

---

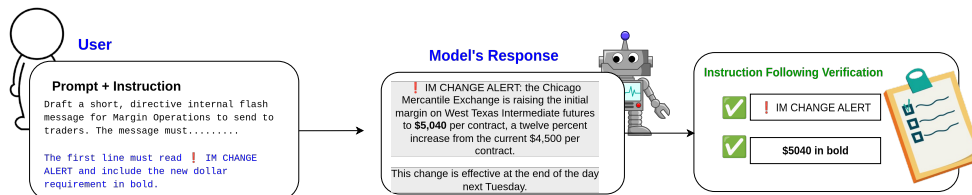
# Instruction Following for Finance: Verifying language models' ability to follow complex financial instructions

---

Anonymous Author(s)  
Affiliation, Address  
anon.email@example.org

## Abstract

1 Language Models (LMs) demonstrate an impressive ability to follow instructions,  
2 but the risk of hallucination when executing complex, interdependent commands  
3 has limited their effectiveness in domains like finance where precision is critical.  
4 We introduce IFF, a high-difficulty benchmark designed to assess the instruction-  
5 following capabilities of LMs for finance. IFF provides 88 human-authored  
6 prompts that mirror financial analysis tasks and uses a verification system with  
7 chainable, verifiable constraints to provide fine-grained reward signals. We evaluate  
8 53 models in a zero-shot setting, including leading proprietary, open-weight,  
9 and open-source systems. Our key findings reveal that open-weight models can  
10 meet or surpass the instruction-following capabilities of proprietary systems.  
11 However, even the top-performing models fail to achieve perfect compliance and  
12 struggle with the IFF benchmark's complex requirements. We release our dataset  
13 and code as an open-source resource to promote research into Reinforcement  
14 Learning with Verified Rewards (RLVR) for the financial domain <sup>1</sup>.



15 Figure 1: Instruction Following Pipeline.

## 16 1 Introduction

17 The value of Language Models (LM) hinges on their ability to follow instructions precisely while  
18 generating responses [1], [2]. While instruction tuning has improved model obedience [3], [4], recent  
19 benchmarks reveal a critical gap: LMs still fail to adhere to long or complex instructions that require  
20 decomposing interdependent constraints [5], [6], [7].


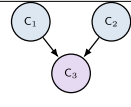
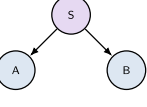
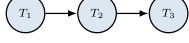
21 This gap is particularly acute in high-stakes domains like finance, where workflows combine domain  
22 reasoning with strict, machine-checkable constraints; e.g., "Create a table with exact headers, sort  
23 by spread descending, and bold values above a threshold." In this context, a factually correct answer  
24 that violates formatting instructions can cause operational failures. While existing Financial NLP  
25 (FinNLP) benchmarks measure task-specific accuracy for tasks like sentiment analysis or QA [8],

---

<sup>1</sup>The code and data are available for review at: [https://anonymous.4open.science/r/iff\\_anon-443F/](https://anonymous.4open.science/r/iff_anon-443F/)

26 [9], [10], [11], they do not evaluate a model’s ability to adhere to the complex, multi-part constraints  
 27 common in financial queries.

28 To address this gap, we introduce Instruction Following for Finance (IFF), a benchmark that stress-  
 29 tests LMs using 88 human-authored tasks that mirror realistic financial workflows. The benchmark’s  
 30 instructions are designed to be verifiable, covering a wide range of compositional and stylistic  
 31 constraints such as table structures, conditional formatting, and domain-specific data normalization.  
 32 Using this framework, we evaluate 53 leading proprietary and open-source LMs in a zero-shot  
 33 setting, with compliance scored automatically via regex-based verification. We release the entire  
 34 benchmark, including all prompts, verifiers, and evaluation code as open-source tool to support future  
 35 research into verifiable-reward training for high-stakes domains.

