao, Lucio Dery, Rolf Jagerman, Jing Lu, Eric Ge, Vaibhav Aggarwal, Arjun Khare, Vinh Tran, Oded Elyada, Ferran Alet, James Rubin, Ian Chou, David Tian, Libin Bai, Lawrence Chan, Lukasz Lew, Karolis Misiunas, Taylan Bilal, Aniket Ray, Sindhu Raghuram, Alex Castro-Ros, Viral Carpenter, CJ Zheng, Michael Kilgore, Josef Broder, Emily Xue, Praveen Kallakuri, Dheeru Dua, Nancy Yuen, Steve Chien, John Schultz, Saurabh Agrawal, Reut Tsarfaty, Jingcao Hu, Ajay Kannan, Dror Marcus, Nisarg Kothari, Baochen Sun, Ben Horn, Matko Bošnjak, Ferjad Naeem, Dean Hirsch, Lewis Chiang, Boya Fang, Jie Han, Qifei Wang, Ben Hora, Antoine He, Mario Lučić, Beer Changpinyo, Anshuman Tripathi, John Youssef, Chester Kwak, Philippe Schlattner, Cat Graves, Rémi Leblond, Wenjun Zeng, Anders Andreassen, Gabriel Rasskin, Yue Song, Eddie Cao, Junhyuk Oh, Matt Hoffman, Wojtek Skut, Yichi Zhang, Jon Stritar, Xingyu Cai, Saarthak Khanna, Kathie Wang, Shriya Sharma, Christian Reisswig, Younghoon Jun, Aman Prasad, Ta-

704

705

706

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

739

740

741

742

743

744

745

746

747

748

749

750

751

752

754

755

tiana Sholokhova, Preeti Singh, Adi Gerzi Rosenthal, Anian Ruoss, Françoise Beaufays, Sean Kirmani, Dongkai Chen, Johan Schalkwyk, Jonathan Herzig, Been Kim, Josh Jacob, Damien Vincent, Adrian N Reyes, Ivana Balazevic, Léonard Hussenot, Jon Schneider, Parker Barnes, Luis Castro, Spandana Raj Babbula, Simon Green, Serkan Cabi, Nico Duduta, Danny Driess, Rich Galt, Noam Velan, Junjie Wang, Hongyang Jiao, Matthew Mauger, Du Phan, Miteyan Patel, Vlado Galić, Jerry Chang, Eyal Marcus, Matt Harvey, Julian Salazar, Elahe Dabir, Suraj Satishkumar Sheth, Amol Mandhane, Hanie Sedghi, Jeremiah Willcock, Amir Zandieh, Shruthi Prabhakara, Aida Amini, Antoine Miech, Victor Stone, Massimo Nicosia, Paul Niemczyk, Ying Xiao, Lucy Kim, Sławek Kwasiborski, Vikas Verma, Ada Maksutaj Oflazer, Christoph Hirnschall, Peter Sung, Lu Liu, Richard Everett, Michiel Bakker, Ágoston Weisz, Yufei Wang, Vivek Sampathkumar, Uri Shaham, Bibo Xu, Yasemin Altun, Mingqiu Wang, Takaaki Saeki, Guanjie Chen, Emanuel Taropa, Shanthal Vasanth, Sophia Austin, Lu Huang, Goran Petrovic, Qingyun Dou, Daniel Golovin, Grigory Rozhdestvenskiy, Allie Culp, Will Wu, Motoki Sano, Divya Jain, Julia Proskurnia, Sébastien Cevey, Alejandro Cruzado Ruiz, Piyush Patil, Mahdi Mirzazadeh, Eric Ni, Javier Snaider, Lijie Fan, Alexandre Fréchette, AJ Pierigiovanni, Shariq Iqbal, Kenton Lee, Claudio Fantacci, Jinwei Xing, Lisa Wang, Alex Irpan, David Raposo, Yi Luan, Zhuoyuan Chen, Harish Ganapathy, Kevin Hui, Jiazhong Nie, Isabelle Guyon, Heming Ge, Roopali Vij, Hui Zheng, Dayeong Lee, Alfonso Castaño, Khuslen Baatarsukh, Gabriel Ibagon, Alexandra Chronopoulou, Nicholas FitzGerald, Shashank Viswanadha, Safeen Huda, Rivka Moroshko, Georgi Stoyanov, Prateek Kolhar, Alain Vaucher, Ishaan Watts, Adhi Kuncoro, Henryk Michalewski, Satish Kambala, Bat-Orgil Batsaikhan, Alek Andreev, Irina Jurenka, Maigo Le, Qihang Chen, Wael Al Jishi, Sarah Chakera, Zhe Chen, Aditya Kini, Vikas Yadav, Aditya Siddhant, Ilia Labzovsky, Balaji Lakshminarayanan, Carrie Grimes Bostock, Pankil Botadra, Ankesh Anand, Colton Bishop, Sam Conway-Rahman, Mohit Agarwal, Yani Donchev, Achintya Singhal, Félix de Chaumont Quitry, Natalia Ponomareva, Nishant Agrawal, Bin Ni, Kalpesh Krishna, Masha Samsikova, John Karro, Yilun Du, Tamara von Glehn, Caden Lu, Christopher A. Choquette-Choo, Zhen Qin, Tingnan Zhang, Sicheng Li, Divya Tyam, Swaroop Mishra, Wing Lowe, Colin Ji, Weiyi Wang, Manaal Faruqui, Ambrose Slone, Valentin Dalibard, Arunachalam Narayanaswamy, John Lambert, Pierre-Antoine Manzagol, Dan Karliner, Andrew Bolt, Ivan Lobov, Aditya Kusupati, Chang Ye, Xuan Yang, Heiga Zen, Nelson George, Mukul Bhutani, Olivier Lacombe, Robert Riachi, Gagan Bansal, Rachel Soh, Yue Gao, Yang Yu, Adams Yu, Emily Nottage, Tania Rojas-Esponda, James Noraky, Manish Gupta, Ragha Kotikalapudi, Jichuan Chang, Sanja Deur, Dan Graur, Alex Mossin, Erin Farnese, Ricardo Figueira, Alexandre Moufarek, Austin Huang, Patrik Zochbauer, Ben Ingram, Tongzhou Chen, Zelin Wu, Adrià Puigdomènech, Leland Rechis, Da Yu, Sri Gayatri Sundara Padmanabhan, Rui Zhu, Chu ling Ko, Andrea Banino, Samira Daruki, Aarush Selvan, Dhruva Bhaswar, Daniel Hernandez Diaz, Chen Su, Salvatore Scellato, Jennifer Brennan, Woohyun Han, Grace Chung, Priyanka Agrawal, Urvashi Khandelwal, Khe Chai Sim, Morgane Lustman, Sam Ritter, Kelvin Guu, Jiawei Xia, Prateek Jain, Emma Wang, Tyrone Hill, Mirko Rossini, Marija Kostelac, Tautvydas Misiunas, Amit Sabne, Kyuyeun Kim, Ahmet Iscen, Congchao Wang, José Leal, Ashwin Sreevatsa, Utku Evci, Manfred Warmuth, Saket Joshi, Daniel Suo, James Lottes, Garrett Honke, Brendan Jou, Stefani Karp, Jieru Hu, Himanshu Sahni, Adrien Ali Taïga, William Kong, Samrat Ghosh, Renshen Wang, Jay Pavagadhi, Natalie Axelsson, Nikolai Grigorev, Patrick Siegler, Rebecca Lin, Guohui Wang, Emilio Parisotto, Sharath Maddineni, Krishan Subudhi, Eyal Ben-David, Elena Pochernina, Orgad Keller, Thi Avrahami, Zhe Yuan, Pulkit Mehta, Jialu Liu, Sherry Yang, Wendy Kan, Katherine Lee, Tom Funkhouser, Derek Cheng, Hongzhi Shi, Archit Sharma, Joe Kelley, Matan Eyal, Yury Malkov, Corentin Tallec, Yuval Bahat, Shen Yan, Xintian, Wu, David Lindner, Chengda Wu, Avi Caciularu, Xiyang Luo, Rodolphe Jenatton, Tim Zaman, Yingying Bi, Ilya Kornakov, Ganesh Mallya, Daisuke Ikeda, Itay Karo, Anima Singh, Colin Evans, Praneeth Netrapalli, Vincent Nallatamby, Isaac Tian, Yannis Assael, Vikas Raunak, Victor Carbune, Ioana Bica, Lior Madmoni, Dee Cattle, Snchit Grover, Krishna Somandepalli, Sid Lall, Amelio Vázquez-Reina, Riccardo Patana, Jiaqi Mu, Pranav Talluri, Maggie Tran, Rajeev Aggarwal, RJ Skerry-Ryan, Jun Xu, Mike Burrows, Xiaoyue Pan, Edouard Yvinec, Di Lu, Zhiying Zhang, Duc Dung Nguyen, Hairong Mu, Gabriel Barcik, Helen Ran, Lauren Beltrone, Krzysztof Choromanski, Dia Kharrat, Samuel Albanie, Sean Purser-haskell, David Bieber, Carrie Zhang, Jing Wang, Tom Hudson, Zhiyuan Zhang, Han Fu, Johannes Mauerer, Mohammad Hossein Bateni, AJ Maschinot, Bing Wang, Muye Zhu, Arjun Pillai, Tobias Weyand, Shuang Liu, Oscar Akerlund, Fred Bertsch, Vittal Premachandran, Alicia Jin, Vincent Roulet, Peter de Boursac, Shubham Mittal, Ndaba Ndebele, Georgi Karadzhov, Sahra Ghalebikesabi, Ricky Liang, Allen Wu, Yale Cong, Nimesh Ghelani, Sumeet Singh, Ba-

758

759

760

761

762

764

765

766

767

768

769

770

771

774

775

776

777

778

780

781

782

783

784

785

786

787

788

789

790

793

794

796

798

799

800

801

802

804

har Fatemi, Warren, Chen, Charles Kwong, Alexey Kolganov, Steve Li, Richard Song, Chenkai Kuang, Sobhan Miryoosefi, Dale Webster, James Wendt, Arkadiusz Socala, Guolong Su, Artur Mendonça, Abhinav Gupta, Xiaowei Li, Tomy Tsai, Qiong, Hu, Kai Kang, Angie Chen, Sertan Girgin, Yongqin Xian, Andrew Lee, Nolan Ramsden, Leslie Baker, Madeleine Clare Elish, Varvara Krayvanova, Rishabh Joshi, Jiri Simsa, Yao-Yuan Yang, Piotr Ambroszczyk, Dipankar Ghosh, Arjun Kar, Yuan Shangguan, Yumeya Yamamori, Yaroslav Akulov, Andy Brock, Haotian Tang, Siddharth Vashishtha, Rich Munoz, Andreas Steiner, Kalyan Andra, Daniel Eppens, Qixuan Feng, Hayato Kobayashi, Sasha Goldshtein, Mona El Mahdy, Xin Wang, Jilei, Wang, Richard Killam, Tom Kwiatkowski, Kavya Kopparapu, Serena Zhan, Chao Jia, Alexei Bendebury, Sheryl Luo, Adrià Recasens, Timothy Knight, Jing Chen, Mohak Patel, YaGuang Li, Ben Withbroe, Dean Weesner, Kush Bhatia, Jie Ren, Danielle Eisenbud, Ebrahim Songhori, Yanhua Sun, Travis Choma, Tasos Kementsietsidis, Lucas Manning, Brian Roark, Wael Farhan, Jie Feng, Susheel Tatineni, James Cobon-Kerr, Yunjie Li, Lisa Anne Hendricks, Isaac Noble, Chris Breaux, Nate Kushman, Liqian Peng, Fuzhao Xue, Taylor Tobin, Jamie Rogers, Josh Lipschultz, Chris Alberti, Alexey Vlaskin, Mostafa Dehghani, Roshan Sharma, Tris Warkentin, Chen-Yu Lee, Benigno Uria, Da-Cheng Juan, Angad Chandorkar, Hila Sheftel, Ruibo Liu, Elnaz Davoodi, Borja De Balle Pigem, Kedar Dhamdhere, David Ross, Jonathan Hoech, Mahdis Mahdieh, Li Liu, Qiujia Li, Liam McCafferty, Chenxi Liu, Markus Mircea, Yunting Song, Omkar Savant, Alaa Saade, Colin Cherry, Vincent Hellendoorn, Siddharth Goyal, Paul Pucciarelli, David Vilar Torres, Zohar Yahav, Hyo Lee, Lars Lowe Sjoesund, Christo Kirov, Bo Chang, Deepanway Ghoshal, Lu Li, Gilles Baechler, Sébastien Pereira, Tara Sainath, Anudhyan Boral, Dominik Grewe, Afief Halumi, Nguyet Minh Phu, Tianxiao Shen, Marco Tulio Ribeiro, Dhriti Varma, Alex Kaskasoli, Vlad Feinberg, Navneet Potti, Jarrod Kahn, Matheus Wisniewski, Shakir Mohamed, Arnar Mar Hrafnkelsson, Bobak Shahriari, Jean-Baptiste Lespiau, Lisa Patel, Legg Yeung, Tom Paine, Lantao Mei, Alex Ramirez, Rakesh Shivanna, Li Zhong, Josh Woodward, Guilherme Tubone, Samira Khan, Heng Chen, Elizabeth Nielsen, Catalin Ionescu, Utsav Prabhu, Mingcen Gao, Qingze Wang, Sean Augenstein, Neesha Subramaniam, Jason Chang, Fotis Iliopoulos, Jiaming Luo, Myriam Khan, Weicheng Kuo, Denis Teplyashin, Florence Perot, Logan Kilpatrick, Amir Globerson, Hongkun Yu, Anfal Siddiqui, Nick Sukhanov, Arun Kandoor, Umang Gupta, Marco Andreetto, Moran Ambar, Donnie Kim, Paweł Wesołowski, Sarah Perrin, Ben Limonchik, Wei Fan, Jim Stephan, Ian Stewart-Binks, Ryan Kappedal, Tong He, Sarah Cogan, Romina Datta, Tong Zhou, Jiayu Ye, Leandro Kieliger, Ana Ramalho, Kyle Kastner, Fabian Mentzer, Wei-Jen Ko, Arun Suggala, Tianhao Zhou, Shiraz Butt, Hana Strejček, Lior Belenki, Subhashini Venugopalan, Mingyang Ling, Evgenii Eltyshev, Yunxiao Deng, Geza Kovacs, Mukund Raghavachari, Hanjun Dai, Tal Schuster, Steven Schwarcz, Richard Nguyen, Arthur Nguyen, Gavin Buttimore, Shrestha Basu Mallick, Sudeep Gandhe, Seth Benjamin, Michal Jastrzebski, Le Yan, Sugato Basu, Chris Apps, Isabel Edkins, James Allingham, Immanuel Odisho, Tomas Kocisky, Jewel Zhao, Linting Xue, Apoorv Reddy, Chrysovalantis Anastasiou, Aviel Atias, Sam Redmond, Kieran Milan, Nicolas Heess, Herman Schmit, Allan Dafoe, Daniel Andor, Tynan Gangwani, Anca Dragan, Sheng Zhang, Ashyana Kachra, Gang Wu, Siyang Xue, Kevin Aydin, Siqi Liu, Yuxiang Zhou, Mahan Malihi, Austin Wu, Siddharth Gopal, Candice Schumann, Peter Stys, Alek Wang, Mirek Olšák, Dangyi Liu, Christian Schallhart, Yiran Mao, Demetra Brady, Hao Xu, Tomas Mery, Chawin Sitawarin, Siva Velusamy, Tom Cobley, Alex Zhai, Christian Walder, Nitzan Katz, Ganesh Jawahar, Chinmay Kulkarni, Antoine Yang, Adam Paszke, Yinan Wang, Bogdan Damoc, Zalán Borsos, Ray Smith, Jinning Li, Mansi Gupta, Andrei Kapishnikov, Sushant Prakash, Florian Luisier, Rishabh Agarwal, Will Grathwohl, Kuangyuan Chen, Kehang Han, Nikhil Mehta, Andrew Over, Shekoofeh Azizi, Lei Meng, Niccolò Dal Santo, Kelvin Zheng, Jane Shapiro, Igor Petrovski, Jeffrey Hui, Amin Ghafouri, Jasper Snoek, James Qin, Mandy Jordan, Caitlin Sikora, Jonathan Malmaud, Yuheng Kuang, Aga Świetlik, Ruoxin Sang, Chongyang Shi, Leon Li, Andrew Rosenberg, Shubin Zhao, Andy Crawford, Jan-Thorsten Peter, Yun Lei, Xavier Garcia, Long Le, Todd Wang, Julien Amelot, Dave Orr, Praneeth Kacham, Dana Alon, Gladys Tyen, Abhinav Arora, James Lyon, Alex Kurakin, Mimi Ly, Theo Guidroz, Zhipeng Yan, Rina Panigrahy, Pingmei Xu, Thais Kagohara, Yong Cheng, Eric Noland, Jinhyuk Lee, Jonathan Lee, Cathy Yip, Maria Wang, Efrat Nehoran, Alexander Bykovsky, Zhihao Shan, Ankit Bhagatwala, Chaochao Yan, Jie Tan, Guillermo Garrido, Dan Ethier, Nate Hurley, Grace Vesom, Xu Chen, Siyuan Qiao, Abhishek Nayyar, Julian Walker, Paramjit Sandhu, Mihaela Rosca, Danny Swisher, Mikhail Dektiarev, Josh Dillon, George-Cristian Muraru, Manuel Tragut, Artiom Myaskovsky, David Reid, Marko Velic, Owen Xiao, Jasmine George, Mark Brand, Jing Li, Wenhao Yu, Shane Gu, Xiang Deng, François-Xavier Aubet, Soheil Hassas Yeganeh, Fred Alcober, Celine Smith, Trevor Cohn,