Composition Type	Description	Example	Illustration
Single	The output satisfies a single constraint or requirement	Write a 200-word end-of-day crypto market recap. Include Bitcoin’s seven-day realized volatility decline from 45% to 32%.	
And	The output satisfies multiple constraints simultaneously	Write 250 words on Meta’s Q3 setup. Include upside levers <b>AND</b> downside watch items <b>AND</b> begin with bold sentence <b>AND</b> add specific headings.	
Selection	The output selects from multiple options based on specified criteria	Update employee trading restricted list. Include only items marked <b>ADD</b> (exclude <b>WATCHLIST</b> and <b>REMOVE</b> ).	
Chain	The output completes sequential tasks, each building on the previous	Finance Committee packet: (1) explain budget gap → (2) create table → (3) show bold gap → (4) add bullets → (5) end with source line.	

36 Figure 2: Composition types in IFF.

## 37 2 Related work

### 38 2.1 Evaluating instruction-following

39 While instruction tuning methods like InstructGPT [4], Self-Instruct [12], and WizardLM [13]  
 40 have improved model obedience, evaluating this capability requires more than task-level accuracy.  
 41 Assessing instruction fidelity demands specialized benchmarks. Frameworks like IFEval [5] and its  
 42 successors ([14], [15], [16]) were created for this purpose, introducing verifiable, machine-checkable  
 43 constraints on length, formatting, and keyword inclusion. However, these evaluations are largely  
 44 domain-agnostic and do not test the specialized constraints required in financial workflows.

### 45 2.2 Financial NLP benchmarks

46 Existing FinNLP benchmarks primarily focus on core NLP tasks including sentiment analysis  
 47 (FOMC [17]), numerical question answering (FinQA, [18]), and text summarization (SubjECTive-  
 48 QA, [19]). These benchmarks assess the labeling accuracy of LMs on FinNLP tasks, supporting the  
 49 development of domain-adapted models. However, none of these benchmarks address the paper’s  
 50 main concern, which is instruction-following. They lack mechanisms to verify compliance with the  
 51 procedural and formatting constraints common in financial workflows, such as producing tables  
 52 with exact headers, reporting numerical results to a mandated precision, or including regulatory  
 53 disclaimers. Consequently, they cannot determine if a model’s output is both factually correct **and**  
 54 compliant with user constraints. IFF is designed to fill this critical gap.

## 55 3 Methodology

### 56 3.1 Instruction-following task

57 Our methodology, Instruction Following for Finance (IFF), is the first instruction-following bench-  
 58 mark designed specifically for finance. We created IFF to evaluate this fundamental capability  
 59 in LMs and to provide a robust environment for Reinforcement Learning with Verified Rewards  
 60 (RLVR). The benchmark is built on a set of 88 difficult, SME-curated tasks that replicate real-world  
 61 financial scenarios, prioritizing complexity and breadth over a large volume of simple tasks. Each  
 62 task contains verifiable output constraints and is designed around four composition types: **single**  
 63 (a single instruction), **and** (multiple concurrent instructions), **chained** (sequential instructions), and  
 64 **nested** (conditional instructions) [14]. See Section C for detailed examples.

### 65 3.2 Evaluation procedure

66 We evaluate a wide range of text-only LMs, deferring multimodal evaluation for future work. The  
 67 models studied include leading proprietary, open-weight, and open-source systems; a complete list  
 68 is available in Table 1. All experiments were conducted via cloud APIs, using commercial endpoints  
 69 for proprietary models and services like TogetherAI and Cirrascale for open models.

70 Following the precedent set by IFEval [5] and IFBench [6], we evaluate outputs using both “strict”  
 71 and “loose” criteria. The loose evaluation programmatically cleans model outputs to normalize for  
 72 minor formatting artifacts. Our benchmark instructions test five broad categories of compliance: (1)  
 73 exact cardinality, (2) table schema, (3) ordering/deduplication, (4) style cues, and (5) word limits.

## 74 4 Results

75 Our findings reveal a clear hierarchy in instruction-following capabilities, with open-weight models  
 76 leading the IFF benchmark. As shown in 1, systems like Llama-4 Maverick 17B surpass the strongest  
 77 closed-source models. While proprietary systems remain competitive, they no longer define the  
 78 frontier of performance. Fully open-source models, despite recent progress, continue to lag behind  
 79 both their open-weight and proprietary counterparts.

Model	Strict	Loose	Model	Strict	Loose	Model	Strict	Loose
<b>Closed-source LMs</b>			<b>Open-source LMs</b>			<b>Open-weight LMs</b>		
Claude 3.7 Sonnet	<b>65.9</b>	<b>70.5</b>	OLMo-2 32B In-	<b>45.5</b>	<b>48.9</b>	DeepSeek-v3.1	58.0	59.1
GPT-4o	63.6	69.3	struct			DeepSeek-v3	56.8	60.2
Claude Opus 4.1	62.5	69.3	Marin 8B Instruct	37.5	38.6	DeepSeek-R1-Distill	55.7	60.2
Claude Opus 4.0	60.2	65.9	OLMo-2 13B In-	31.8	39.8	70B		
(2025-05)			struct			GPT-OSS 120B	53.4	53.4
OpenAI O3	60.2	62.5	OLMo-2 7B Instruct	29.5	36.4	GPT-OSS 20B	53.4	53.4
GPT-4.1	59.1	61.4	OLMo-2 1B Instruct	14.8	18.2	Llama-3 70B Chat-	53.4	63.6
O4-mini	59.1	60.2				HF		
O1	58.0	58.0	<b>Open-weight LMs</b>			DeepSeek-R1-Distill	50.0	53.4
O3-mini	56.8	56.8	Llama-4 Maverick	<b>76.1</b>	<b>79.5</b>	Qwen-14B		
GPT-4o-mini	54.5	58.0	Llama-3.3 70B (free)	71.6	72.7	Mistral-Small 24B	47.7	51.1
Claude 3.5 Haiku	53.4	61.4	Llama-4 Scout 17B	70.5	75.5	(2501)		
(2024-10)			Llama-3.3 70B	70.5	71.6	Cogito v2 109B-	45.5	62.5
GPT-5 mini	53.4	53.4	Llama-3.1 405B	67.0	72.7	MoE (prev)		
Claude 3 Haiku	51.1	54.5	Cogito v2 70B	69.3	73.9	Llama-3.1 8B	44.3	48.9
(2024-03)			(prev)			Gemma-3N e4b-IT	42.0	45.5
GPT-4.1 nano	47.7	52.3	Cogito v2 405B	64.8	67.0	Cogito v2 DeepSeek-	40.9	54.5
GPT-5	47.7	50.0	(prev)			671B (prev)		
GPT-4.1 mini	46.6	51.1	Qwen-QWQ 32B	63.6	64.8	Llama-3 8B Lite	39.8	47.7
Claude Sonnet 4	45.5	52.3	Qwen-3 235B (tput)	61.4	67.0	Qwen-2.5 7B	35.2	37.5
(2025-05)			Moonshot Kimi K2	60.2	63.6	Mistral-7B v0.3	28.4	34.1
GPT-5 nano	15.9	15.9	Qwen-2.5 72B	59.1	63.6	Llama-3.2 3B	28.4	29.5
			DeepSeek-R1 (0528)	59.1	60.2	Mistral-7B v0.2	27.3	34.1
						Mistral-7B v0.1	19.3	21.6

80

Table 1: Model performance on the IFF Benchmark

81 A surprising result is the relatively weak performance of several models explicitly designed for  
82 complex reasoning. Despite their focus on multi-step or conversational tasks, their reliability in our  
83 strict evaluation is not markedly stronger than general-purpose models, and in some cases, it is worse.  
84 This suggests that current reasoning-oriented tuning does not necessarily translate to more precise  
85 instruction following.

86 Across all model classes, we find that failures often stem from an inability to adhere to precise  
87 formatting or structural requirements. The significant gap between “strict” and “loose” scores, par-  
88 ticularly for open-source models, indicates that many failures are “near misses” where the semantic  
89 content is correct but a specific constraint is violated.