811

812

813

814

815

816

817

818

819

820

821

822

823

824

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

858

861

862

Kay McKinney, Michael Tschannen, Ramesh Sampath, Gowoon Cheon, Liangchen Luo, Luyang Liu, Jordi Orbay, Hui Peng, Gabriela Botea, Xiaofan Zhang, Charles Yoon, Cesar Magalhaes, Paweł Stradomski, Ian Mackinnon, Steven Hemingray, Kumaran Venkatesan, Rhys May, Jaeyoun Kim, Alex Druinsky, Jingchen Ye, Zheng Xu, Terry Huang, Jad Al Abdallah, Adil Dostmohamed, Rachana Fellinger, Tsendsuren Munkhdalai, Akanksha Maurya, Peter Garst, Yin Zhang, Maxim Krikun, Simon Bucher, Aditya Srikanth Veerubhotla, Yaxin Liu, Sheng Li, Nishesh Gupta, Jakub Adamek, Hanwen Chen, Bernett Orlando, Aleksandr Zaks, Joost van Amersfoort, Josh Camp, Hui Wan, HyunJeong Choe, Zhichun Wu, Kate Olszewska, Weiren Yu, Archita Vadali, Martin Scholz, Daniel De Freitas, Jason Lin, Amy Hua, Xin Liu, Frank Ding, Yichao Zhou, Boone Severson, Katerina Tsihlas, Samuel Yang, Tammo Spalink, Varun Yerram, Helena Pankov, Rory Blevins, Ben Vargas, Sarthak Jauhari, Matt Miecnikowski, Ming Zhang, Sandeep Kumar, Clement Farabet, Charline Le Lan, Sebastian Flennerhag, Yonatan Bitton, Ada Ma, Arthur Bražinskas, Eli Collins, Niharika Ahuja, Sneha Kudugunta, Anna Bortsova, Minh Giang, Wanzheng Zhu, Ed Chi, Scott Lundberg, Alexey Stern, Subha Puttagunta, Jing Xiong, Xiao Wu, Yash Pande, Amit Jhindal, Daniel Murphy, Jon Clark, Marc Brockschmidt, Maxine Deines, Kevin R. McKee, Dan Bahir, Jiajun Shen, Minh Truong, Daniel McDuff, Andrea Gesmundo, Edouard Rosseel, Bowen Liang, Ken Caluwaerts, Jessica Hamrick, Joseph Kready, Mary Cassin, Rishikesh Ingale, Li Lao, Scott Pollom, Yifan Ding, Wei He, Lizzetth Bellot, Joana Iljazi, Ramya Sree Boppana, Shan Han, Tara Thompson, Amr Khalifa, Anna Bulanova, Blagoj Mitrevski, Bo Pang, Emma Cooney, Tian Shi, Rey Coaguila, Tamar Yakar, Marc'aurelio Ranzato, Nikola Momchev, Chris Rawles, Zachary Charles, Young Maeng, Yuan Zhang, Rishabh Bansal, Xiaokai Zhao, Brian Albert, Yuan Yuan, Sudheendra Vijayanarasimhan, Roy Hirsch, Vinay Ramasesh, Kiran Vodrahalli, Xingyu Wang, Arushi Gupta, DJ Strouse, Jianmo Ni, Roma Patel, Gabe Taubman, Zhouyuan Huo, Dero Gharibian, Marianne Monteiro, Hoi Lam, Shobha Vasudevan, Aditi Chaudhary, Isabela Albuquerque, Kilol Gupta, Sebastian Riedel, Chaitra Hegde, Avraham Ruderman, András György, Marcus Wainwright, Ashwin Chaugule, Burcu Karagol Ayan, Tomer Levinboim, Sam Shleifer, Yogesh Kalley, Vahab Mirrokni, Abhishek Rao, Prabakar Radhakrishnan, Jay Hartford, Jialin Wu, Zhenhai Zhu, Francesco Bertolini, Hao Xiong, Nicolas Serrano, Hamish Tomlinson, Myle Ott, Yifan Chang, Mark Graham, Jian Li, Marco Liang, Xiangzhu Long, Sebastian Borgeaud, Yanif Ahmad, Alex Grills, Diana Mincu, Martin Izzard, Yuan Liu, Jinyu Xie, Louis O'Bryan, Sameera Ponda, Simon Tong, Michelle Liu, Dan Malkin, Khalid Salama, Yuankai Chen, Rohan Anil, Anand Rao, Rigel Swavely, Misha Bilenko, Nina Anderson, Tat Tan, Jing Xie, Xing Wu, Lijun Yu, Oriol Vinyals, Andrey Ryabtsev, Rumen Dangovski, Kate Baumli, Daniel Keysers, Christian Wright, Zoe Ashwood, Betty Chan, Artem Shtefan, Yaohui Guo, Ankur Bapna, Radu Soricut, Steven Pecht, Sabela Ramos, Rui Wang, Jiahao Cai, Trieu Trinh, Paul Barham, Linda Friso, Eli Stickgold, Xiangzhuo Ding, Siamak Shakeri, Diego Ardila, Eleftheria Briakou, Phil Culliton, Adam Raveret, Jingyu Cui, David Saxton, Subhrajit Roy, Javad Azizi, Pengcheng Yin, Lucia Loher, Andrew Bunner, Min Choi, Faruk Ahmed, Eric Li, Yin Li, Shengyang Dai, Michael Elabd, Sriram Ganapathy, Shivani Agrawal, Yiqing Hua, Paige Kunkle, Sujeevan Rajayogam, Arun Ahuja, Arthur Conmy, Alex Vasiloff, Parker Beak, Christopher Yew, Jayaram Mudigonda, Bartek Wydrowski, Jon Blanton, Zhengdong Wang, Yann Dauphin, Zhuo Xu, Martin Polacek, Xi Chen, Hexiang Hu, Pauline Sho, Markus Kunesch, Mehdi Hafezi Manshadi, Eliza Rutherford, Bo Li, Sissie Hsiao, Iain Barr, Alex Tudor, Matija Kecman, Arsha Nagrani, Vladimir Pchelin, Martin Sundermeyer, Aishwarya P S, Abhijit Karmarkar, Yi Gao, Grishma Chole, Olivier Bachem, Isabel Gao, Arturo BC, Matt Dibb, Mauro Verzetti, Felix Hernandez-Campos, Yana Lunts, Matthew Johnson, Julia Di Trapani, Raphael Koster, Idan Brusilovsky, Binbin Xiong, Megha Mohabey, Han Ke, Joe Zou, Tea Sabolić, Víctor Campos, John Palowitch, Alex Morris, Linhai Qiu, Pranavaraj Ponnuramu, Fangtao Li, Vivek Sharma, Kiranbir Sodhia, Kaan Tekelioglu, Aleksandr Chuklin, Madhavi Yenugula, Erika Gemzer, Theofilos Strinopoulos, Sam El-Husseini, Huiyu Wang, Yan Zhong, Edouard Leurent, Paul Natsev, Weijun Wang, Dre Mahaarachchi, Tao Zhu, Songyou Peng, Sami Alabed, Cheng-Chun Lee, Anthony Brohan, Arthur Szlam, GS Oh, Anton Kovsharov, Jenny Lee, Renee Wong, Megan Barnes, Gregory Thornton, Felix Gimeno, Omer Levy, Martin Sevenich, Melvin Johnson, Jonathan Mallinson, Robert Dadashi, Ziyue Wang, Qingchun Ren, Preethi Lahoti, Arka Dhar, Josh Feldman, Dan Zheng, Thatcher Ulrich, Liviu Panait, Michiel Blokzijl, Cip Baetu, Josip Matak, Jitendra Harlalka, Maulik Shah, Tal Marian, Daniel von Dincklage, Cosmo Du, Ruy Ley-Wild, Bethanie Brownfield, Max Schumacher, Yury Stuken, Shadi Noghabi, Sonal Gupta, Xiaoqi Ren, Eric Malmi, Felix Weissenberger, Blanca Huergo, Maria Bauza, Thomas Lampe, Arthur Douillard, Mojtaba Seyedhosseini, Roy Frostig, Zoubin Ghahramani, Kelvin Nguyen, Kashyap Krishnakumar, Chengxi Ye, Rahul Gupta, Alireza Nazari, Robert Geirhos, Pete Shaw, Ahmed

865

866

867

868

870

871

872

873

874

875

876

877

878

879

880

883

885

889

890

891

892

893

894

895

897

899

900

901

902

903 904

905

906

907

908

909

910

911

912

913

914

915

916

917

Eleryan, Dima Damen, Jennimaria Palomaki, Ted Xiao, Qiyin Wu, Quan Yuan, Phoenix Meadowlark, Matthew Bilotti, Raymond Lin, Mukund Sridhar, Yannick Schroecker, Da-Woon Chung, Jincheng Luo, Trevor Strohman, Tianlin Liu, Anne Zheng, Jesse Emond, Wei Wang, Andrew Lampinen, Toshiyuki Fukuzawa, Folawiyo Campbell-Ajala, Monica Roy, James Lee-Thorp, Lily Wang, Iftekhar Naim, Tony, Nguy ên, Guy Bensky, Aditya Gupta, Dominika Rogozińska, Justin Fu, Thanumalayan Sankaranarayana Pillai, Petar Veličković, Shahar Drath, Philipp Neubeck, Vaibhav Tulsyan, Arseniy Klimovskiy, Don Metzler, Sage Stevens, Angel Yeh, Junwei Yuan, Tianhe Yu, Kelvin Zhang, Alec Go, Vincent Tsang, Ying Xu, Andy Wan, Isaac Galatzer-Levy, Sam Sobell, Abodunrinwa Toki, Elizabeth Salesky, Wenlei Zhou, Diego Antognini, Sholto Douglas, Shimu Wu, Adam Lelkes, Frank Kim, Paul Cavallaro, Ana Salazar, Yuchi Liu, James Besley, Tiziana Refice, Yiling Jia, Zhang Li, Michal Sokolik, Arvind Kannan, Jon Simon, Jo Chick, Avia Aharon, Meet Gandhi, Mayank Daswani, Keyvan Amiri, Vighnesh Birodkar, Abe Ittycheriah, Peter Grabowski, Oscar Chang, Charles Sutton, Zhixin, Lai, Umesh Telang, Susie Sargsyan, Tao Jiang, Raphael Hoffmann, Nicole Brichtova, Matteo Hessel, Jonathan Halcrow, Sammy Jerome, Geoff Brown, Alex Tomala, Elena Buchatskaya, Dian Yu, Sachit Menon, Pol Moreno, Yuguo Liao, Vicky Zayats, Luming Tang, SQ Mah, Ashish Shenoy, Alex Siegman, Majid Hadian, Okwan Kwon, Tao Tu, Nima Khajehnouri, Ryan Foley, Parisa Haghani, Zhongru Wu, Vaishakh Keshava, Khyatti Gupta, Tony Bruguier, Rui Yao, Danny Karmon, Luisa Zintgraf, Zhicheng Wang, Enrique Piqueras, Junehyuk Jung, Jenny Brennan, Diego Machado, Marissa Giustina, MH Tessler, Kamyu Lee, Qiao Zhang, Joss Moore, Kaspar Daugaard, Alexander Frömmgen, Jennifer Beattie, Fred Zhang, Daniel Kasenberg, Ty Geri, Danfeng Qin, Gaurav Singh Tomar, Tom Ouyang, Tianli Yu, Luowei Zhou, Rajiv Mathews, Andy Davis, Yaoyiran Li, Jai Gupta, Damion Yates, Linda Deng, Elizabeth Kemp, Ga-Young Joung, Sergei Vassilvitskii, Mandy Guo, Pallavi LV, Dave Dopson, Sami Lachgar, Lara McConnaughey, Himadri Choudhury, Dragos Dena, Aaron Cohen, Joshua Ainslie, Sergey Levi, Parthasarathy Gopavarapu, Polina Zablotskaia, Hugo Vallet, Sanaz Bahargam, Xiaodan Tang, Nenad Tomasev, Ethan Dyer, Daniel Balle, Hongrae Lee, William Bono, Jorge Gonzalez Mendez, Vadim Zubov, Shentao Yang, Ivor Rendulic, Yanyan Zheng, Andrew Hogue, Golan Pundak, Ralph Leith, Avishkar Bhoopchand, Michael Han, Mislav Žanić, Tom Schaul, Manolis Delakis, Tejas Iyer, Guanyu Wang, Harman Singh, Abdelrahman Abdelhamed, Tara Thomas, Siddhartha Brahma, Hilal Dib, Naveen Kumar, Wenxuan Zhou, Liang Bai, Pushkar Mishra, Jiao Sun, Valentin Anklin, Roykrong Sukkerd, Lauren Agubuzu, Anton Briukhov, Anmol Gulati, Maximilian Sieb, Fabio Pardo, Sara Nasso, Junquan Chen, Kexin Zhu, Tiberiu Sosea, Alex Goldin, Keith Rush, Spurthi Amba Hombaiah, Andreas Noever, Allan Zhou, Sam Haves, Mary Phuong, Jake Ades, Yi ting Chen, Lin Yang, Joseph Pagadora, Stan Bileschi, Victor Cotruta, Rachel Saputro, Arijit Pramanik, Sean Ammirati, Dan Garrette, Kevin Villela, Tim Blyth, Canfer Akbulut, Neha Jha, Alban Rrustemi, Arissa Wongpanich, Chirag Nagpal, Yonghui Wu, Morgane Rivière, Sergey Kishchenko, Pranesh Srinivasan, Alice Chen, Animesh Sinha, Trang Pham, Bill Jia, Tom Hennigan, Anton Bakalov, Nithya Attaluri, Drew Garmon, Daniel Rodriguez, Dawid Wegner, Wenhao Jia, Evan Senter, Noah Fiedel, Denis Petek, Yuchuan Liu, Cassidy Hardin, Harshal Tushar Lehri, Joao Carreira, Sara Smoot, Marcel Prasetya, Nami Akazawa, Anca Stefanoiu, Chia-Hua Ho, Anelia Angelova, Kate Lin, Min Kim, Charles Chen, Marcin Sieniek, Alice Li, Tongfei Guo, Sorin Baltateanu, Pouya Tafti, Michael Wunder, Nadav Olmert, Divyansh Shukla, Jingwei Shen, Neel Kovelamudi, Balaji Venkatraman, Seth Neel, Romal Thoppilan, Jerome Connor, Frederik Benzing, Axel Stjerngren, Golnaz Ghiasi, Alex Polozov, Joshua Howland, Theophane Weber, Justin Chiu, Ganesh Poomal Girirajan, Andreas Terzis, Pidong Wang, Fangda Li, Yoav Ben Shalom, Dinesh Tewari, Matthew Denton, Roee Aharoni, Norbert Kalb, Heri Zhao, Junlin Zhang, Angelos Filos, Matthew Rahtz, Lalit Jain, Connie Fan, Vitor Rodrigues, Ruth Wang, Richard Shin, Jacob Austin, Roman Ring, Mariella Sanchez-Vargas, Mehadi Hassen, Ido Kessler, Uri Alon, Gufeng Zhang, Wenhu Chen, Yenai Ma, Xiance Si, Le Hou, Azalia Mirhoseini, Marc Wilson, Geoff Bacon, Becca Roelofs, Lei Shu, Gautam Vasudevan, Jonas Adler, Artur Dwornik, Tayfun Terzi, Matt Lawlor, Harry Askham, Mike Bernico, Xuanyi Dong, Chris Hidey, Kevin Kilgour, Gaël Liu, Surya Bhupatiraju, Luke Leonhard, Siqi Zuo, Partha Talukdar, Qing Wei, Aliaksei Severyn, Vít Listík, Jong Lee, Aditya Tripathi, SK Park, Yossi Matias, Hao Liu, Alex Ruiz, Rajesh Jayaram, Jackson Tolins, Pierre Marcenac, Yiming Wang, Bryan Seybold, Henry Prior, Deepak Sharma, Jack Weber, Mikhail Sirotenko, Yunhsuan Sung, Dayou Du, Ellie Pavlick, Stefan Zinke, Markus Freitag, Max Dylla, Montse Gonzalez Arenas, Natan Potikha, Omer Goldman, Connie Tao, Rachita Chhaparia, Maria Voitovich, Pawan Dogra, Andrija Ražnatović, Zak Tsai, Chong You, Oleaser Johnson, George Tucker, Chenjie Gu, Jae Yoo, Maryam Majzoubi, Valentin Gabeur, Bahram Raad, Rocky Rhodes,

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

Kashyap Kolipaka, Heidi Howard, Geta Sampemane, Benny Li, Chulayuth Asawaroengchai, Duy Nguyen, Chiyuan Zhang, Timothee Cour, Xinxin Yu, Zhao Fu, Joe Jiang, Po-Sen Huang, Gabriela Surita, Iñaki Iturrate, Yael Karov, Michael Collins, Martin Baeuml, Fabian Fuchs, Shilpa Shetty, Swaroop Ramaswamy, Sayna Ebrahimi, Qiuchen Guo, Jeremy Shar, Gabe Barth-Maron, Sravanti Addepalli, Bryan Richter, Chin-Yi Cheng, Eugénie Rives, Fei Zheng, Johannes Griesser, Nishanth Dikkala, Yoel Zeldes, Ilkin Safarli, Dipanjan Das, Himanshu Srivastava, Sadh MNM Khan, Xin Li, Aditya Pandey, Larisa Markeeva, Dan Belov, Qiqi Yan, Mikołaj Rybiński, Tao Chen, Megha Nawhal, Michael Quinn, Vineetha Govindaraj, Sarah York, Reed Roberts, Roopal Garg, Namrata Godbole, Jake Abernethy, Anil Das, Lam Nguyen Thiet, Jonathan Tompson, John Nham, Neera Vats, Ben Caine, Wesley Helmholz, Francesco Pongetti, Yeongil Ko, James An, Clara Huiyi Hu, Yu-Cheng Ling, Julia Pawar, Robert Leland, Keisuke Kinoshita, Waleed Khawaja, Marco Selvi, Eugene Ie, Danila Sinopalnikov, Lev Proleev, Nilesh Tripuraneni, Michele Bevilacqua, Seungji Lee, Clayton Sanford, Dan Suh, Dustin Tran, Jeff Dean, Simon Baumgartner, Jens Heitkaemper, Sagar Gubbi, Kristina Toutanova, Yichong Xu, Chandu Thekkath, Keran Rong, Palak Jain, Annie Xie, Yan Virin, Yang Li, Lubo Litchev, Richard Powell, Tarun Bharti, Adam Kraft, Nan Hua, Marissa Ikonomidis, Ayal Hitron, Sanjiv Kumar, Loic Matthey, Sophie Bridgers, Lauren Lax, Ishaan Malhi, Ondrej Skopek, Ashish Gupta, Jiawei Cao, Mitchelle Rasquinha, Siim Põder, Wojciech Stokowiec, Nicholas Roth, Guowang Li, Michaël Sander, Joshua Kessinger, Vihan Jain, Edward Loper, Wonpyo Park, Michal Yarom, Liqun Cheng, Guru Guruganesh, Kanishka Rao, Yan Li, Catarina Barros, Mikhail Sushkov, Chun-Sung Ferng, Rohin Shah, Ophir Aharoni, Ravin Kumar, Tim McConnell, Peiran Li, Chen Wang, Fernando Pereira, Craig Swanson, Fayaz Jamil, Yan Xiong, Anitha Vijayakumar, Prakash Shroff, Kedar Soparkar, Jindong Gu, Livio Baldini Soares, Eric Wang, Kushal Majmundar, Aurora Wei, Kai Bailey, Nora Kassner, Chizu Kawamoto, Goran Žužić, Victor Gomes, Abhirut Gupta, Michael Guzman, Ishita Dasgupta, Xinyi Bai, Zhufeng Pan, Francesco Piccinno, Hadas Natalie Vogel, Octavio Ponce, Adrian Hutter, Paul Chang, Pan-Pan Jiang, Ionel Gog, Vlad Ionescu, James Manyika, Fabian Pedregosa, Harry Ragan, Zach Behrman, Ryan Mullins, Coline Devin, Aroonalok Pyne, Swapnil Gawde, Martin Chadwick, Yiming Gu, Sasan Tavakkol, Andy Twigg, Naman Goyal, Ndidi Elue, Anna Goldie, Srinivasan Venkatachary, Hongliang Fei, Ziqiang Feng, Marvin Ritter, Isabel Leal, Sudeep Dasari, Pei Sun, Alif Raditya Rochman, Brendan O'Donoghue, Yuchen Liu, Jim Sproch, Kai Chen, Natalie Clay, Slav Petrov, Sailesh Sidhwani, Ioana Mihailescu, Alex Panagopoulos, AJ Piergiovanni, Yunfei Bai, George Powell, Deep Karkhanis, Trevor Yacovone, Petr Mitrichev, Joe Kovac, Dave Uthus, Amir Yazdanbakhsh, David Amos, Steven Zheng, Bing Zhang, Jin Miao, Bhuvana Ramabhadran, Soroush Radpour, Shantanu Thakoor, Josh Newlan, Oran Lang, Orion Jankowski, Shikhar Bharadwaj, Jean-Michel Sarr, Shereen Ashraf, Sneha Mondal, Jun Yan, Ankit Singh Rawat, Sarmishta Velury, Greg Kochanski, Tom Eccles, Franz Och, Abhanshu Sharma, Ethan Mahintorabi, Alex Gurney, Carrie Muir, Vered Cohen, Saksham Thakur, Adam Bloniarz, Asier Mujika, Alexander Pritzel, Paul Caron, Altaf Rahman, Fiona Lang, Yasumasa Onoe, Petar Sirkovic, Jay Hoover, Ying Jian, Pablo Duque, Arun Narayanan, David Soergel, Alex Haig, Loren Maggiore, Shyamal Buch, Josef Dean, Ilya Figotin, Igor Karpov, Shaleen Gupta, Denny Zhou, Muhuan Huang, Ashwin Vaswani, Christopher Semturs, Kaushik Shivakumar, Yu Watanabe, Vinodh Kumar Rajendran, Eva Lu, Yanhan Hou, Wenting Ye, Shikhar Vashishth, Nana Nti, Vytenis Sakenas, Darren Ni, Doug DeCarlo, Michael Bendersky, Sumit Bagri, Nacho Cano, Elijah Peake, Simon Tokumine, Varun Godbole, Carlos Guía, Tanya Lando, Vittorio Selo, Seher Ellis, Danny Tarlow, Daniel Gillick, Alessandro Epasto, Siddhartha Reddy Jonnalagadda, Meng Wei, Meiyan Xie, Ankur Taly, Michela Paganini, Mukund Sundararajan, Daniel Toyama, Ting Yu, Dessie Petrova, Aneesh Pappu, Rohan Agrawal, Senaka Buthpitiya, Justin Frye, Thomas Buschmann, Remi Crocker, Marco Tagliasacchi, Mengchao Wang, Da Huang, Sagi Perel, Brian Wieder, Hideto Kazawa, Weiyue Wang, Jeremy Cole, Himanshu Gupta, Ben Golan, Seojin Bang, Nitish Kulkarni, Ken Franko, Casper Liu, Doug Reid, Sid Dalmia, Jay Whang, Kevin Cen, Prasha Sundaram, Johan Ferret, Berivan Isik, Lucian Ionita, Guan Sun, Anna Shekhawat, Muqthar Mohammad, Philip Pham, Ronny Huang, Karthik Raman, Xingyi Zhou, Ross Mcilroy, Austin Myers, Sheng Peng, Jacob Scott, Paul Covington, Sofia Erell, Pratik Joshi, João Gabriel Oliveira, Natasha Noy, Tajwar Nasir, Jake Walker, Vera Axelrod, Tim Dozat, Pu Han, Chun-Te Chu, Eugene Weinstein, Anand Shukla, Shreyas Chandrakaladharan, Petra Poklukar, Bonnie Li, Ye Jin, Prem Eruvbetine, Steven Hansen, Avigail Dabush, Alon Jacovi, Samrat Phatale, Chen Zhu, Steven Baker, Mo Shomrat, Yang Xiao, Jean Pouget-Abadie, Mingyang Zhang, Fanny Wei, Yang Song, Helen King, Yiling Huang, Yun Zhu, Ruoxi Sun, Juliana Vicente Franco, Chu-Cheng Lin, Sho Arora, Hui, Li, Vivian Xia, Luke Vilnis, Mariano Schain, Kaiz Alarakyia, Laurel Prince, Aaron Phillips, Caleb

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1019

1020

1023

1024

1025

Habtegebriel, Luyao Xu, Huan Gui, Santiago Ontanon, Lora Aroyo, Karan Gill, Peggy Lu, Yash Katariya, Dhruv Madeka, Shankar Krishnan, Shubha Srinivas Raghvendra, James Freedman, Yi Tay, Gaurav Menghani, Peter Choy, Nishita Shetty, Dan Abolafia, Doron Kukliansky, Edward Chou, Jared Lichtarge, Ken Burke, Ben Coleman, Dee Guo, Larry Jin, Indro Bhattacharya, Victoria Langston, Yiming Li, Suyog Kotecha, Alex Yakubovich, Xinyun Chen, Petre Petrov, Tolly Powell, Yanzhang He, Corbin Quick, Kanav Garg, Dawsen Hwang, Yang Lu, Srinadh Bhojanapalli, Kristian Kjems, Ramin Mehran, Aaron Archer, Hado van Hasselt, Ashwin Balakrishna, JK Kearns, Meiqi Guo, Jason Riesa, Mikita Sazanovich, Xu Gao, Chris Sauer, Chengrun Yang, XiangHai Sheng, Thomas Jimma, Wouter Van Gansbeke, Vitaly Nikolaev, Wei Wei, Katie Millican, Ruizhe Zhao, Justin Snyder, Levent Bolelli, Maura O'Brien, Shawn Xu, Fei Xia, Wentao Yuan, Arvind Neelakantan, David Barker, Sachin Yadav, Hannah Kirkwood, Farooq Ahmad, Joel Wee, Jordan Grimstad, Boyu Wang, Matthew Wiethoff, Shane Settle, Miaosen Wang, Charles Blundell, Jingjing Chen, Chris Duvarney, Grace Hu, Olaf Ronneberger, Alex Lee, Yuanzhen Li, Abhishek Chakladar, Alena Butryna, Georgios Evangelopoulos, Guillaume Desjardins, Jonni Kanerva, Henry Wang, Averi Nowak, Nick Li, Alyssa Loo, Art Khurshudov, Laurent El Shafey, Nagabhushan Baddi, Karel Lenc, Yasaman Razeghi, Tom Lieber, Amer Sinha, Xiao Ma, Yao Su, James Huang, Asahi Ushio, Hanna Klimczak-Plucińska, Kareem Mohamed, JD Chen, Simon Osindero, Stav Ginzburg, Lampros Lamprou, Vasilisa Bashlovkina, Duc-Hieu Tran, Ali Khodaei, Ankit Anand, Yixian Di, Ramy Eskander, Manish Reddy Vuyyuru, Jasmine Liu, Aishwarya Kamath, Roman Goldenberg, Mathias Bellaiche, Juliette Pluto, Bill Rosgen, Hassan Mansoor, William Wong, Suhas Ganesh, Eric Bailey, Scott Baird, Dan Deutsch, Jinoo Baek, Xuhui Jia, Chansoo Lee, Abe Friesen, Nathaniel Braun, Kate Lee, Amayika Panda, Steven M. Hernandez, Duncan Williams, Jianqiao Liu, Ethan Liang, Arnaud Autef, Emily Pitler, Deepali Jain, Phoebe Kirk, Oskar Bunyan, Jaume Sanchez Elias, Tongxin Yin, Machel Reid, Aedan Pope, Nikita Putikhin, Bidisha Samanta, Sergio Guadarrama, Dahun Kim, Simon Rowe, Marcella Valentine, Geng Yan, Alex Salcianu, David Silver, Gan Song, Richa Singh, Shuai Ye, Hannah DeBalsi, Majd Al Merey, Eran Ofek, Albert Webson, Shibl Mourad, Ashwin Kakarla, Silvio Lattanzi, Nick Roy, Evgeny Sluzhaev, Christina Butterfield, Alessio Tonioni, Nathan Waters, Sudhindra Kopalle, Jason Chase, James Cohan, Girish Ramchandra Rao, Robert Berry, Michael Voznesensky, Shuguang Hu, Kristen Chiafullo, Sharat Chikkerur, George Scrivener, Ivy Zheng, Jeremy Wiesner, Wolfgang Macherey, Timothy Lillicrap, Fei Liu, Brian Walker, David Welling, Elinor Davies, Yangsibo Huang, Lijie Ren, Nir Shabat, Alessandro Agostini, Mariko Iinuma, Dustin Zelle, Rohit Sathyanarayana, Andrea D'olimpio, Morgan Redshaw, Matt Ginsberg, Ashwin Murthy, Mark Geller, Tatiana Matejovicova, Ayan Chakrabarti, Ryan Julian, Christine Chan, Qiong Hu, Daniel Jarrett, Manu Agarwal, Jeshwanth Challagundla, Tao Li, Sandeep Tata, Wen Ding, Maya Meng, Zhuyun Dai, Giulia Vezzani, Shefali Garg, Jannis Bulian, Mary Jasarevic, Honglong Cai, Harish Rajamani, Adam Santoro, Florian Hartmann, Chen Liang, Bartek Perz, Apoorv Jindal, Fan Bu, Sungyong Seo, Ryan Poplin, Adrian Goedeckemeyer, Badih Ghazi, Nikhil Khadke, Leon Liu, Kevin Mather, Mingda Zhang, Ali Shah, Alex Chen, Jinliang Wei, Keshav Shivam, Yuan Cao, Donghyun Cho, Angelo Scorza Scarpati, Michael Moffitt, Clara Barbu, Ivan Jurin, Ming-Wei Chang, Hongbin Liu, Hao Zheng, Shachi Dave, Christine Kaeser-Chen, Xiaobin Yu, Alvin Abdagic, Lucas Gonzalez, Yanping Huang, Peilin Zhong, Cordelia Schmid, Bryce Petrini, Alex Wertheim, Jifan Zhu, Hoang Nguyen, Kaiyang Ji, Yanqi Zhou, Tao Zhou, Fangxiaoyu Feng, Regev Cohen, David Rim, Shubham Milind Phal, Petko Georgiev, Ariel Brand, Yue Ma, Wei Li, Somit Gupta, Chao Wang, Pavel Dubov, Jean Tarbouriech, Kingshuk Majumder, Huijian Li, Norman Rink, Apurv Suman, Yang Guo, Yinghao Sun, Arun Nair, Xiaowei Xu, Mohamed Elhawaty, Rodrigo Cabrera, Guangxing Han, Julian Eisenschlos, Junwen Bai, Yuqi Li, Yamini Bansal, Thibault Sellam, Mina Khan, Hung Nguyen, Justin Mao-Jones, Nikos Parotsidis, Jake Marcus, Cindy Fan, Roland Zimmermann, Yony Kochinski, Laura Graesser, Feryal Behbahani, Alvaro Caceres, Michael Riley, Patrick Kane, Sandra Lefdal, Rob Willoughby, Paul Vicol, Lun Wang, Shujian Zhang, Ashleah Gill, Yu Liang, Gautam Prasad, Soroosh Mariooryad, Mehran Kazemi, Zifeng Wang, Kritika Muralidharan, Paul Voigtlaender, Jeffrey Zhao, Huanjie Zhou, Nina D'Souza, Aditi Mavalankar, Séb Arnold, Nick Young, Obaid Sarvana, Chace Lee, Milad Nasr, Tingting Zou, Seokhwan Kim, Lukas Haas, Kaushal Patel, Neslihan Bulut, David Parkinson, Courtney Biles, Dmitry Kalashnikov, Chi Ming To, Aviral Kumar, Jessica Austin, Alex Greve, Lei Zhang, Megha Goel, Yeqing Li, Sergey Yaroshenko, Max Chang, Abhishek Jindal, Geoff Clark, Hagai Taitelbaum, Dale Johnson, Ofir Roval, Jeongwoo Ko, Anhad Mohananey, Christian Schuler, Shenil Dodhia, Ruichao Li, Kazuki Osawa, Claire Cui, Peng Xu, Rushin Shah, Tao Huang, Ela Gruzewska, Nathan Clement, Mudit Verma, Olcan Sercinoglu, Hai Qian, Viral Shah,

Masa Yamaguchi, Abhinit Modi, Takahiro Kosakai, Thomas Strohmann, Junhao Zeng, Beliz Gunel, Jun Qian, Austin Tarango, Krzysztof Jastrzebski, Robert David, Jyn Shan, Parker Schuh, Kunal Lad, Willi Gierke, Mukundan Madhavan, Xinyi Chen, Mark Kurzeja, Rebeca Santamaria-Fernandez, Dawn Chen, Alexandra Cordell, Yuri Chervonyi, Frankie Garcia, Nithish Kannen, Vincent Perot, Nan Ding, Shlomi Cohen-Ganor, Victor Lavrenko, Junru Wu, Georgie Evans, Cicero Nogueira dos Santos, Madhavi Sewak, Ashley Brown, Andrew Hard, Joan Puigcerver, Zeyu Zheng, Yizhong Liang, Evgeny Gladchenko, Reeve Ingle, Uri First, Pierre Sermanet, Charlotte Magister, Mihajlo Velimirović, Sashank Reddi, Susanna Ricco, Eirikur Agustsson, Hartwig Adam, Nir Levine, David Gaddy, Dan Holtmann-Rice, Xuanhui Wang, Ashutosh Sathe, Abhijit Guha Roy, Blaž Bratanič, Alen Carin, Harsh Mehta, Silvano Bonacina, Nicola De Cao, Mara Finkelstein, Verena Rieser, Xinyi Wu, Florent Altché, Dylan Scandinaro, Li Li, Nino Vieillard, Nikhil Sethi, Garrett Tanzer, Zhi Xing, Shibo Wang, Parul Bhatia, Gui Citovsky, Thomas Anthony, Sharon Lin, Tianze Shi, Shoshana Jakobovits, Gena Gibson, Raj Apte, Lisa Lee, Mingqing Chen, Arunkumar Byravan, Petros Maniatis, Kellie Webster, Andrew Dai, Pu-Chin Chen, Jiaqi Pan, Asya Fadeeva, Zach Gleicher, Thang Luong, and Niket Kumar Bhumihar. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities, 2025. URL https://arxiv.org/abs/2507.06261.

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1057

1058

1061

1062

1063

1064

1067

1068

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1039

1040

DeepSeek-AI, Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Haowei Zhang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, Junxiao Song, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Lei Xu, Leyi Xia, Liang Zhao, Litong Wang, Liyue Zhang, Meng Li, Miaojun Wang, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Mingming Li, Ning Tian, Panpan Huang, Peiyi Wang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu, Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou, Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun, W. L. Xiao, Wangding Zeng, Wanjia Zhao, Wei An, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, X. Q. Li, Xiangyue Jin, Xianzu Wang, Xiao Bi, Xiaodong Liu, Xiaohan Wang, Xiaojin Shen, Xiaokang Chen, Xiaokang Zhang, Xiaosha Chen, Xiaotao Nie, Xiaowen Sun, Xiaoxiang Wang, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xingkai Yu, Xinnan Song, Xinxia Shan, Xinyi Zhou, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. X. Zhu, Yang Zhang, Yanhong Xu, Yanhong Xu, Yanping Huang, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Li, Yaohui Wang, Yi Yu, Yi Zheng, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Ying Tang, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yu Wu, Yuan Ou, Yuchen Zhu, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yukun Zha, Yunfan Xiong, Yunxian Ma, Yuting Yan, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Z. F. Wu, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu, Zhen Huang, Zhen Zhang, Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhibin Gou, Zhicheng Ma, Zhigang Yan, Zhihong Shao, Zhipeng Xu, Zhiyu Wu, Zhongyu Zhang, Zhuoshu Li, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Ziyi Gao, and Zizheng Pan. Deepseek-v3 technical report, 2025. URL https://arxiv.org/abs/2412.19437.

1069 1070 1071

Chao Deng, Jiale Yuan, Pi Bu, Peijie Wang, Zhong-Zhi Li, Jian Xu, Xiao-Hui Li, Yuan Gao, Jun Song, Bo Zheng, and Cheng-Lin Liu. LongDocURL: a comprehensive multimodal long document benchmark integrating understanding, reasoning, and locating. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1135–1159, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.57. URL https://aclanthology.org/2025.acl-long.57/.

1077 1078 1079

1074

1075

Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani, Gautier Viaud, CELINE HUDELOT, and Pierre Colombo. Colpali: Efficient document retrieval with vision language models. In

- The Thirteenth International Conference on Learning Representations, 2025. URL https://openreview.net/forum?id=ogjBpZ8uSi.
  - Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Qianyu Guo, Meng Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A survey. *ArXiv*, abs/2312.10997, 2023. URL https://api.semanticscholar.org/CorpusID:266359151.
  - Yingqiang Ge, Wenyue Hua, Kai Mei, jianchao ji, Juntao Tan, Shuyuan Xu, Zelong Li, and Yongfeng Zhang. Openagi: When Ilm meets domain experts. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), Advances in Neural Information Processing Systems, volume 36, pp. 5539-5568. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper\_files/paper/2023/file/1190733f217404edc8a7f4e15a57f301-Paper-Datasets\_and\_Benchmarks.pdf.
  - Daming Guo, Dongdong Yang, Hongyi Zhang, et al. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645:633–638, 2025. doi: 10.1038/s41586-025-09422-z.
  - Yichen He, Guanhua Huang, Peiyuan Feng, Yuan Lin, Yuchen Zhang, Hang Li, and Weinan E. Pasa: An Ilm agent for comprehensive academic paper search, 2025. URL https://arxiv.org/abs/2501.10120.
  - Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learning, 2022. URL https://arxiv.org/abs/2112.09118.
  - Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. NV-embed: Improved techniques for training LLMs as generalist embedding models. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=lgsyLSsDRe.
  - Lei Li, Yuqi Wang, Runxin Xu, Peiyi Wang, Xiachong Feng, Lingpeng Kong, and Qi Liu. Multimodal ArXiv: A dataset for improving scientific comprehension of large vision-language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 14369–14387, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.775. URL https://aclanthology.org/2024.acl-long.775/.
  - Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. Llava-next: Improved reasoning, ocr, and world knowledge, January 2024. URL https://llava-vl.github.io/blog/2024-01-30-llava-next/.
  - Renze Lou, Hanzi Xu, Sijia Wang, Jiangshu Du, Ryo Kamoi, Xiaoxin Lu, Jian Xie, Yuxuan Sun, Yusen Zhang, Jihyun Janice Ahn, Hongchao Fang, Zhuoyang Zou, Wenchao Ma, Xi Li, Kai Zhang, Congying Xia, Lifu Huang, and Wenpeng Yin. AAAR-1.0: Assessing AI's potential to assist research. In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=RHAWcjIy12.
  - Yubo Ma, Yuhang Zang, Liangyu Chen, Meiqi Chen, Yizhu Jiao, Xinze Li, Xinyuan Lu, Ziyu Liu, Yan Ma, Xiaoyi Dong, Pan Zhang, Liangming Pan, Yu-Gang Jiang, Jiaqi Wang, Yixin Cao, and Aixin Sun. MMLONGBENCH-DOC: Benchmarking long-context document understanding with visualizations. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL https://openreview.net/forum?id=loJMlacwzf.
  - Meta. Llama 4 | model cards and prompt formats. https://www.llama.com/docs/model-cards-and-prompt-formats/llama4/, 2025. Official model card and prompt format documentation.
- Meredith Ringel Morris, Jascha Sohl-Dickstein, Noah Fiedel, Tris Warkentin, Allan Dafoe, Aleksandra Faust, Clement Farabet, and Shane Legg. Position: Levels of AGI for operationalizing progress on the path to AGI. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id=0ofzEysK2D.