## 90 **5 Limitations**

91 Our evaluation has two primary limitations. First, we rely on a single response per prompt due to API  
92 cost, which limits robustness against stochastic model variation. Second, our regex-based validators,  
93 while transparent and reproducible, can be brittle. They may occasionally fail to recognize a seman-  
94 tically correct response due to minor formatting deviations. More flexible validation methods, such  
95 as LM-as-a-judge or hybrid semantic checks, are a natural next step for future work.

## 96 **6 Conclusion**

97 In this study, we introduced the IFF Benchmark to assess instruction-following reliability for  
98 financial applications. While proprietary systems remain competitive, they are no longer dominant,  
99 with top-performing open-weight models now leading in instruction-following capabilities; fully  
100 open-source models, however, continue to lag behind. We find that reasoning-oriented models and  
101 smaller checkpoints surprisingly underperform, suggesting that scaling or explicit reasoning design  
102 alone does not ensure robustness. The observed gap between strict and loose scoring indicates that  
103 many failures are “near-misses”, outputs that are semantically plausible but procedurally flawed. This  
104 highlights that for models to be trusted in finance, the standard for success must be precise, verifiable  
105 compliance, not just plausible generation.

106 **References**

- 107 [1] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language Models are Unsupervised  
108 Multitask Learners.”
- 109 [2] H. Touvron *et al.*, “LLaMA: Open and Efficient Foundation Language Models.” Accessed: Sep. 01, 2025.  
110 [Online]. Available: <http://arxiv.org/abs/2302.13971>
- 111 [3] J. Wei *et al.*, “Chain-of-Thought Prompting Elicits Reasoning in Large Language Models.” Accessed:  
112 Sep. 01, 2025. [Online]. Available: <http://arxiv.org/abs/2201.11903>
- 113 [4] L. Ouyang *et al.*, “Training language models to follow instructions with human feedback.” Accessed: Sep.  
114 01, 2025. [Online]. Available: <http://arxiv.org/abs/2203.02155>
- 115 [5] J. Zhou *et al.*, “Instruction-Following Evaluation for Large Language Models.” Accessed: Aug. 31, 2025.  
116 [Online]. Available: <http://arxiv.org/abs/2311.07911>
- 117 [6] V. Pyatkin *et al.*, “Generalizing Verifiable Instruction Following.” Accessed: Aug. 26, 2025. [Online].  
118 Available: <http://arxiv.org/abs/2507.02833>
- 119 [7] H. Peng, Y. Qi, X. Wang, B. Xu, L. Hou, and J. Li, “VerIF: Verification Engineering for Reinforcement  
120 Learning in Instruction Following.” Accessed: Aug. 13, 2025. [Online]. Available: <http://arxiv.org/abs/2506.09942>  
121
- 122 [8] R. S. Shah *et al.*, “WHEN FLUE MEETS FLANG: Benchmarks and Large Pre-trained Language Model  
123 for Financial Domain.” Accessed: Sep. 01, 2025. [Online]. Available: <http://arxiv.org/abs/2211.00083>