1138 1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158 1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

OpenAI. Gpt-5 system card. Technical report, OpenAI, August 2025. URL https://cdn.openai.com/gpt-5-system-card.pdf. Version: August 13, 2025.

OpenAI. gpt-4.1 — openai api documentation, 2025. URL https://platform.openai.com/docs/models/gpt-4.1. Accessed: 2025-09-25.

OpenAI, :, Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K. Arora, Yu Bai, Bowen Baker, Haiming Bao, Boaz Barak, Ally Bennett, Tyler Bertao, Nivedita Brett, Eugene Brevdo, Greg Brockman, Sebastien Bubeck, Che Chang, Kai Chen, Mark Chen, Enoch Cheung, Aidan Clark, Dan Cook, Marat Dukhan, Casey Dvorak, Kevin Fives, Vlad Fomenko, Timur Garipov, Kristian Georgiev, Mia Glaese, Tarun Gogineni, Adam Goucher, Lukas Gross, Katia Gil Guzman, John Hallman, Jackie Hehir, Johannes Heidecke, Alec Helyar, Haitang Hu, Romain Huet, Jacob Huh, Saachi Jain, Zach Johnson, Chris Koch, Irina Kofman, Dominik Kundel, Jason Kwon, Volodymyr Kyrylov, Elaine Ya Le, Guillaume Leclerc, James Park Lennon, Scott Lessans, Mario Lezcano-Casado, Yuanzhi Li, Zhuohan Li, Ji Lin, Jordan Liss, Lily, Liu, Jiancheng Liu, Kevin Lu, Chris Lu, Zoran Martinovic, Lindsay McCallum, Josh McGrath, Scott McKinney, Aidan McLaughlin, Song Mei, Steve Mostovoy, Tong Mu, Gideon Myles, Alexander Neitz, Alex Nichol, Jakub Pachocki, Alex Paino, Dana Palmie, Ashley Pantuliano, Giambattista Parascandolo, Jongsoo Park, Leher Pathak, Carolina Paz, Ludovic Peran, Dmitry Pimenov, Michelle Pokrass, Elizabeth Proehl, Huida Qiu, Gaby Raila, Filippo Raso, Hongyu Ren, Kimmy Richardson, David Robinson, Bob Rotsted, Hadi Salman, Suvansh Sanjeev, Max Schwarzer, D. Sculley, Harshit Sikchi, Kendal Simon, Karan Singhal, Yang Song, Dane Stuckey, Zhiqing Sun, Philippe Tillet, Sam Toizer, Foivos Tsimpourlas, Nikhil Vyas, Eric Wallace, Xin Wang, Miles Wang, Olivia Watkins, Kevin Weil, Amy Wendling, Kevin Whinnery, Cedric Whitney, Hannah Wong, Lin Yang, Yu Yang, Michihiro Yasunaga, Kristen Ying, Wojciech Zaremba, Wenting Zhan, Cyril Zhang, Brian Zhang, Eddie Zhang, and Shengjia Zhao. gpt-oss-120b gpt-oss-20b model card, 2025. URL https://arxiv.org/abs/2508.10925.

Long Phan, Alice Gatti, Ziwen Han, Nathaniel Li, Josephina Hu, Hugh Zhang, Chen Bo Calvin Zhang, Mohamed Shaaban, John Ling, Sean Shi, Michael Choi, Anish Agrawal, Arnay Chopra, Adam Khoja, Ryan Kim, Richard Ren, Jason Hausenloy, Oliver Zhang, Mantas Mazeika, Dmitry Dodonov, Tung Nguyen, Jaeho Lee, Daron Anderson, Mikhail Doroshenko, Alun Cennyth Stokes, Mobeen Mahmood, Oleksandr Pokutnyi, Oleg Iskra, Jessica P. Wang, John-Clark Levin, Mstyslav Kazakov, Fiona Feng, Steven Y. Feng, Haoran Zhao, Michael Yu, Varun Gangal, Chelsea Zou, Zihan Wang, Serguei Popov, Robert Gerbicz, Geoff Galgon, Johannes Schmitt, Will Yeadon, Yongki Lee, Scott Sauers, Alvaro Sanchez, Fabian Giska, Marc Roth, Søren Riis, Saiteja Utpala, Noah Burns, Gashaw M. Goshu, Mohinder Maheshbhai Naiya, Chidozie Agu, Zachary Giboney, Antrell Cheatom, Francesco Fournier-Facio, Sarah-Jane Crowson, Lennart Finke, Zerui Cheng, Jennifer Zampese, Ryan G. Hoerr, Mark Nandor, Hyunwoo Park, Tim Gehrunger, Jiaqi Cai, Ben McCarty, Alexis C Garretson, Edwin Taylor, Damien Sileo, Qiuyu Ren, Usman Qazi, Lianghui Li, Jungbae Nam, John B. Wydallis, Pavel Arkhipov, Jack Wei Lun Shi, Aras Bacho, Chris G. Willcocks, Hangrui Cao, Sumeet Motwani, Emily de Oliveira Santos, Johannes Veith, Edward Vendrow, Doru Cojoc, Kengo Zenitani, Joshua Robinson, Longke Tang, Yuqi Li, Joshua Vendrow, Natanael Wildner Fraga, Vladyslav Kuchkin, Andrey Pupasov Maksimov, Pierre Marion, Denis Efremov, Jayson Lynch, Kaiqu Liang, Aleksandar Mikov, Andrew Gritsevskiy, Julien Guillod, Gözdenur Demir, Dakotah Martinez, Ben Pageler, Kevin Zhou, Saeed Soori, Ori Press, Henry Tang, Paolo Rissone, Sean R. Green, Lina Brüssel, Moon Twayana, Aymeric Dieuleveut, Joseph Marvin Imperial, Ameya Prabhu, Jinzhou Yang, Nick Crispino, Arun Rao, Dimitri Zvonkine, Gabriel Loiseau, Mikhail Kalinin, Marco Lukas, Ciprian Manolescu, Nate Stambaugh, Subrata Mishra, Tad Hogg, Carlo Bosio, Brian P Coppola, Julian Salazar, Jaehyeok Jin, Rafael Sayous, Stefan Ivanov, Philippe Schwaller, Shaipranesh Senthilkuma, Andres M Bran, Andres Algaba, Kelsey Van den Houte, Lynn Van Der Sypt, Brecht Verbeken, David Noever, Alexei Kopylov, Benjamin Myklebust, Bikun Li, Lisa Schut, Evgenii Zheltonozhskii, Qiaochu Yuan, Derek Lim, Richard Stanley, Tong Yang, John Maar, Julian Wykowski, Martí Oller, Anmol Sahu, Cesare Giulio Ardito, Yuzheng Hu, Ariel Ghislain Kemogne Kamdoum, Alvin Jin, Tobias Garcia Vilchis, Yuexuan Zu, Martin Lackner, James Koppel, Gongbo Sun, Daniil S. Antonenko, Steffi Chern, Bingchen Zhao, Pierrot Arsene, Joseph M Cavanagh, Daofeng Li, Jiawei Shen, Donato Crisostomi, Wenjin Zhang, Ali Dehghan, Sergey Ivanov, David Perrella, Nurdin Kaparov, Allen Zang, Ilia Sucholutsky, Arina Kharlamova, Daniil Orel, Vladislav Poritski, Shalev Ben-David, Zachary Berger, Parker Whitfill, Michael Foster, Daniel Munro, Linh

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1201

1202

1203

1205

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1236

1237

1239

1240

Ho, Shankar Sivarajan, Dan Bar Hava, Aleksey Kuchkin, David Holmes, Alexandra Rodriguez-Romero, Frank Sommerhage, Anji Zhang, Richard Moat, Keith Schneider, Zakayo Kazibwe, Don Clarke, Dae Hyun Kim, Felipe Meneguitti Dias, Sara Fish, Veit Elser, Tobias Kreiman, Victor Efren Guadarrama Vilchis, Immo Klose, Ujjwala Anantheswaran, Adam Zweiger, Kaivalya Rawal, Jeffery Li, Jeremy Nguyen, Nicolas Daans, Haline Heidinger, Maksim Radionov, Václav Rozhoň, Vincent Ginis, Christian Stump, Niv Cohen, Rafał Poświata, Josef Tkadlec, Alan Goldfarb, Chenguang Wang, Piotr Padlewski, Stanislaw Barzowski, Kyle Montgomery, Ryan Stendall, Jamie Tucker-Foltz, Jack Stade, T. Ryan Rogers, Tom Goertzen, Declan Grabb, Abhishek Shukla, Alan Givré, John Arnold Ambay, Archan Sen, Muhammad Fayez Aziz, Mark H Inlow, Hao He, Ling Zhang, Younesse Kaddar, Ivar Ängquist, Yanxu Chen, Harrison K Wang, Kalyan Ramakrishnan, Elliott Thornley, Antonio Terpin, Hailey Schoelkopf, Eric Zheng, Avishy Carmi, Ethan D. L. Brown, Kelin Zhu, Max Bartolo, Richard Wheeler, Martin Stehberger, Peter Bradshaw, JP Heimonen, Kaustubh Sridhar, Ido Akov, Jennifer Sandlin, Yury Makarychev, Joanna Tam, Hieu Hoang, David M. Cunningham, Vladimir Goryachev, Demosthenes Patramanis, Michael Krause, Andrew Redenti, David Aldous, Jesyin Lai, Shannon Coleman, Jiangnan Xu, Sangwon Lee, Ilias Magoulas, Sandy Zhao, Ning Tang, Michael K. Cohen, Orr Paradise, Jan Hendrik Kirchner, Maksym Ovchynnikov, Jason O. Matos, Adithya Shenoy, Michael Wang, Yuzhou Nie, Anna Sztyber-Betley, Paolo Faraboschi, Robin Riblet, Jonathan Crozier, Shiv Halasyamani, Shreyas Verma, Prashant Joshi, Eli Meril, Ziqiao Ma, Jérémy Andréoletti, Raghav Singhal, Jacob Platnick, Volodymyr Nevirkovets, Luke Basler, Alexander Ivanov, Seri Khoury, Nils Gustafsson, Marco Piccardo, Hamid Mostaghimi, Qijia Chen, Virendra Singh, Tran Quoc Khánh, Paul Rosu, Hannah Szlyk, Zachary Brown, Himanshu Narayan, Aline Menezes, Jonathan Roberts, William Alley, Kunyang Sun, Arkil Patel, Max Lamparth, Anka Reuel, Linwei Xin, Hanmeng Xu, Jacob Loader, Freddie Martin, Zixuan Wang, Andrea Achilleos, Thomas Preu, Tomek Korbak, Ida Bosio, Fereshteh Kazemi, Ziye Chen, Biró Bálint, Eve J. Y. Lo, Jiaqi Wang, Maria Inês S. Nunes, Jeremiah Milbauer, M Saiful Bari, Zihao Wang, Behzad Ansarinejad, Yewen Sun, Stephane Durand, Hossam Elgnainy, Guillaume Douville, Daniel Tordera, George Balabanian, Hew Wolff, Lynna Kvistad, Hsiaoyun Milliron, Ahmad Sakor, Murat Eron, Andrew Favre D. O., Shailesh Shah, Xiaoxiang Zhou, Firuz Kamalov, Sherwin Abdoli, Tim Santens, Shaul Barkan, Allison Tee, Robin Zhang, Alessandro Tomasiello, G. Bruno De Luca, Shi-Zhuo Looi, Vinh-Kha Le, Noam Kolt, Jiayi Pan, Emma Rodman, Jacob Drori, Carl J Fossum, Niklas Muennighoff, Milind Jagota, Ronak Pradeep, Honglu Fan, Jonathan Eicher, Michael Chen, Kushal Thaman, William Merrill, Moritz Firsching, Carter Harris, Stefan Ciobâcă, Jason Gross, Rohan Pandey, Ilya Gusev, Adam Jones, Shashank Agnihotri, Pavel Zhelnov, Mohammadreza Mofayezi, Alexander Piperski, David K. Zhang, Kostiantyn Dobarskyi, Roman Leventov, Ignat Soroko, Joshua Duersch, Vage Taamazyan, Andrew Ho, Wenjie Ma, William Held, Ruicheng Xian, Armel Randy Zebaze, Mohanad Mohamed, Julian Noah Leser, Michelle X Yuan, Laila Yacar, Johannes Lengler, Katarzyna Olszewska, Claudio Di Fratta, Edson Oliveira, Joseph W. Jackson, Andy Zou, Muthu Chidambaram, Timothy Manik, Hector Haffenden, Dashiell Stander, Ali Dasouqi, Alexander Shen, Bita Golshani, David Stap, Egor Kretov, Mikalai Uzhou, Alina Borisovna Zhidkovskaya, Nick Winter, Miguel Orbegozo Rodriguez, Robert Lauff, Dustin Wehr, Colin Tang, Zaki Hossain, Shaun Phillips, Fortuna Samuele, Fredrik Ekström, Angela Hammon, Oam Patel, Faraz Farhidi, George Medley, Forough Mohammadzadeh, Madellene Peñaflor, Haile Kassahun, Alena Friedrich, Rayner Hernandez Perez, Daniel Pyda, Taom Sakal, Omkar Dhamane, Ali Khajegili Mirabadi, Eric Hallman, Kenchi Okutsu, Mike Battaglia, Mohammad Maghsoudimehrabani, Alon Amit, Dave Hulbert, Roberto Pereira, Simon Weber, Handoko, Anton Peristyy, Stephen Malina, Mustafa Mehkary, Rami Aly, Frank Reidegeld, Anna-Katharina Dick, Cary Friday, Mukhwinder Singh, Hassan Shapourian, Wanyoung Kim, Mariana Costa, Hubeyb Gurdogan, Harsh Kumar, Chiara Ceconello, Chao Zhuang, Haon Park, Micah Carroll, Andrew R. Tawfeek, Stefan Steinerberger, Daattavya Aggarwal, Michael Kirchhof, Linjie Dai, Evan Kim, Johan Ferret, Jainam Shah, Yuzhou Wang, Minghao Yan, Krzysztof Burdzy, Lixin Zhang, Antonio Franca, Diana T. Pham, Kang Yong Loh, Joshua Robinson, Abram Jackson, Paolo Giordano, Philipp Petersen, Adrian Cosma, Jesus Colino, Colin White, Jacob Votava, Vladimir Vinnikov, Ethan Delaney, Petr Spelda, Vit Stritecky, Syed M. Shahid, Jean-Christophe Mourrat, Lavr Vetoshkin, Koen Sponselee, Renas Bacho, Zheng-Xin Yong, Florencia de la Rosa, Nathan Cho, Xiuyu Li, Guillaume Malod, Orion Weller, Guglielmo Albani, Leon Lang, Julien Laurendeau, Dmitry Kazakov, Fatimah Adesanya, Julien Portier, Lawrence Hollom, Victor Souza, Yuchen Anna Zhou, Julien Degorre, Yiğit Yalın, Gbenga Daniel Obikoya, Rai, Filippo Bigi, M. C. Boscá, Oleg Shumar, Kaniuar Bacho, Gabriel Recchia, Mara Popescu, Nikita Shulga, Ngefor Mildred Tanwie,

1243

1244

1245

1246

1247

1248

1249

1250

1251

1252

1253

1255

1256

1257

1259

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1276

1278

1279

1280

1281

1282

1283

1284

1285

1286

1290 1291

1293

1294

1295

Thomas C. H. Lux, Ben Rank, Colin Ni, Matthew Brooks, Alesia Yakimchyk, Huanxu, Liu, Stefano Cavalleri, Olle Häggström, Emil Verkama, Joshua Newbould, Hans Gundlach, Leonor Brito-Santana, Brian Amaro, Vivek Vajipey, Rynaa Grover, Ting Wang, Yosi Kratish, Wen-Ding Li, Sivakanth Gopi, Andrea Caciolai, Christian Schroeder de Witt, Pablo Hernández-Cámara, Emanuele Rodolà, Jules Robins, Dominic Williamson, Vincent Cheng, Brad Raynor, Hao Qi, Ben Segev, Jingxuan Fan, Sarah Martinson, Erik Y. Wang, Kaylie Hausknecht, Michael P. Brenner, Mao Mao, Christoph Demian, Peyman Kassani, Xinyu Zhang, David Avagian, Eshawn Jessica Scipio, Alon Ragoler, Justin Tan, Blake Sims, Rebeka Plecnik, Aaron Kirtland, Omer Faruk Bodur, D. P. Shinde, Yan Carlos Leyva Labrador, Zahra Adoul, Mohamed Zekry, Ali Karakoc, Tania C. B. Santos, Samir Shamseldeen, Loukmane Karim, Anna Liakhovitskaia, Nate Resman, Nicholas Farina, Juan Carlos Gonzalez, Gabe Maayan, Earth Anderson, Rodrigo De Oliveira Pena, Elizabeth Kelley, Hodjat Mariji, Rasoul Pouriamanesh, Wentao Wu, Ross Finocchio, Ismail Alarab, Joshua Cole, Danyelle Ferreira, Bryan Johnson, Mohammad Safdari, Liangti Dai, Siriphan Arthornthurasuk, Isaac C. McAlister, Alejandro José Moyano, Alexey Pronin, Jing Fan, Angel Ramirez-Trinidad, Yana Malysheva, Daphiny Pottmaier, Omid Taheri, Stanley Stepanic, Samuel Perry, Luke Askew, Raúl Adrián Huerta Rodríguez, Ali M. R. Minissi, Ricardo Lorena, Krishnamurthy Iyer, Arshad Anil Fasiludeen, Ronald Clark, Josh Ducey, Matheus Piza, Maja Somrak, Eric Vergo, Juehang Qin, Benjámin Borbás, Eric Chu, Jack Lindsey, Antoine Jallon, I. M. J. McInnis, Evan Chen, Avi Semler, Luk Gloor, Tej Shah, Marc Carauleanu, Pascal Lauer, Tran uc Huy, Hossein Shahrtash, Emilien Duc, Lukas Lewark, Assaf Brown, Samuel Albanie, Brian Weber, Warren S. Vaz, Pierre Clavier, Yiyang Fan, Gabriel Poesia Reis e Silva, Long, Lian, Marcus Abramovitch, Xi Jiang, Sandra Mendoza, Murat Islam, Juan Gonzalez, Vasilios Mavroudis, Justin Xu, Pawan Kumar, Laxman Prasad Goswami, Daniel Bugas, Nasser Heydari, Ferenc Jeanplong, Thorben Jansen, Antonella Pinto, Archimedes Apronti, Abdallah Galal, Ng Ze-An, Ankit Singh, Tong Jiang, Joan of Arc Xavier, Kanu Priya Agarwal, Mohammed Berkani, Gang Zhang, Zhehang Du, Benedito Alves de Oliveira Junior, Dmitry Malishev, Nicolas Remy, Taylor D. Hartman, Tim Tarver, Stephen Mensah, Gautier Abou Loume, Wiktor Morak, Farzad Habibi, Sarah Hoback, Will Cai, Javier Gimenez, Roselynn Grace Montecillo, Jakub Łucki, Russell Campbell, Asankhaya Sharma, Khalida Meer, Shreen Gul, Daniel Espinosa Gonzalez, Xavier Alapont, Alex Hoover, Gunjan Chhablani, Freddie Vargus, Arunim Agarwal, Yibo Jiang, Deepakkumar Patil, David Outevsky, Kevin Joseph Scaria, Rajat Maheshwari, Abdelkader Dendane, Priti Shukla, Ashley Cartwright, Sergei Bogdanov, Niels Mündler, Sören Möller, Luca Arnaboldi, Kunvar Thaman, Muhammad Rehan Siddiqi, Prajvi Saxena, Himanshu Gupta, Tony Fruhauff, Glen Sherman, Mátyás Vincze, Siranut Usawasutsakorn, Dylan Ler, Anil Radhakrishnan, Innocent Enyekwe, Sk Md Salauddin, Jiang Muzhen, Aleksandr Maksapetyan, Vivien Rossbach, Chris Harjadi, Mohsen Bahaloohoreh, Claire Sparrow, Jasdeep Sidhu, Sam Ali, Song Bian, John Lai, Eric Singer, Justine Leon Uro, Greg Bateman, Mohamed Sayed, Ahmed Menshawy, Darling Duclosel, Dario Bezzi, Yashaswini Jain, Ashley Aaron, Murat Tiryakioglu, Sheeshram Siddh, Keith Krenek, Imad Ali Shah, Jun Jin, Scott Creighton, Denis Peskoff, Zienab EL-Wasif, Ragavendran P V, Michael Richmond, Joseph McGowan, Tejal Patwardhan, Hao-Yu Sun, Ting Sun, Nikola Zubić, Samuele Sala, Stephen Ebert, Jean Kaddour, Manuel Schottdorf, Dianzhuo Wang, Gerol Petruzella, Alex Meiburg, Tilen Medved, Ali ElSheikh, S Ashwin Hebbar, Lorenzo Vaquero, Xianjun Yang, Jason Poulos, Vilém Zouhar, Sergey Bogdanik, Mingfang Zhang, Jorge Sanz-Ros, David Anugraha, Yinwei Dai, Anh N. Nhu, Xue Wang, Ali Anil Demircali, Zhibai Jia, Yuyin Zhou, Juncheng Wu, Mike He, Nitin Chandok, Aarush Sinha, Gaoxiang Luo, Long Le, Mickaël Noyé, Michał Perełkiewicz, Ioannis Pantidis, Tianbo Qi, Soham Sachin Purohit, Letitia Parcalabescu, Thai-Hoa Nguyen, Genta Indra Winata, Edoardo M. Ponti, Hanchen Li, Kaustubh Dhole, Jongee Park, Dario Abbondanza, Yuanli Wang, Anupam Nayak, Diogo M. Caetano, Antonio A. W. L. Wong, Maria del Rio-Chanona, Dániel Kondor, Pieter Francois, Ed Chalstrey, Jakob Zsambok, Dan Hoyer, Jenny Reddish, Jakob Hauser, Francisco-Javier Rodrigo-Ginés, Suchandra Datta, Maxwell Shepherd, Thom Kamphuis, Qizheng Zhang, Hyunjun Kim, Ruiji Sun, Jianzhu Yao, Franck Dernoncourt, Satyapriya Krishna, Sina Rismanchian, Bonan Pu, Francesco Pinto, Yingheng Wang, Kumar Shridhar, Kalon J. Overholt, Glib Briia, Hieu Nguyen, David, Soler Bartomeu, Tony CY Pang, Adam Wecker, Yifan Xiong, Fanfei Li, Lukas S. Huber, Joshua Jaeger, Romano De Maddalena, Xing Han Lù, Yuhui Zhang, Claas Beger, Patrick Tser Jern Kon, Sean Li, Vivek Sanker, Ming Yin, Yihao Liang, Xinlu Zhang, Ankit Agrawal, Li S. Yifei, Zechen Zhang, Mu Cai, Yasin Sonmez, Costin Cozianu, Changhao Li, Alex Slen, Shoubin Yu, Hyun Kyu Park, Gabriele Sarti, Marcin Briański, Alessandro Stolfo, Truong An Nguyen, Mike Zhang, Yotam Perlitz, Jose Hernandez-Orallo, Runjia Li, Amin Shabani, Felix Juefei-Xu, Shikhar Dhingra,

Orr Zohar, My Chiffon Nguyen, Alexander Pondaven, Abdurrahim Yilmaz, Xuandong Zhao, Chuanyang Jin, Muyan Jiang, Stefan Todoran, Xinyao Han, Jules Kreuer, Brian Rabern, Anna Plassart, Martino Maggetti, Luther Yap, Robert Geirhos, Jonathon Kean, Dingsu Wang, Sina Mollaei, Chenkai Sun, Yifan Yin, Shiqi Wang, Rui Li, Yaowen Chang, Anjiang Wei, Alice Bizeul, Xiaohan Wang, Alexandre Oliveira Arrais, Kushin Mukherjee, Jorge Chamorro-Padial, Jiachen Liu, Xingyu Qu, Junyi Guan, Adam Bouyamourn, Shuyu Wu, Martyna Plomecka, Junda Chen, Mengze Tang, Jiaqi Deng, Shreyas Subramanian, Haocheng Xi, Haoxuan Chen, Weizhi Zhang, Yinuo Ren, Haoqin Tu, Sejong Kim, Yushun Chen, Sara Vera Marjanović, Junwoo Ha, Grzegorz Luczyna, Jeff J. Ma, Zewen Shen, Dawn Song, Cedegao E. Zhang, Zhun Wang, Gaël Gendron, Yunze Xiao, Leo Smucker, Erica Weng, Kwok Hao Lee, Zhe Ye, Stefano Ermon, Ignacio D. Lopez-Miguel, Theo Knights, Anthony Gitter, Namkyu Park, Boyi Wei, Hongzheng Chen, Kunal Pai, Ahmed Elkhanany, Han Lin, Philipp D. Siedler, Jichao Fang, Ritwik Mishra, Károly Zsolnai-Fehér, Xilin Jiang, Shadab Khan, Jun Yuan, Rishab Kumar Jain, Xi Lin, Mike Peterson, Zhe Wang, Aditya Malusare, Maosen Tang, Isha Gupta, Ivan Fosin, Timothy Kang, Barbara Dworakowska, Kazuki Matsumoto, Guangyao Zheng, Gerben Sewuster, Jorge Pretel Villanueva, Ivan Rannev, Igor Chernyavsky, Jiale Chen, Deepayan Banik, Ben Racz, Wenchao Dong, Jianxin Wang, Laila Bashmal, Duarte V. Gonçalves, Wei Hu, Kaushik Bar, Ondrej Bohdal, Atharv Singh Patlan, Shehzaad Dhuliawala, Caroline Geirhos, Julien Wist, Yuval Kansal, Bingsen Chen, Kutay Tire, Atak Talay Yücel, Brandon Christof, Veerupaksh Singla, Zijian Song, Sanxing Chen, Jiaxin Ge, Kaustubh Ponkshe, Isaac Park, Tianneng Shi, Martin Q. Ma, Joshua Mak, Sherwin Lai, Antoine Moulin, Zhuo Cheng, Zhanda Zhu, Ziyi Zhang, Vaidehi Patil, Ketan Jha, Qiutong Men, Jiaxuan Wu, Tianchi Zhang, Bruno Hebling Vieira, Alham Fikri Aji, Jae-Won Chung, Mohammed Mahfoud, Ha Thi Hoang, Marc Sperzel, Wei Hao, Kristof Meding, Sihan Xu, Vassilis Kostakos, Davide Manini, Yueying Liu, Christopher Toukmaji, Jay Paek, Eunmi Yu, Arif Engin Demircali, Zhiyi Sun, Ivan Dewerpe, Hongsen Qin, Roman Pflugfelder, James Bailey, Johnathan Morris, Ville Heilala, Sybille Rosset, Zishun Yu, Peter E. Chen, Woongyeong Yeo, Eeshaan Jain, Ryan Yang, Sreekar Chigurupati, Julia Chernyavsky, Sai Prajwal Reddy, Subhashini Venugopalan, Hunar Batra, Core Francisco Park, Hieu Tran, Guilherme Maximiano, Genghan Zhang, Yizhuo Liang, Hu Shiyu, Rongwu Xu, Rui Pan, Siddharth Suresh, Ziqi Liu, Samaksh Gulati, Songyang Zhang, Peter Turchin, Christopher W. Bartlett, Christopher R. Scotese, Phuong M. Cao, Aakaash Nattanmai, Gordon McKellips, Anish Cheraku, Asim Suhail, Ethan Luo, Marvin Deng, Jason Luo, Ashley Zhang, Kavin Jindel, Jay Paek, Kasper Halevy, Allen Baranov, Michael Liu, Advaith Avadhanam, David Zhang, Vincent Cheng, Brad Ma, Evan Fu, Liam Do, Joshua Lass, Hubert Yang, Surya Sunkari, Vishruth Bharath, Violet Ai, James Leung, Rishit Agrawal, Alan Zhou, Kevin Chen, Tejas Kalpathi, Ziqi Xu, Gavin Wang, Tyler Xiao, Erik Maung, Sam Lee, Ryan Yang, Roy Yue, Ben Zhao, Julia Yoon, Sunny Sun, Aryan Singh, Ethan Luo, Clark Peng, Tyler Osbey, Taozhi Wang, Daryl Echeazu, Hubert Yang, Timothy Wu, Spandan Patel, Vidhi Kulkarni, Vijaykaarti Sundarapandiyan, Ashley Zhang, Andrew Le, Zafir Nasim, Srikar Yalam, Ritesh Kasamsetty, Soham Samal, Hubert Yang, David Sun, Nihar Shah, Abhijeet Saha, Alex Zhang, Leon Nguyen, Laasya Nagumalli, Kaixin Wang, Alan Zhou, Aidan Wu, Jason Luo, Anwith Telluri, Summer Yue, Alexandr Wang, and Dan Hendrycks. Humanity's last exam, 2025. URL https://arxiv.org/abs/2501.14249.

1334 1335 1336

1337

1338

1339

1342

1331

1332

1333

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

1307

1309

1310

1311

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324

1325

Shraman Pramanick, Rama Chellappa, and Subhashini Venugopalan. SPIQA: A dataset for multimodal question answering on scientific papers. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024. URL https://openreview.net/forum?id=h3lddsY5nf.

1340 1341

Stephen E. Robertson, Steve Walker, Susan Jones, Micheline Hancock-Beaulieu, and Mike Gatford. Okapi at trec-3. In *Text Retrieval Conference*, 1994. URL https://api.semanticscholar.org/CorpusID:41563977.

1343 1344

R. Smith. An overview of the tesseract ocr engine. In *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, volume 2, pp. 629–633, 2007. doi: 10.1109/ICDAR. 2007.4376991.

1347 1348

1349

Rubèn Tito, Minesh Mathew, C. V. Jawahar, Ernest Valveny, and Dimosthenis Karatzas. ICDAR 2021 competition on document visualquestion answering. *CoRR*, abs/2111.05547, 2021. URL https://arxiv.org/abs/2111.05547.

- Qiuchen Wang, Ruixue Ding, Yu Zeng, Zehui Chen, Lin Chen, Shihang Wang, Pengjun Xie, Fei Huang, and Feng Zhao. Vrag-rl: Empower vision-perception-based rag for visually rich information understanding via iterative reasoning with reinforcement learning, 2025. URL https://arxiv.org/abs/2505.22019.
- Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, Alexis Chevalier, Sanjeev Arora, and Danqi Chen. Charxiv: Charting gaps in realistic chart understanding in multimodal llms, 2024. URL https://arxiv.org/abs/2406.18521.
  - xAI. Grok 4 fast model card. Technical report, xAI, September 2025. URL https://data.x.ai/2025-09-19-grok-4-fast-model-card.pdf. Last updated: September 19, 2025.
  - Dawei Yan, Yang Li, Qing-Guo Chen, Weihua Luo, Peng Wang, Haokui Zhang, and Chunhua Shen. Mmcr: Advancing visual language model in multimodal multi-turn contextual reasoning, 2025. URL https://arxiv.org/abs/2503.18533.
  - An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.
  - Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *CoRR*, abs/1809.09600, 2018. URL http://arxiv.org/abs/1809.09600.
  - Shi Yu, Chaoyue Tang, Bokai Xu, Junbo Cui, Junhao Ran, Yukun Yan, Zhenghao Liu, Shuo Wang, Xu Han, Zhiyuan Liu, and Maosong Sun. VisRAG: Vision-based retrieval-augmented generation on multi-modality documents. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=zG459X3Xge.
  - Xiang Yue, Yuansheng Ni, Tianyu Zheng, Kai Zhang, 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 Wenhu Chen. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 9556–9567, 2024. doi: 10.1109/CVPR52733.2024.00913.
  - Yilun Zhao, Yitao Long, Hongjun Liu, Ryo Kamoi, Linyong Nan, Lyuhao Chen, Yixin Liu, Xiangru Tang, Rui Zhang, and Arman Cohan. DocMath-eval: Evaluating math reasoning capabilities of LLMs in understanding long and specialized documents. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 16103–16120, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.852. URL https://aclanthology.org/2024.acl-long.852/.
  - Ge Zheng, Bin Yang, Jiajin Tang, Hong-Yu Zhou, and Sibei Yang. Ddcot: Duty-distinct chain-of-thought prompting for multimodal reasoning in language models. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 5168–5191. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper\_files/paper/2023/file/108030643e640ac050e0ed5e6aace48f-Paper-Conference.pdf.
  - Ruiyang Zhou, Lu Chen, and Kai Yu. Is LLM a reliable reviewer? a comprehensive evaluation of LLM on automatic paper reviewing tasks. In Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue (eds.), *Proceedings of the 2024*

#### Not Search, But Scan: Benchmarking MLLMs on Scan-Oriented Academic Paper Reasoning Supplementary Material **Table of Contents in Appendix** A Use of LLMs **B** Prompts **Examples from Existing Datasets** Example from DocVQA **Dataset Annotation and Construction** D.3.3Common Failure Cases of MLLMs

E.2

E.4

	E.5	MO (Measurement & Operationalization)
	E.6	DHP (Data Handling & Preprocessing)
	E.7	CF (Computation & Formulae)
	E.8	IC (Inference & Conclusions)
	E.9	LE (Language & Expression)
	2.,	22 (2miguage de 2mpression) * * * * * * * * * * * * * * * * * * *
$\mathbf{F}$	Hun	nan-Machine Consistency Evaluation
		·

# A USE OF LLMS

Large language models (LLMs) were used solely to assist in language editing and stylistic refinement during manuscript preparation. All technical content, experiments, dataset construction, evaluation protocols, and analysis were conceived, implemented, and validated entirely by the authors. No LLMs were involved in the generation of benchmark data, research methodology design, or result interpretation. The use of LLMs did not influence the scientific conclusions of this paper.

## B PROMPTS

1620

1621 1622

16231624

1625 1626

1627

1628

1629

1630

1631

1632

1633

1634

1635

1636

1637

1638

1640

1641

1642

1643

1644

1645

1646

1647

1648

1649

1650

1651

1652

1653

1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666

1667

1668

1669

1670

1671 1672 1673

#### B.1 WITHIN-GENERATE PROMPT

Within-Generate Prompt You will receive a high-quality, already accepted scientific paper as a PDF. Working only with the PDF itself (and any appendix embedded in the same PDF), edit specific textual spans to inject one or more errors chosen only from the taxonomy below, such that the errors are hard yet clearly identifiable by a professional reviewer reading the PDF alone. Error Type (fixed): Research Ouestion & Definitions Definition: The core construct/hypothesis/variable is insufficiently or inconsistently defined (conceptual vs operational), leaving the estimand ambiguous. Design & Identifiability Definition: Given a clear estimand, the design violates structural identification conditions so the effect is not identifiable even with infinite data and perfect measurement. Sampling & Generalizability Definition: The sampling frame/process/composition or cluster/power setup does not support valid or stable sample → population claims. Measurement & Operationalization Definition: Measures/manipulations lack feasibility/ reliability/validity/timing, so observed variables systematically diverge from the intended construct/ treatment. Data Handling & Preprocessing Definition: Pipeline choices in missing handling, joins/ keys, temporal splitting, feature construction, or partitioning introduce bias (incl. leakage or unit/ scale conflicts). Computation & Formulae Definition: Arithmetic/algebra/notation errors (totals/ ratios, unit conversion, CI vs point estimate, p-value vs label, symbol reuse, undefined variables, dimension mismatch). Inference & Conclusions Definition: Interpretations or causal statements exceed what methods/data support, or contradict the shown statistics/tables/captions. Referential and Citation Alignment Definition: Contradictions about the same quantity/term across text, tables, captions, or appendix within the paper. Language & Expression Definition: Terminology/capitalization/grammar ambiguities that affect meaning or domain-critical term

consistency (not cosmetic typos).