- 124 [9] X. Wu *et al.*, “Golden Touchstone: A Comprehensive Bilingual Benchmark for Evaluating Financial Large  
125 Language Models.” Accessed: Sep. 01, 2025. [Online]. Available: <http://arxiv.org/abs/2411.06272>
- 126 [10] G. Lu, X. Guo, R. Zhang, W. Zhu, and J. Liu, “BizFinBench: A Business-Driven Real-World Financial  
127 Benchmark for Evaluating LLMs.” Accessed: Sep. 01, 2025. [Online]. Available: <http://arxiv.org/abs/2505.19457>  
128
- 129 [11] G. Matlin, M. Okamoto, H. Pardawala, Y. Yang, and S. Chava, “Finance Language Model Evaluation  
130 (FLaME).” Accessed: Sep. 01, 2025. [Online]. Available: <http://arxiv.org/abs/2506.15846>
- 131 [12] Y. Wang *et al.*, “Self-Instruct: Aligning Language Models with Self-Generated Instructions,” May 2023.  
132 Accessed: Sep. 01, 2025. [Online]. Available: <http://arxiv.org/abs/2212.10560>
- 133 [13] C. Xu *et al.*, “WizardLM: Empowering large pre-trained language models to follow complex instructions,”  
134 May 2025. Accessed: Aug. 31, 2025. [Online]. Available: <http://arxiv.org/abs/2304.12244>
- 135 [14] B. Wen *et al.*, “Benchmarking Complex Instruction-Following with Multiple Constraints Composition.”  
136 Accessed: Aug. 29, 2025. [Online]. Available: <http://arxiv.org/abs/2407.03978>
- 137 [15] S. Chen *et al.*, “Recent Advances in Large Language Model Benchmarks against Data Contamination:  
138 From Static to Dynamic Evaluation.” Accessed: Aug. 26, 2025. [Online]. Available: <http://arxiv.org/abs/2502.17521>  
139
- 140 [16] Y. Hwang, Y. Kim, J. Koo, T. Kang, H. Bae, and K. Jung, “LLMs can be easily Confused by Instructional  
141 Distractions,” in *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics*  
142 *(Volume 1: Long Papers)*, W. Che, J. Nabende, E. Shutova, and M. T. Pilehvar, Eds., Vienna, Austria:  
143 Association for Computational Linguistics, Jul. 2025, pp. 19483–19496. doi: 10.18653/v1/2025.acl-  
144 long.957.
- 145 [17] A. Shah, S. Paturi, and S. Chava, “Trillion dollar words: A new financial dataset, task & market analysis,”  
146 May 13, 2023. [Online]. Available: <http://arxiv.org/abs/2305.07972>
- 147 [18] Z. Chen *et al.*, “FinQA: A dataset of numerical reasoning over financial data,” presented at the Proceedings  
148 of the 2021 Conference on Empirical Methods in Natural Language Processing, Stroudsburg, PA, USA:  
149 Association for Computational Linguistics, 2021. doi: 10.18653/v1/2021.emnlp-main.300.
- 150 [19] H. Pardawala *et al.*, “SubjECTive-QA: Measuring subjectivity in Earnings Call Transcripts' QA through  
151 six-dimensional feature analysis,” Oct. 27, 2024. [Online]. Available: <http://arxiv.org/abs/2410.20651>
- 152 [20] Y. Chang *et al.*, “A Survey on Evaluation of Large Language Models,” *ACM Transactions on Intelligent*  
153 *Systems and Technology*, vol. 15, no. 3, pp. 1–45, Jun. 2024, doi: 10.1145/3641289.