1674 Within-Generate Prompt (Continued) 1675 1676 Global constraints (must comply) 1677 1. Each error must map to exactly one primary category in the 1678 taxonomy. Do not mix causes. 1679 2. Each error must involve more than 2 micro-edits (each edit 1680 ≤ 20 English words) spread across distinct pages or 1681 paragraphs. 1682 3. If an edit would create an immediate contradiction in the 1683 same sentence/paragraph/caption, you may add shadow patch 1684 (es) for the same error to keep the text natural (still 1685 counted as edit locations). 1686 4. Independence across errors (per-copy generation) Generate each error on a separate copy of the original PDF 1687 . Different errors must be logically and operationally 1688 independent: 1689 No progression or variant relations: an error must not be 1690 a stricter/looser version, superset/subset, or minor 1691 wording variant of another error. 1692 No anchor reuse: do not target the same sentence/caption/ 1693 table cell or reuse the same old str (or a near-1694 duplicate paraphrase) across different errors. 1695 Applying any single error in isolation to the original PDF 1696 must still yield a detectable, clearly categorizable 1697 error according to the taxonomy. 5. Every error must be supportable using text inside the PDF. 1698 Do not rely on external supplementary files or prior 1699 knowledge. 1700 6. Design as difficult as possible but clean errors. Prefer 1701 edits that force cross-checking between two spots (e.g., 1702 Methods vs Results). Avoid trivialities. Edits must 1703 remain locally plausible and not advertise themselves via 1704 obviously artificial phrases (e.g., avoid contrived 1705 tokens purely added to be detectable). 1706 7. ''No cosmetic issues'' applies except for I (Language & 1707 Expression). For I, edits must affect meaning or domain-1708 critical terminology (e.g., ambiguous phrasing, inconsistent technical terms). Pure typos, punctuation 1709 tweaks, or layout nits are not allowed. 1710 8. Do not edit titles, author lists, bibliography entries, 1711 equation numbering, figure images, or add new figures/ 1712 tables/references. 1713 9. Frame each question as a neutral imperative that asks for 1714

1715

1716

1717

1718

a decision about a specific condition, using (but not

suggestive intensifiers (e.g., clearly/obviously/likely/

limited to) Decide/Determine/Judge/Evaluate/Assess

whether.... Do not presuppose an outcome or use

suspicious as examples).

```
1728
                             Within-Generate Prompt (Continued)
1729
1730
         10. Output English-only and strictly follow the JSON schema
1731
            below. Do not include any additional text outside the
1732
1733
         [
1734
1735
            "id":"1-based integer as string",
1736
            "modify":[
1737
              {
1738
                "location": "Page number + short unique nearby quote (
1739
                   \leq15 tokens).",
1740
                "old_str": "Exact original text from the PDF (verbatim)
                   .",
1741
                "new_str":"Edited text after your change."
1742
              }
1743
              /* Add 1-2 more locations; each location ≤ 20 words
1744
                 changed.
1745
                Shadow patches for local coherence count as locations.
1746
                    */
1747
1748
            "question": "One neutral audit-style task (1-25 words).",
1749
            "explanation": "Explain in 2-4 sentences why a reviewer can
1750
                 detect this error from the edited PDF alone.",
1751
            "Type": "Name the primary category (e.g., Inference &
                Conclusions).",
1752
1753
           /* More Errors */
1754
         1
1755
1756
1757
1758
1759
```

1782 B.2 WITHIN-SAMPLE PROMPT 1783 1784 Within-Sample Prompt 1785 1786 You will receive a paper PDF and the weaknesses mentioned in 1787 its peer-review comments. Your task is, based only on the 1788 content of that PDF, to sample from the review comments 1789 and verify possible errors related to the categories 1790 below, and for each confirmed or highly plausible error, 1791 generate one question and one explanation. 1792 1793 Error Type (fixed): Research Question & Definitions 1794 Definition: The core construct/hypothesis/variable is 1795 insufficiently or inconsistently defined (conceptual 1796 vs operational), leaving the estimand ambiguous. 1797 Design & Identifiability 1798 Definition: Given a clear estimand, the design violates 1799 structural identification conditions so the effect is 1800 not identifiable even with infinite data and perfect 1801 measurement. 1802 Sampling & Generalizability 1803 Definition: The sampling frame/process/composition or 1804 cluster/power setup does not support valid or stable 1805 sample → population claims. Measurement & Operationalization 1806 Definition: Measures/manipulations lack feasibility/ 1807 reliability/validity/timing, so observed variables 1808 systematically diverge from the intended construct/ 1809 treatment. 1810 Data Handling & Preprocessing 1811 Definition: Pipeline choices in missing handling, joins/ 1812 keys, temporal splitting, feature construction, or 1813 partitioning introduce bias (incl. leakage or unit/ 1814 scale conflicts). 1815 Computation & Formulae Definition: Arithmetic/algebra/notation errors (totals/ 1816 ratios, unit conversion, CI vs point estimate, p-value 1817 vs label, symbol reuse, undefined variables, 1818 dimension mismatch). 1819 Inference & Conclusions 1820 Definition: Interpretations or causal statements exceed 1821 what methods/data support, or contradict the shown 1822 statistics/tables/captions. 1823 Referential and Citation Alignment; 1824 Definition: Contradictions about the same quantity/term 1825 across text, tables, captions, or appendix within the 1826 paper. Language & Expression 1827 Definition: Terminology/capitalization/grammar ambiguities 1828 that affect meaning or domain-critical term 1829

34

consistency (not cosmetic typos).

#### Within-Sample Prompt (Continued) Global constraints (must comply) Output only the specified categories; even if other error types appear in the reviews, do not output them. Sample first, then verify: extract candidates from the review comments, then confirm them in the PDF. If you cannot locate supporting anchors in the PDF (page number plus phrase/label), do not output that candidate. Questions must be neutral and non-leading: use an "audit task + decision" style, avoiding yes/no bias. Independence: each question must target a different figure or different textual anchor; no minor variants of the same issue. Evidence first: the explanation must cite locatable anchors in the PDF (page number + original phrase/caption). You may mention a key short phrase from the review as a clue, but write the question and explanation in your own words Language & format: both question and explanation must be in English; output JSON only, with no extra text. Quantity: sort by evidence strength and output up to 5 items; if none qualify, output an empty array []. Example output [ "id": "1", "question": "Audit y-axis baselines and possible axis breaks in Figure 2; decide presence/absence and cite evidence.", "explanation": "The review flags possible exaggeration in Fig. 2. In the PDF (p.6, caption 'Performance vs baseline'), the y-axis starts at 0.85 with a break, magnifying small differences; panels use different ranges." "Type": "Visualization & Presentation Bias"

**B.3** Extractor Prompt

1890 1891 1892 Extractor Prompt 1893 1894 You will receive three inputs: 1895 Q: the open-ended question; 1896 E: the gold explanation (describes exactly one error; extra 1897 details still belong to the same single error); 1898 A: the model's answer to be evaluated. 1899 Your job is to extract counts only and output a single JSON 1900 object with the exact schema below. Do not compute any 1901 scores. Do not add fields. 1902 Core selection rule (multiple errors in A) 1903 1. Parse E into a single gold error (the "target error"). 1904 2. From A, identify how many distinct error claims are made. 1905 Cluster together mentions that support the same error ( 1906 multiple locations for one error are still one error). 1907 3. Existence decision (binary correctness only): 1908 Let the gold existence be 1 if E asserts an error exists, 1909 else 0. 1910 Let the predicted existence be 1 if A asserts any error, else 1911 0 (e.g., states no error). 1912 Set existance = 1 if predicted existence equals gold 1913 existence; otherwise set existance = 0. 4. If existance = 0: set contains\_target\_error = 0; set all 1914 location and reasoning counts to 0; and set 1915 unrelated\_errors to the total number of distinct error 1916 claims in A. Then output the JSON. 1917 5. If existance = 1: 1918 If the gold existence is 1: determine whether A contains the 1919 target error (match by the main error idea in E: category 1920 /intent/scope; treat E's subpoints as the same error). 1921 If yes, set contains\_target\_error = 1 and compute location 1922 and reasoning only for the target error. Count all 1923 other error claims in A as unrelated\_errors. If no, set contains\_target\_error = 0; set all location and 1924 reasoning counts to 0; set unrelated\_errors to the 1925 total number of distinct error claims in A. 1926 If the gold existence is 0: set contains\_target\_error = 0; 1927 set all location and reasoning counts to 0; set 1928 unrelated\_errors to the total number of distinct error 1929 claims in A. (These negative items are for binary 1930 accuracy only; they are not used for detailed scoring.) 1931 1932 Matching quidance (A error  $\leftrightarrow$  target error): match by the 1933 main error idea in E (category/intent/scope), not by 1934 wording. Treat E's subpoints as part of the same single error. Prefer the best-matching cluster in A; if ties, 1935 choose the one with stronger alignment to E's core claim. 1936 1937

```
1944
                               Extractor Prompt (Continued)
1945
1946
         Counting rules
1947
         Location (for the target error only when existance=1 and
1948
            contains_target_error=1):
1949
         gold_steps: number of unique error locations described in E (
1950
            after normalization and deduplication).
1951
        hit_steps: number of predicted locations in A that match any
1952
            gold location for the target error.
1953
         extra_steps: number of predicted locations in A for the
1954
            target error that do not match any gold location.
1955
1956
        Reasoning (for the target error only when existance=1 and
            contains_target_error=1):
1957
         Convert E into a canonical set or ordered chain of reasoning
1958
            steps for the target error.
1959
         gold_steps: total number of such steps.
1960
         reached steps:
1961
            single-chain tasks: length of the longest valid prefix of
1962
               A along the gold chain;
1963
            multi-path/parallel tasks: size of the intersection
1964
               between A's steps and the gold step set (or the
1965
               maximum across gold paths if multiple are defined).
1966
        missing_steps: gold_steps - reached_steps (non-negative
1967
            integer).
         Unrelated errors:
1968
         unrelated_errors: number of distinct error claims in A that
1969
            are not the target error (0 if none).
1970
         Output schema (return exactly this JSON; integers only)
1971
1972
          "existance": 0,
1973
          "contains_target_error": 0,
1974
          "location": {
1975
            "gold steps": 0,
1976
            "hit_steps": 0,
1977
            "extra steps": 0
1978
          },
          "reasoning": {
1979
            "gold_steps": 0,
1980
            "reached_steps": 0,
1981
            "missing_steps": 0
1982
          },
1983
          "unrelated_errors": 0
1984
         }
1985
1986
```

B.4 SYSTEM PROMPT

System Prompt

You are a neutral, careful academic reviewer. You will receive an open-ended question and the paper content. The paper may or may not have issues related to the question Do not assume there are errors. If the question is about citations, you will be given a citing paper and a cited paper; evaluate only the citing paper for possible issues and use the cited paper only as the reference for comparison. Write in natural prose with no fixed template

## Rules:

Speak only when sure. State an error only if you are confident it is a real error (not a mere weakness).

Stay on scope. Discuss only what the question asks about. Evidence completeness. For every error you state, list all distinct evidence cues you are confident about from the PDF. Include plain identifiers (figure/table/section/equation/citation) or quotes. Avoid redundant repeats of the exact same instance; include all distinct locations needed to support the error.

Be clear and brief. Use short, direct sentences.

No metaphors. No fancy wording. No guesses or outside sources.

Do not invent figures, tables, equations, citations, or results.

Report as many distinct, well-supported errors as you can within scope. If none are clear, write exactly: "No clear issue relevant to the question." and nothing else.

## C EXAMPLES FROM EXISTING DATASETS

## C.1 EXAMPLE FROM DOCMATH-EVAL

## One Example from DocMath-Eval

Question\_ID: complong-testmini-30

Question: What is the percentage of total offering cost on the total amount raised in the IPO if the total offering cost is \$14,528,328 and each unit sold is \$10?

## **Context Modalities: Texts Documents**

- 1. Offering costs consist of legal, accounting and other costs incurred through the balance sheet date that are directly related to the Initial Public Offering. Offering costs amounting to \$14,528,328 were charged to shareholders' equity upon the completion of the Initial Public Offering.
- 2. Pursuant to the Initial Public Offering on July 20, 2020, the Company sold 25,300,000 Units, which includes the full exercise by the underwriter of its option to purchase an additional 3,300,000 Units, at a purchase price of \$10.00 per Unit. Each Unit consists of one Class A ordinary share and one-half of one redeemable warrant ("Public Warrant"). Each whole Public Warrant entitles the holder to purchase one Class A ordinary share at an exercise price of \$11.50 per whole share (see Note 7).

## **NO Multi-Modal Documents Context**

## Covered areas:

## **Focus Only On the Field of Mathematics**

## **Cross-evidence Reasoning:**

Focusing on solving mathematical problems requires integrating evidence such as mathematical formulas, question stem conditions, and chart data from different positions in the document.

## Task Paradigm: search

## **Search-oriented**

## C.2 EXAMPLE FROM SLIDEVQA

2107 2108 2109

## 2110 Ouestion\_ID: 1

Question: How much difference in INR is there between the average order value of CY2013 and that of CY2012?

One Example from SlideVQA

2112 2113 2114

2115

2116

2117

2118

2125

2126

2127

2128

2129

2130

2135

2136

2137

2138

2111

## Context Modalities: Multi-Modal Documents and Texts **Executive summary** Key findings

Increasing average order value INR 3,600

2119 2120 INR 1.860 2121 INR 1,080 25% CAGR 2122 2123 CY2013 CY2012 2124

\$278 N

## CY2016P sories GMV \$2,811 M

\$559 M

## **KEY FINDING**

1. Last year there was a significant jump in average order value as there was a penetration of new

Average order values climbing up rapidly

categories like jewellery, home décor etc. 2. Also, users are becoming more comfortable buying higher priced items online.

## Fashion category doubled last year

1. Last year was the rise of the fashion category fashion e-commerce GMV doubled since 2012. 2. Given the young demographic which is shopping for latest looks online and increasing choice online - we estimate that this category will see 400% growth in the next 3 years and rival electronics and mobile category in GMV.

Accel estimates and Industry sources

## Covered areas:

The documents cover core technical research fields such as visual question answering and machine reading comprehension, as well as industry application fields including education and scientific research, finance and commerce, and healthcare (with derivative adaptation to pathological slice analysis), and also involves derivative technical fields like retrieval-augmented generation.

2139 2140 2141

2142 2143

## **Cross-evidence Reasoning:**

Simple question types only require a single piece of evidence

**Not Cross-evidence Reasoning** 

2144 2145 2146

## Task Paradigm: search

2147

2148 2149 2150

## **Search-oriented**

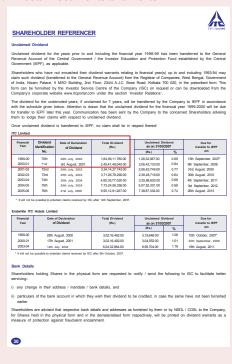
## C.3 EXAMPLE FROM MMLONGBENCH-DOC

One Example from MMLongBench-Doc

Question\_ID:

**Question**: How much higher was the proposed dividend paid (Rupees in lacs) in 2002 compared to 2001?

## Context Modalities: Multi-Modal Documents and Texts



## Covered areas:

The documents cover 7 diverse fields such as scientific research reports, business financial reports, and technical manuals.

## **Cross-evidence Reasoning:**

33% of the questions are cross-page questions, which require integrating different types of evidence such as texts, tables, and charts from multi-page documents

Task Paradigm: search

## C.4 EXAMPLE FROM LONGDOCURL

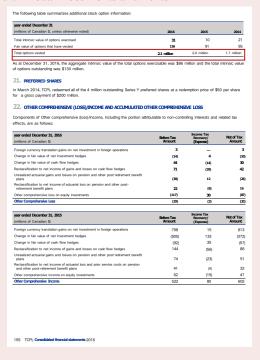
## 

One Example from LongDocURL

**Question\_ID**: free\_gemini15\_pro\_4061601\_47\_71\_8

Question: What was the total fair value of options that vested in 2016, 2015, and 2014, in millions of Canadian dollars?

## Context Modalities: Multi-Modal Documents and Texts



### Covered areas:

The document types of LongDocURL cover 8 major categories such as research reports, user manuals, and books.

## **Cross-evidence Reasoning:**

Most questions require integrating evidence across chapters and elements

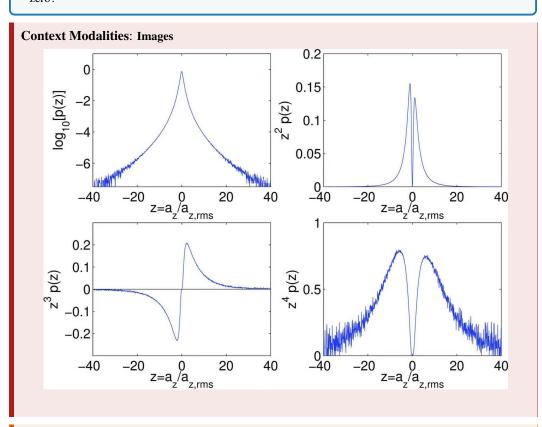
## Task Paradigm: search

## C.5 EXAMPLE FROM ARXIVQA

## One Example from ArXivQA

**Question\_ID**: physics-8049

**Question**: Based on the top-right graph, how would you describe the behavior of P(z) as z approaches zero?



## Covered areas:

The document includes arXiv academic papers in various fields such as physics and mathematics.

## **Covers Few Areas**

## **Cross-evidence Reasoning:**

Only focus on a single element.

## **Not Cross-evidence Reasoning**

Task Paradigm: search

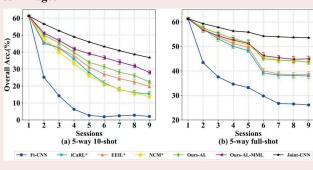
## C.6 EXAMPLE FROM CHARXIV

## One Example from Charxiv

**Question\_ID**: 2004.10956

**Question**: Which model shows a greater decline in accuracy from Session 1 to Session 9 in the 5-way full-shot scenario?

## **Context Modalities: Images**



## Covered areas:

The document type consists of multi-type charts and graphs from 2323 papers in 8 disciplines, namely physics, computer science, mathematics, biology, chemistry, statistics, engineering, and economics, which are derived from the arXiv platform.

## **Cross-evidence Reasoning:**

Only focus on a single element.

## **Not Cross-evidence Reasoning**

Task Paradigm: search

## C.7 EXAMPLE FROM AAAR

## **Question\_ID**: 1902.00751

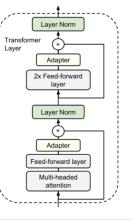
Question: What experiments do you suggest doing? Why do you suggest these experiments?