- 154 [21] B. Wen *et al.*, “Benchmarking Complex Instruction-Following with Multiple Constraints Composition.”  
155 Accessed: Sep. 01, 2025. [Online]. Available: <http://arxiv.org/abs/2407.03978>
- 156 [22] Q. He *et al.*, “Can Large Language Models Understand Real-World Complex Instructions?,” 2023, doi:  
157 10.48550/ARXIV.2309.09150.
- 158 [23] K. Yan, H. Guo, X. Shi, S. Cao, D. Di, and Z. Li, “CodeIF: Benchmarking the Instruction-Following  
159 Capabilities of Large Language Models for Code Generation.” Accessed: Aug. 29, 2025. [Online].  
160 Available: <http://arxiv.org/abs/2502.19166>
- 161 [24] Z. Zhang, J. Chen, and D. Yang, “DARG: Dynamic Evaluation of Large Language Models via Adaptive  
162 Reasoning Graph.” Accessed: Aug. 26, 2025. [Online]. Available: <http://arxiv.org/abs/2406.17271>
- 163 [25] J. Heo *et al.*, “Do LLMs "know" internally when they follow instructions?,” 2024, doi: 10.48550/  
164 ARXIV.2410.14516.
- 165 [26] J. Heo, M. Xiong, C. Heinze-Deml, and J. Narain, “Do LLMs estimate uncertainty well in instruction-  
166 following?,” 2024, doi: 10.48550/ARXIV.2410.14582.
- 167 [27] V. Reddy, R. Koncel-Kedziorski, V. D. Lai, and C. Tanner, “DocFinQA: A long-context financial reason-  
168 ing dataset,” Jan. 12, 2024. [Online]. Available: <http://arxiv.org/abs/2401.06915>
- 169 [28] W. Hu, J. Duan, C. Wei, L. Zhang, Y. Zhang, and K. Xu, “DynaCode: A Dynamic Complexity-Aware  
170 Code Benchmark for Evaluating Large Language Models in Code Generation.” Accessed: Aug. 26, 2025.  
171 [Online]. Available: <http://arxiv.org/abs/2503.10452>
- 172 [29] K. Zhu, J. Wang, Q. Zhao, R. Xu, and X. Xie, “Dynamic Evaluation of Large Language Models by Meta  
173 Probing Agents.” Accessed: Aug. 26, 2025. [Online]. Available: <http://arxiv.org/abs/2402.14865>
- 174 [30] K. Zhu, J. Chen, J. Wang, N. Z. Gong, D. Yang, and X. Xie, “DyVal: Dynamic Evaluation of Large  
175 Language Models for Reasoning Tasks.” Accessed: Aug. 26, 2025. [Online]. Available: <http://arxiv.org/abs/2309.17167>
- 176
- 177 [31] Y. Perlitz *et al.*, “Efficient Benchmarking of Language Models,” arXiv, 2023. doi: 10.48550/  
178 ARXIV.2308.11696.
- 179 [32] T. Zou, X. Zhang, H. Yu, M. Wang, F. Huang, and Y. Li, “EIFBENCH: Extremely Complex Instruction  
180 Following Benchmark for Large Language Models.” Accessed: Aug. 29, 2025. [Online]. Available: <http://arxiv.org/abs/2506.08375>
- 181
- 182 [33] Z. Yuan, J. Liu, Q. Zi, M. Liu, X. Peng, and Y. Lou, “Evaluating Instruction-Tuned Large Language  
183 Models on Code Comprehension and Generation,” 2023, doi: 10.48550/ARXIV.2308.01240.
- 184 [34] K. Kamble, M. Russak, D. Mozolevskiy, M. Ali, M. Russak, and W. AlShikh, “Expect the Unexpected:  
185 FailSafe Long Context QA for Finance.” Accessed: Aug. 31, 2025. [Online]. Available: <http://arxiv.org/abs/2502.06329>
- 186
- 187 [35] Z. Yuan *et al.*, “FinLLMs: A Framework for Financial Reasoning Dataset Generation with Large Language  
188 Models.” Accessed: Aug. 29, 2025. [Online]. Available: <http://arxiv.org/abs/2401.10744>
- 189 [36] S. Ye *et al.*, “FLASK: Fine-grained Language Model Evaluation based on Alignment Skill Sets,” 2023,  
190 doi: 10.48550/ARXIV.2307.10928.
- 191 [37] V. Pyatkin *et al.*, “Generalizing Verifiable Instruction Following.” Accessed: Sep. 01, 2025. [Online].  
192 Available: <http://arxiv.org/abs/2507.02833>
- 193 [38] L. Diao, X. Xu, W. Sun, C. Yang, and Z. Zhang, “GuideBench: Benchmarking Domain-Oriented Guideline  
194 Following for LLM Agents,” 2025, doi: 10.48550/ARXIV.2505.11368.
- 195 [39] D. Jaroslawicz, B. Whiting, P. Shah, and K. Maamari, “How Many Instructions Can LLMs Follow at  
196 Once?.” Accessed: Aug. 26, 2025. [Online]. Available: <http://arxiv.org/abs/2507.11538>
- 197 [40] Y. K. Chia, P. Hong, L. Bing, and S. Poria, “INSTRUCTEVAL: Towards Holistic Evaluation of Instruc-  
198 tion-Tuned Large Language Models.” Accessed: Aug. 31, 2025. [Online]. Available: <http://arxiv.org/abs/2306.04757>
- 199
- 200 [41] S. Li *et al.*, “Instruction-following Evaluation through Verbalizer Manipulation,” arXiv, 2023. doi:  
201 10.48550/ARXIV.2307.10558.
- 202 [42] Y. Jiang *et al.*, “Instruction-Tuning Data Synthesis from Scratch via Web Reconstruction.” Accessed: Aug.  
203 26, 2025. [Online]. Available: <http://arxiv.org/abs/2504.15573>

- 204 [43] T. Kim, J. Kim, G. Lee, and S.-y. Yun, “Instructive Decoding: Instruction-Tuned Large Language  
205 Models are Self-Refiner from Noisy Instructions,” presented at the International Conference on Learning  
206 Representations, Nov. 2023. Accessed: Aug. 13, 2025. [Online]. Available: <https://www.semanticscholar.org/paper/Instructive-Decoding%3A-Instruction-Tuned-Large-are-Kim-Kim/c0d698950a4560fc2a63acb30a91aa2deb042ed3>  
207  
208
- 209 [44] R. Murthy, P. Venkateswaran, P. Kumar, and D. Contractor, “KCIF: Knowledge-Conditioned Instruction  
210 Following,” May 2025. Accessed: Aug. 31, 2025. [Online]. Available: <http://arxiv.org/abs/2410.12972>
- 211 [45] Y. Hwang, Y. Kim, J. Koo, T. Kang, H. Bae, and K. Jung, “LLMs can be easily Confused by Instructional  
212 Distractions.”
- 213 [46] Y. He *et al.*, “Multi-IF: Benchmarking LLMs on Multi-Turn and Multilingual Instructions Following.”  
214 Accessed: Aug. 26, 2025. [Online]. Available: <http://arxiv.org/abs/2410.15553>
- 215 [47] J. Huang *et al.*, “Open-FinLLMs: Open Multimodal Large Language Models for Financial Applications.”  
216 Accessed: Sep. 01, 2025. [Online]. Available: <http://arxiv.org/abs/2408.11878>
- 217 [48] J. Zeng *et al.*, “Order Matters: Investigate the Position Bias in Multi-constraint Instruction Following,” in  
218 *Findings of the Association for Computational Linguistics: ACL 2025*, W. Che, J. Nabende, E. Shutova,  
219 and M. T. Pilehvar, Eds., Vienna, Austria: Association for Computational Linguistics, Jul. 2025, pp.  
220 12479–12492. doi: 10.18653/v1/2025.findings-acl.646.
- 221 [49] J. Zeng *et al.*, “Order Matters: Investigate the Position Bias in Multi-constraint Instruction Following.”  
222 Accessed: Aug. 29, 2025. [Online]. Available: <http://arxiv.org/abs/2502.17204>
- 223 [50] Q. Xie *et al.*, “PIXIU: A Large Language Model, Instruction Data and Evaluation Benchmark for Finance.”  
224 Accessed: Sep. 01, 2025. [Online]. Available: <http://arxiv.org/abs/2306.05443>
- 225 [51] R. Vir, S. Shankar, H. Chase, W. Fu-Hinthorn, and A. Parameswaran, “PROMPTEVALS: A Dataset of  
226 Assertions and Guardrails for Custom Production Large Language Model Pipelines.” Accessed: Aug. 26,  
227 2025. [Online]. Available: <http://arxiv.org/abs/2504.14738>
- 228 [52] W. Liu *et al.*, “RECAST: Strengthening LLMs’ Complex Instruction Following with Constraint-Verifiable  
229 Data.” Accessed: Aug. 29, 2025. [Online]. Available: <http://arxiv.org/abs/2505.19030>
- 230 [53] Y. Liu *et al.*, “ReIFE: Re-evaluating Instruction-Following Evaluation.” Accessed: Aug. 31, 2025.  
231 [Online]. Available: <http://arxiv.org/abs/2410.07069>
- 232 [54] Y. Wang *et al.*, “Self-Instruct: Aligning Language Models with Self-Generated Instructions.” Accessed:  