One Example from AAAR

## **Context Modalities: Multi-Modal Documents**

## Adapter Layer Feedforward up-project Nonlinearity

Feedforward



## Covered areas:

The document types cover two core categories: one is the AAAR-1.0 benchmark dataset documents, which are used to evaluate the research capabilities of LLMs and contain annotated data for 4 types of research tasks such as equation inference; the other is the documents related to the academic organization operation of the American Association for Aerosol Research (AAAR).

## **Covers Few Areas**

## **Cross-evidence Reasoning:**

It is necessary to integrate textual evidence across paragraphs and chapters.

Task Paradigm: search

## C.8 EXAMPLE FROM MMCR

## One Example from MMCR

## Question\_ID: 1

Question: Which module's weights are frozen?

## Context Modalities: Multi-Modal Documents and Texts

## | Victude Cas | Southeast University | Casy chandrows with a continuence of the Case | Southeast University | Casy chandrows with a continuence of the Case | Case

## Covered areas:

Its document type focuses on multimodal information fusion and clinical semantic understanding in medical scenarios.

## **Focus Only On the Field of Medicine**

## **Cross-evidence Reasoning:**

It is necessary to forcibly integrate medical imaging evidence (such as abnormal areas in CT images) with clinical report text evidence

## Task Paradigm: search

## C.9 EXAMPLE FROM DOCVQA

## One Example from DocVQA

Question\_ID: 24581

**Question**: What is name of university?

## **Context Modalities: Multi-Modal Documents**

UNIVER	RSITY OF CALIFORNIA, SAN DIEGO
То	Jack
Date	11/30/P2 Time 9:040M
	WHILE YOU WERE OUT
ďΩ Mr. Ms.	Wilson 455-8056
From	Janger Clinic
	ephoned   Will phone again   Please phone ne to see you   Will come again   Rush
w.	
No 11	MESSAGE
16. [] Kas	rogram Committee
16 D Kan	rogram Committee
le D Kas Green We	log Idn. It will
He Sold Second	vagrang committee
Massey	lieg Idn. It will, orbelg be fet or 3 - 2 they than latter half.
Resolution of the Phone por Taken by	rogram committee ley For It will bibly be 1st or I a ck in morch (1983) They than latter half:

## Covered areas:

Including invoices, resumes, academic papers, financial reports, manuals, etc. in formats such as scanned copies, PDFs, and screenshots.

## **Cross-evidence Reasoning:**

Simple question types (such as "invoice amount") only require evidence from a single location, while complex question types (such as "judging device compatibility based on parameter tables and explanatory texts across multiple pages of a manual") require integrating evidence across elements and locations.

## **Not All Cross-evidence Reasoning**

Task Paradigm: search

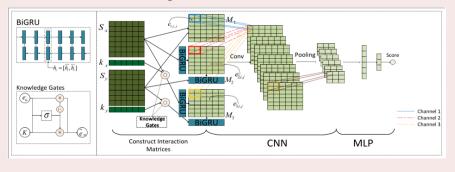
## C.10 EXAMPLE FROM SPIQA

## One Example from SPIQA

**Question\_ID**: 1611.04684v1

Question: What is the role of the knowledge gates in the KEHNN architecture?

## Context Modalities: Multi-Modal Figures and Charts



## Covered areas:

The document type of SPIQA originates from academic papers in fields such as computer science and physics.

## **Covered Few Areas**

## **Cross-evidence Reasoning:**

Only focus on a single element.

## **Not Cross-evidence Reasoning**

Task Paradigm: search

## D DATASET ANNOTATION AND CONSTRUCTION

## D.1 HUMAN ANNOTATOR GUIDELINES

The defective academic papers in our dataset are curated from three primary sources: (1) We synthetically inject 9 types of errors into papers accepted at ICLR and Nature Communications. (2) For the papers rejected by ICLR, we identified the shortcomings in the papers based on the reviewers' comments and categorized them into 9 error types.(3) For accepted ICLR papers, we generate consistency-related errors by cross-referencing their content against cited literature. To ensure the quality of each error, all entries undergo a rigorous, multistage validation protocol executed by human annotators. For synthetically generated errors, annotators manually embed them into the source papers following this protocol:

- **Credibility Validation**: Each error must be logically sound and verifiable. For generated errors, annotators first confirm their logical coherence and unambiguity. Flawed error descriptions are revised whenever possible; only irrepairable cases are discarded.
- Evidence Verification: All evidence substantiating an error must be either directly traceable to the source document or grounded in established domain-specific knowledge. Annotators are required to meticulously verify the origin and accuracy of all supporting data and background information.
- Category Classification: Each error must be accurately classified into one of the 9 predefined categories according to their formal definitions. Annotators verify the correctness of the assigned category and reclassify it if necessary.
- Paper Revision: Upon successful validation, annotators embed the generated error into the
  original manuscript by adding, deleting, or modifying relevant text segments as dictated by
  the error's specification.

This unified and standardized annotation protocol enables the creation of a high-quality dataset of academic papers with curated errors, providing a robust benchmark for evaluating the document sacnning and error detection capabilities of Large Multimodal Models.

## D.2 ANNOTATION STATISTICS

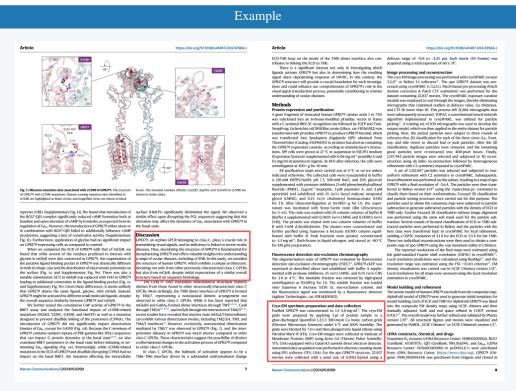
Initially, we generated or sampled a pool of 3,500 academic paper instances containing potential errors. During the manual annotation phase, following the protocol described above, we discarded 1,700 instances to ensure the logical rigor of the errors, the accuracy of the evidence, and a balanced distribution of categories.

Of the remaining 1,800 instances, 1,541 (85.6%)underwent manual revision. The distribution of these modifications is as follows:

- 535 questions were rewritten to eliminate ambiguity or to increase their retrieval and reasoning difficulty.
- 1,207 explanations were revised to correct erroneous evidence references and resolve logical flaws.
- 1,141 instances underwent category reclassification or manual paper editing. This process served to fix classifications that were inconsistent with our definitions and, for errors generated, to manually inject them into the source papers to create the flawed documents.

## D.3 EXAMPLES OF ANNOTATION

## D.3.1 CASE 1: DISCARD DIRECTLY



**Question**: Assess whether the conclusions drawn about the protein's functional state and therapeutic applicability are supported by the presented methods and results.

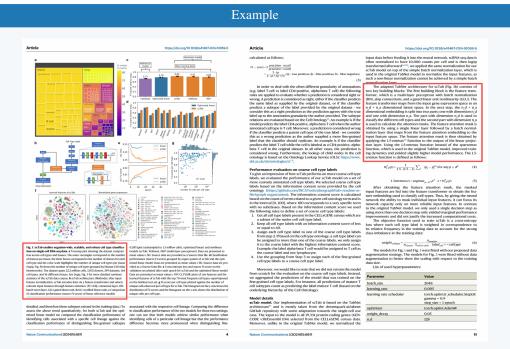
**Explanation**: Edits in the abstract and discussion claim the paper presents an active-state structure that reveals the activation mechanism and provides a roadmap for drug design. This overstates the findings, as the paper repeatedly describes solving the 'apo' (unbound) structure and explicitly states the activating ligand is unknown (p.6). To make the error subtle, a contradictory sentence was added to the methods (p.9) claiming a stabilizing agonist was used, but this is falsified by the numerous, unmodified mentions of the 'apo GPR179' structure throughout the results and methods.

**Error Type**: IC (Inference & Conclusions)

## **Decision**: Discard

**Analysis:** Based on the modifications, the revised abstract and conclusion claim that the paper elucidates the protein's 'active-state' structure and provides a roadmap for drug design. However, the original text repeatedly states (e.g., on pages 5 and 9) that it is the 'apo' (inactive) structure that was resolved, and critically notes on page 6 that the 'activating ligand is still unknown'. This constitutes a clear RCA-type error, defined by the inconsistent description of a concept within the article. Yet, the large model misclassifies this as an IC-type (Inference & Conclusions) error, which is a significant mistake. Considering that the inconsistency regarding the 'active-state' description is overly superficial and obvious—a type of error almost never encountered in actual academic literature—it lacks practical value. Even reclassifying it as an H-type question would be of little significance. Therefore, we have decided to delete this instance.

## D.3.2 CASE 2: MODIFY QUESTION



**Question:** Judge whether the mathematical description of the scTab model architecture presents any dimensional mismatches.

**Explanation**: The Methods section on page 11 (and the Figure 1b caption) now states that the model's feature transformer creates an embedding of dimension  $n_d$ . However, the same paragraph then describes splitting this embedding into two components of dimension  $n_d$  and  $n_d$  respectively. This is algebraically impossible for any non-zero  $n_a$ , and the hyperparameter table on page 12 confirms that  $n_a$  is set to 64, creating a fundamental dimensional mismatch in the model's description.

Error Type: CF (Computation & Formulae)

## Before:

**question**: Judge whether the mathematical description of the scTab model architecture presents any dimensional mismatches.

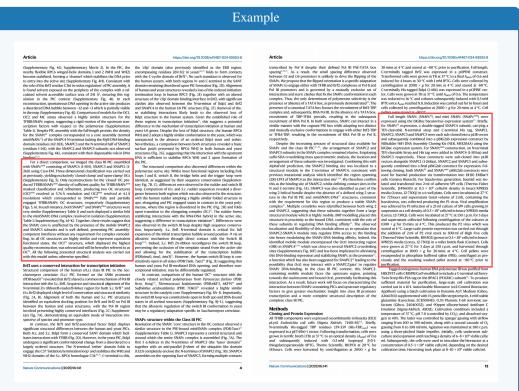
**Decision: Modify** 

### After:

question: Assess the Methods section for Computation & Formulae issues.

**Analysis:** Based on the error information and the text, the modified model description states that a vector of dimension  $n_d$  is split into two parts: one of dimension  $n_d$  and another of dimension  $n_a$ . This is algebraically impossible, as the total dimension  $(n_d)$  cannot equal the dimension of one of its parts  $(n_d)$  plus another non-zero part  $(n_a)$  is set to 64). This constitutes a clear dimensional mismatch, rendering the model's architectural description logically invalid. The original question was overly specific, as it explicitly prompted an assessment of whether the mathematical description of the scTab model architecture contained 'any dimensional mismatches'. This hint was too detailed, reducing the analytical difficulty for the model. To increase the difficulty, we have revised the question's phrasing to ask only whether the mathematical description of the scTab model architecture presents any problems.

## D.3.3 CASE 3: MODIFY EXPLANATION



**Question**: Evaluate if the composition of the SNAPcmini construct is consistently defined throughout the paper.

**Explanation**: The results on page 4 state that the assembled SNAPcmini construct includes the SNAPC2 subunit. However, the methods on page 12 describe the construction of SNAPcmini using only SNAPC4, SNAPC3, and SNAPC1, with SNAPC2 explicitly removed from the cloning description. A third conflicting statement on page 6 implies SNAPC2 was expected to be part of the minimal core, creating conceptual and operational inconsistency regarding this key experimental complex.

Error Type: RQD (Research Question & Definitions)

## Before:

**Explanation**: ...with SNAPC2 explicitly removed from the cloning description. A third conflicting statement on page 6...

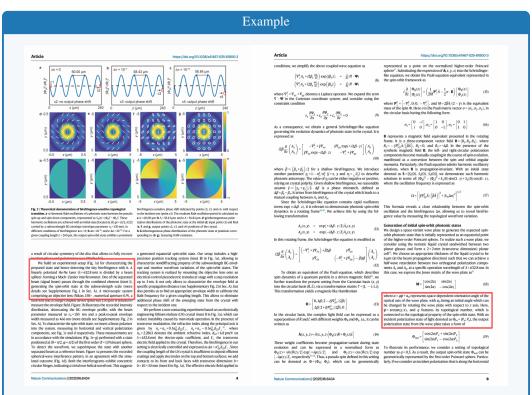
**Decision**: Modify

## After:

**Explanation**: ...with SNAPC2 explicitly removed from the cloning description.

**Analysis:** Based on the modifications, the revised abstract and conclusion claim that the paper elucidates the protein's 'active-state' structure and provides a roadmap for drug design. However, the original text repeatedly states (e.g., on pages 5 and 9) that it is the 'apo' (inactive) structure that was resolved, and critically notes on page 6 that the 'activating ligand is still unknown'. This constitutes a clear RCA-type error, defined by the inconsistent description of a concept within the article. Yet, the large model misclassifies this as an IC-type (Inference & Conclusions) error, which is a significant mistake. Considering that the inconsistency regarding the 'active-state' description is overly superficial and obvious—a type of error almost never encountered in actual academic literature—it lacks practical value. Even reclassifying it as an H-type question would be of little significance. Therefore, we have decided to delete this instance.

## D.3.4 CASE 4: MODIFY CATEGORY



**Question**: Assess whether the conclusions drawn about the protein's functional state and therapeutic applicability are supported by the presented methods and results.

**Explanation**: Edits in the abstract and discussion claim the paper presents an active-state structure that reveals the activation mechanism and provides a roadmap for drug design. This overstates the findings, as the paper repeatedly describes solving the 'apo' (unbound) structure and explicitly states the activating ligand is unknown (p.6). To make the error subtle, a contradictory sentence was added to the methods (p.9) claiming a stabilizing agonist was used, but this is falsified by the numerous, unmodified mentions of the 'apo GPR179' structure throughout the results and methods.

Error Type: MO (Measurement & Operationalization)

Before:

Error Type: MO (Measurement & Operationalization)

**Decision**: Modify

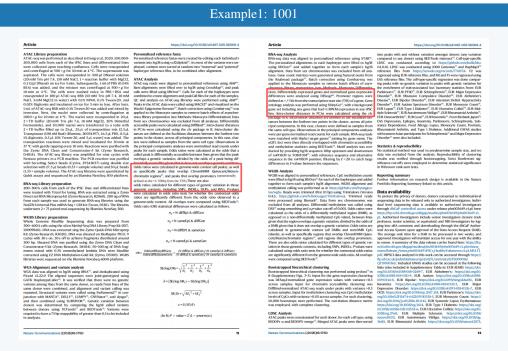
After:

**Error Type**: RCA (Referential and Citation Alignment)

**Analysis:** The introduced error systematically changes the laser wavelength used in the experiment to 532.0 nm. However, the calculation of a key physical quantity (birefringence) continues to use material constants (the electro-optic coefficient) that are only valid at the old wavelength of 632.8 nm. Because the optical properties of materials are wavelength-dependent, this systematic mismatch between experimental conditions and calculation parameters creates a significant contradiction in a core part of the paper. Compared to a Measurement & Operationalization (MO) error, this error is more accurately described as an internal inconsistency. Therefore, we are reclassifying this question from MO to RCA.

## E COMMON FAILURE CASES OF MLLMS

## E.1 RQD (RESEARCH QUESTION & DEFINITIONS)



**Question**: Assess the Methods section for Research Question & Definitions issues.

**Explanation**: The definition of a 'promoter region', a key analytical construct, is inconsistent across the paper, making the estimand ambiguous. The RNA-seq methods (page 12) define it as +/-1kb from the TSS, the ATAC-seq analysis methods (page 11) define it as +/-500bp from the TSS, and the Results section (page 4) defines it as +/-2kb from the TSS. These three conflicting operational definitions mean that analyses involving 'promoters' are not comparable and the construct is insufficiently defined.

Error Type: RQD (Research Question & Definitions)

Type: Within-Generate

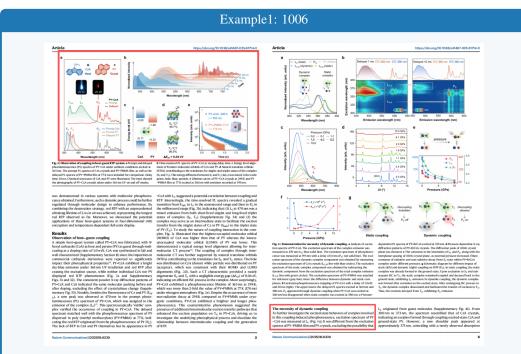
## Example2: 816 $\mathbb{E}\left[I_r \otimes \frac{1}{r}gg^{\top}\right] = \mathbb{E}[H(\text{vec}(\theta))].$ INVESTEE APPROXIMATION WITH SCHULZ S METHOD of (Petković, 1995). To compute the inverse of one matrix $A_t$ t of matrix iteration methods has attracted the attention of many resear genere guarantee (Altman, 1905; Garnett III et al., 1971; Bazán & Bo $X_{t+1} = X_t(I + T_t + T_t^2 + ... + T_t^{g-1}), \quad T_t = I - AX_t$ Question: Scrutinize the Methods section for Research Question & Definitions issues.

Explanation: Lemma 3.1, which is a cornerstone of the paper's theoretical contribution for low-rank Hessian approximation, relies on a strong and insufficiently justified assumption. The lemma states: "Assume that each column of the sample gradient ... is independent and identically distributed random vector with zero mean under the distribution  $p(y|x,\theta)$ ". The paper provides only a brief, hand-wavy justification (p.5, lines 230-232), suggesting it "could stand" in an "ideal case" of model convergence. These critical i.i.d. and zero-mean conditions are not rigorously established or empirically validated for the contexts in which the method is applied. This leaves a core hypothesis of the paper ambiguously defined and justified, which is an error of type Research Question & Definitions.

**Error Type**: RQD (Research Question & Definitions)

Type: Within-Sample

## E.2 DI (DESIGN & IDENTIFIABILITY)



Question: Assess the Experiments section for Design & Identifiability issues.

**Explanation**: The paper's core argument is that it identifies a specific 'dynamic coupling' pathway as essential for RTP, distinct from a 'static coupling' pathway. The edits state that the key experiment (excitation-phosphorescence mapping) cannot distinguish between these two pathways, as the final phosphorescence shows spectral signatures of originating from both. This introduces a structural identification problem: with two potential causal pathways leading to the same outcome and no way to isolate their effects, the claim that the dynamic pathway is the definitive mechanism is not identifiable from the data presented.

**Error Type**: DI (Design & Identifiability)

Type: Within-Generate

## Example2: 724 $W = USV^T$ , $\longrightarrow$ $\tilde{S}_{ii} = \begin{cases} \nu_i & \text{for } i \neq r \\ 0 & \text{else} \end{cases}$ $\longrightarrow$ $\tilde{W} = U\tilde{S}V^T$

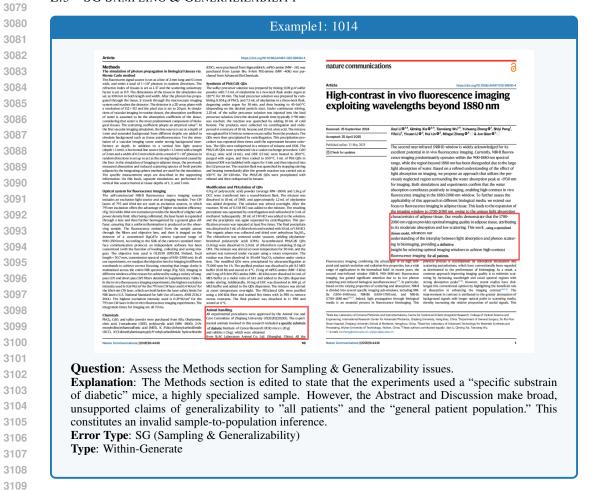
Question: Assess the Methods section for Design & Identifiability issues.

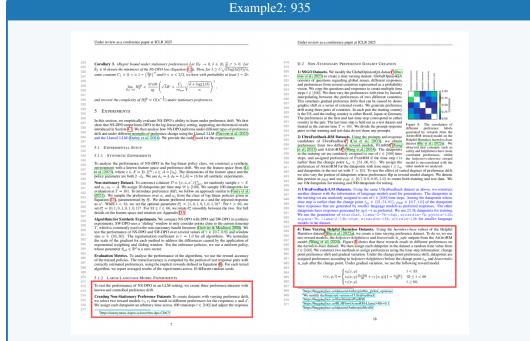
Explanation: A reviewer points out a flaw in the experimental design for the pruning experiments. The paper states on page 9, "We then removed one of the singular value deciles from a specific matrix type in all layers". The reviewer argues this "coarse intervention" constitutes a design flaw because by modifying all layers simultaneously, it becomes impossible to attribute performance changes to specific layers. This confounds the effects, making it difficult to identify where in the model the removed information was critical. This directly undermines the paper's stated goal of "locating information." The design choice violates the conditions for identifying layer-specific contributions, which is an error of type Design & Identifiability.

Error Type: DI (Design & Identifiability)

Type: Within-Sample

## E.3 SG SAMPLING & GENERALIZABILITY





**Question**: Evaluate the Experiments section for Sampling & Generalizability problems.

**Explanation**: The reviewer correctly points out that all experiments are based on synthetic preference drift. The paper's experimental setup, described on pages 7-8 and in Appendix D, involves taking existing datasets (e.g., UltraFeedback) and artificially generating non-stationarity. For instance, preferences are generated by "two different reward models, PAIRRM and ARMORM" (p. 18), and a switch occurs at a predefined change point 'tcp'. Because the core phenomenon being studied—preference drift—is entirely simulated rather than observed organically, the experimental sample does not represent real-world conditions. This limits the generalizability of the paper's findings, as the model's performance on synthetic drift may not translate to its performance on natural, complex preference drift. This is a Sampling & Generalizability (C) issue.

Error Type: SG (Sampling & Generalizability)

**Type**: Within-Sample

## E.4 RCA (REFERENTIAL AND CITATION ALIGNMENT)

## 

**Question**: Scan the errors in cited reference Chen et al.(2021)

**Explanation**: The edited P contains a Type H error by misrepresenting the performance of the cited model. P (p. 8) claims that the NSTPP model from Chen et al. (2021) 'reported performance comparable to a standard Hawkes process baseline'. This contradicts the results in S, where the proposed models (i.e., NSTPP) consistently outperform the Hawkes process baseline, often by a large margin. For example, S (p. 9, Table 1) shows on the BOLD5000 dataset that the 'Attentive CNF' model achieves a temporal log-likelihood of  $5.842 \pm 0.005$ , which is substantially better than the Hawkes Process at  $2.860 \pm 0.050$ .

**Error Type**: RCA (Referential and Citation Alignment)

Type: Cross-Generate

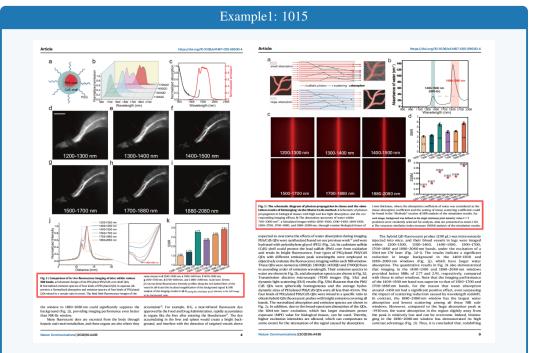
## Example2: 570 **Question**: Evaluate the Results section for Internal Consistency problems.

**Explanation**: A reviewer points out that the paper's reported performance on FB15k-237 is lower than a state-of-the-art method. The paper's main text makes a claim that is directly contradicted by its own table. Specifically, the text states that on the FB15k-237 dataset, their model "still outperforms rule-based and text-based methods" (page 7, 'On the FB15k-237 dataset...methods'). However, Table 1 on the same page presents results for KRACL, a method listed under the "Text-based Methods" category, which achieves higher scores than the proposed model on both MRR (36.0 vs. 35.5) and Hit@1 (26.6 vs. 26.4). This discrepancy between the narrative claim and the tabular data constitutes a clear internal consistency error.

Error Type: RCA (Referential and Citation Alignment)

**Type**: Within-Sample 3271

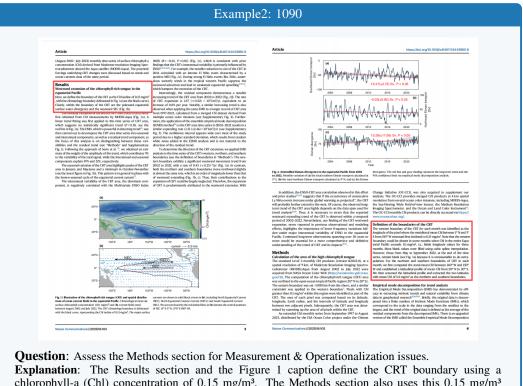
## E.5 MO (MEASUREMENT & OPERATIONALIZATION)



**Question**: Assess the Figures/Tables section for Measurement & Operationalization issues. **Explanation**: The figure captions on pages 3 and 4 have been edited to specify that the background for Signal-to-Background Ratio (SBR) calculations was defined as the single minimum pixel intensity in the image. This is not a valid or reliable operationalization of the "background" construct, as it's highly susceptible to single-point noise or detector artifacts. This flawed measurement procedure systematically undermines all conclusions based on the SBR metric.

Error Type: MO (Measurement & Operationalization)

Type: Within-Generate



**Explanation**: The Results section and the Figure 1 caption define the CRT boundary using a chlorophyll-a (Chl) concentration of 0.15 mg/m³. The Methods section also uses this 0.15 mg/m³ threshold for the western boundary. However, the same Methods section then defines the northern and southern boundaries using a different threshold of 0.1 mg/m³, creating an inconsistent operational definition for the paper's primary construct.

Error Type: MO (Measurement & Operationalization)

**Type**: Within-Generate

## E.6 DHP (DATA HANDLING & PREPROCESSING)

# Under review as a conference paper at ICLR 2005 The likelihood that the comput configuration will usuity the desired constraints. In our case, the configuration of the likelihood that the comput configuration will usuity the desired constraints. In our case, the configuration of the configuration of

Question: Assess the Methods section for Data Handling & Preprocessing issues.

**Explanation**: The reviewer correctly identifies that the authors tuned hyperparameters on the test set. The paper's "Implementation Details" section on page 5 states: "For hyperparameter tuning, we employed Bayesian optimization with the wandb sweep tool (Biewald, 2020), aiming to minimize MPJPE for the S9 and S11 in the H36M dataset and PA-MPJPE for the S8 in the H3WB dataset, following the convention of prior works." According to standard protocols for the H36M dataset, subjects S9 and S11 constitute the test set. Tuning hyperparameters directly on the test set introduces data leakage, leading to an optimistic bias in the reported results and invalidating claims of generalization. This is a critical violation of machine learning best practices and fits the Data Handling & Preprocessing (E) category, as a pipeline choice introduces bias.

**Error Type**: DHP (Data Handling & Preprocessing)

Type: Within-Sample

## Example2: 1566 **Question**: Assess the Methods section for Data Handling & Preprocessing issues. Explanation: The modified text on page 3 states that data for the COVID-19 lockdown period were imputed using pre-pandemic averages. This data handling choice is highly problematic, as it smooths over a major, non-random structural break in the time series rather than modeling or excluding it. The imputation method introduces significant bias and data leakage, as a simple average does not accurately reflect the known, drastic reduction in elective surgeries during that specific period, compromising the validity of the causal model. Error Type: DHP (Data Handling & Preprocessing) **Type**: Within-Generate

## E.7 CF (COMPUTATION & FORMULAE)

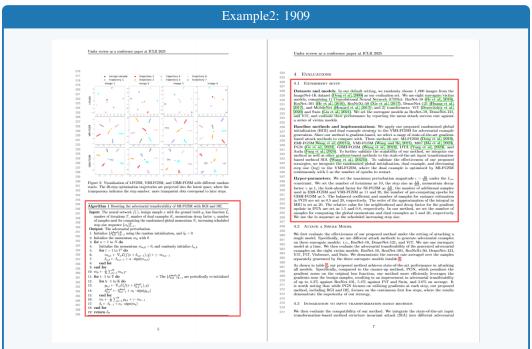
# Under review in a conference paper at RTLE 2025 Under review in a conference paper at RTLE 2025 Under review in a conference paper at RTLE 2025 Figure 2. The conference of the conference of

Question: Check the Methods section for Computation & Formulae errors.

**Explanation**: The Problem Definition section introduces core variables for the mathematical setup, including N and M for the sizes of the labeled and unlabeled datasets. On page 4, line 183, these are defined with the sentence: "And N or M is the total number of image samples." This statement is ambiguous and fails to clearly define N and M individually. A reader cannot determine from this phrase that N is the number of labeled samples and M is the number of unlabeled samples. This notational ambiguity in the definition of variables that are fundamental to the subsequent equations and problem formulation constitutes a Computation & Formulae error, as key variables are left undefined or poorly defined.

**Error Type**: CF (Computation & Formulae)

**Type**: Within-Generate



Question: Scan the Methods section for Computation & Formulae errors.

**Explanation**: Algorithm 1 on page 5 uses the parameter T' in the loop definition on line 5: for t=1 to T' do. This parameter determines the number of iterations for the Randomized Global Initialization phase. However, the value of T' is never specified anywhere in the paper, including the "Hyper-parameters" section (Section 4.1 on page 7). An algorithm cannot be implemented or reproduced with an undefined critical parameter. This fits the Computation & Formulae category as an "undefined variable".

Error Type: CF (Computation & Formulae)

**Type**: Within-Sample

## E.8 IC (Inference & Conclusions)

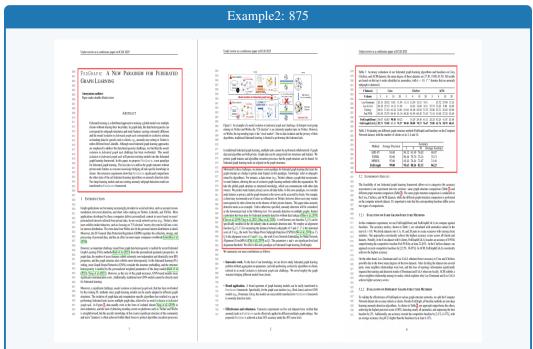


Question: Assess the Discussion section for Inference & Conclusions issues.

**Explanation**: The paper's evidence is based entirely on preclinical models (simulations, mice, rabbits, ex vivo porcine tissue). The edits in the Abstract and Discussion make strong, unhedged claims about setting a "new standard for clinical bioimaging" that is "ready for immediate adoption" in "human surgery." These conclusions are a gross overstatement, as the preclinical data do not support such direct and immediate claims of clinical efficacy and adoption.

**Error Type**: IC (Inference & Conclusions)

Type: Witin-Generate



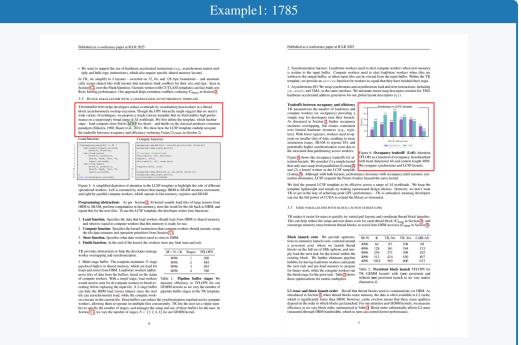
**Question**: Evaluate Abstract, Introduction and Experiment section for issues in Inference & Conclusions.

**Explanation**: The paper's claims of generality are not supported by its evidence. The title and abstract introduce "FedGraph: A New Paradigm for Federated Graph Learning" (page 1), suggesting a broadly applicable framework. However, the methodology is heavily tailored to, and the experiments are exclusively focused on, the single downstream task of anomaly detection. For example, a stated contribution is "Broad application," but this is immediately qualified with "the models are successfully transferred to FEDGRAPH framework in anomaly detection tasks" (page 2). Furthermore, Section 5, "EXPERIMENTS", exclusively reports results on anomaly detection tasks. This discrepancy represents an issue of Inference & Conclusions, as the broad conclusion of having created a new "paradigm" for FGL is an overstatement that exceeds what the narrow experimental results can support.

Error Type: IC (Inference & Conclusions)

**Type**: Within-Sample

## E.9 LE (LANGUAGE & EXPRESSION)

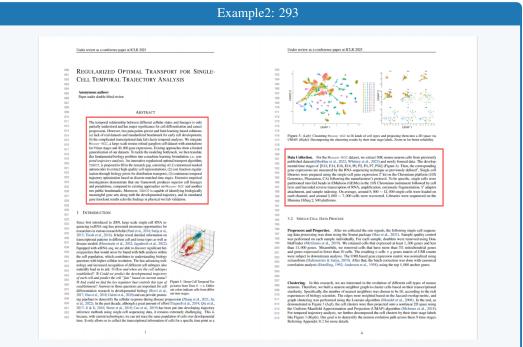


Question: Assess the Methods section for Language & Expression issues.

**Explanation**: The paper introduces a key contribution, the 'load-compute-store-finish' template, and its acronym 'LCSF'. This error introduces inconsistencies in this critical term: it's defined as 'LCS-F' on page 6, called 'LCFS' in a figure title on page 7, and written out in full in the conclusion on page 10, while the original 'LCSF' acronym remains elsewhere. This terminological inconsistency for a central, paper-defined concept creates ambiguity and undermines the paper's precision.

**Error Type**: LE (Language & Expression)

Type: Within-Generate



**Question**: Review the Abstract and Methods sections for Language & Expression problems. **Explanation**: The paper uses the phrase "gene expressions" ambiguously, creating confusion about the size and composition of the dataset. The abstract mentions integrating a dataset with "30, 000 gene expressions" (page 1, line 016), which is repeated on page 2 (contains 30, 000 gene expressions). This phrasing could be misinterpreted as 30,000 unique genes being measured. The Data Collection section later clarifies that the dataset actually consists of "30K mouse neuron cells" (page 4, line 181). This inconsistent terminology affects the meaning of a critical domain quantity (the number of samples), making it a Language & Expression error.

**Error Type**: LE (Language & Expression)

Type: Within-Sample

## F HUMAN-MACHINE CONSISTENCY EVALUATION

As described in Section 4.1, we employ GPT-4.1 to extract detailed information (e.g. evidence sets, reasoning chains) from the responses generated by the models under evaluation. Subsequently, based on the formulas presented in Section 4.1, we calculate  $S_{\rm location}$  and  $S_{\rm reasoning}$ , which are then used to derive  $S_{\rm total}$  for each model's response to the given question.

To evaluate whether GPT-4.1 accurately extracts detailed information from the model responses, we conduct a human-Machine consistency evaluation. We first randomly sampled 200 questions from the dataset. Then, we invited human experts to analyze the corresponding model-generated responses for these questions and to manually extract key information, including evidence sets, reasoning chains, and the number of unrelated errors.

	Stotal	$S_{ m location}$	$S_{ m reasoning}$	$P_{\text{unrelated\_err}}$
Spearman's correlation coefficients	0.841	0.806	0.842	0.954

Table 4: Spearman's correlation coefficients for:  $S_{\text{total}}$ ,  $S_{\text{location}}$ ,  $S_{\text{reasoning}}$ , and  $P_{\text{unrelated\_err.}}$ 

Using the information extracted by the human experts, we perform the following calculations:

- (1) The  $\vec{S}_{\text{location}}$  vector for the 200 questions is calculated based on the evidence sets and Equation 3.
- (2) The  $\vec{S}_{\text{reasoning}}$  vector is computed from the reasoning chains and Equation 4.
- (3) The  $\vec{P}_{\text{unrelated\_err}}$  vector is obtained from the count of unrelated errors.
- (4) The  $\vec{S}_{\text{total}}$  vector is calculated for the 200 questions using Equation 6.

Subsequently, these human-derived vectors ( $\vec{S}_{\text{location}}$ ,  $\vec{S}_{\text{reasoning}}$ ,  $\vec{P}_{\text{unrelated.err}}$ , and  $\vec{S}_{\text{total}}$ ) are compared against their counterparts generated by GPT-4.1. Spearman's correlation coefficient is then calculated for these four metrics. The results are presented in Table 4.

Among the four Spearman correlation coefficients, the metric  $P_{\text{unrelated\_err}}$  exhibits the highest correlation. This indicates that GPT-4.1's extraction of unrelated errors closely aligns with that of human experts, making it the most precise among the three types of extracted information(*i.e.* evidence sets, reasoning chains, and unrelated errors).

Although the correlation coefficients for the *evidence location score* and *reasoning process score* are relatively lower than  $P_{\text{unrelated.err}}$ , they still fall within the range of strong positive correlation. This demonstrates a high degree of consistency in the numerical trends of the scores calculated from GPT-4.1 and human expert extractions, respectively, proving that GPT-4.1 is capable of extracting the majority of effective evidence sets and reasoning chains.

The correlation for the *total score* also lies within the strong positive range and slightly surpasses the correlations for the evidence location score. This also reflects a high level of agreement between GPT-4.1 and human experts.

In summary, GPT-4.1 can extract relevant evidence and reasoning steps with considerable accuracy, leading to precise evaluation scores. This validates the effectiveness of our methodology, which uses GPT-4.1 to parse the responses of the models under evaluation.