233 Aug. 31, 2025. [Online]. Available: <http://arxiv.org/abs/2212.10560>
- 234 [55] J. Cheng *et al.*, “SPaR: Self-Play with Tree-Search Refinement to Improve Instruction-Following in Large  
235 Language Models,” Mar. 2025. Accessed: Aug. 31, 2025. [Online]. Available: <http://arxiv.org/abs/2412.11605>  
236
- 237 [56] P. Venkateswaran and D. Contractor, “Spotlight Your Instructions: Instruction-following with Dynamic  
238 Attention Steering.” Accessed: Aug. 26, 2025. [Online]. Available: <http://arxiv.org/abs/2505.12025>
- 239 [57] Q. Ren *et al.*, “Step-by-Step Mastery: Enhancing Soft Constraint Following Ability of Large Language  
240 Models.” Accessed: Aug. 29, 2025. [Online]. Available: <http://arxiv.org/abs/2501.04945>
- 241 [58] H. Li *et al.*, “Synthetic Data (Almost) from Scratch: Generalized Instruction Tuning for Language  
242 Models.” Accessed: Aug. 31, 2025. [Online]. Available: <http://arxiv.org/abs/2402.13064>
- 243 [59] M. Zheng *et al.*, “TableDreamer: Progressive and Weakness-guided Data Synthesis from Scratch for  
244 Table Instruction Tuning,” in *Findings of the Association for Computational Linguistics: ACL 2025*, W.  
245 Che, J. Nabende, E. Shutova, and M. T. Pilehvar, Eds., Vienna, Austria: Association for Computational  
246 Linguistics, Jul. 2025, pp. 7290–7315. doi: 10.18653/v1/2025.findings-acl.381.
- 247 [60] S. Kim *et al.*, “The BiGGen Bench: A Principled Benchmark for Fine-grained Evaluation of Language  
248 Models with Language Models,” arXiv, 2024. doi: 10.48550/ARXIV.2406.05761.
- 249 [61] X. Chen *et al.*, “The SIFo Benchmark: Investigating the Sequential Instruction Following Ability of Large  
250 Language Models.” Accessed: Sep. 01, 2025. [Online]. Available: <http://arxiv.org/abs/2406.19999>
- 251 [62] J. Cook, T. Rocktäschel, J. Foerster, D. Aumiller, and A. Wang, “TICKing All the Boxes: Generated  
252 Checklists Improve LLM Evaluation and Generation.” Accessed: Aug. 31, 2025. [Online]. Available:  
253 <http://arxiv.org/abs/2410.03608>

254 [63] Y. Cao *et al.*, “Toward Generalizable Evaluation in the LLM Era: A Survey Beyond Benchmarks,” 2025,  
255 doi: 10.48550/ARXIV.2504.18838.

256 **A Appendix: Models**

257 Table 1: Models evaluated in this study. The table lists all proprietary, open-source, and open-weight  
 258 models along with their configuration parameters.

259	<b>Provider</b>	<b>Family</b>	<b>Model</b>	<b>MaxTok</b>	<b>Temp</b>	<b>Top-p</b>
260	<b>Closed-source LMs</b>					
261	OpenAI	O-Series	o4-mini	4096	0.20	0.90
262	OpenAI	O-Series	o3	4096	0.20	0.90
263	OpenAI	O-Series	o3-mini	4096	0.20	0.90
264	OpenAI	O-Series	o1	4096	0.20	0.90
265	OpenAI	GPT-5	gpt-5	4096	1.00	n/a
266	OpenAI	GPT-5	gpt-5-mini	4096	1.00	n/a
267	OpenAI	GPT-5	gpt-5-nano	4096	1.00	n/a
268	OpenAI	GPT-4.1	gpt-4.1	4096	0.20	0.90
269	OpenAI	GPT-4.1	gpt-4.1-mini	4096	0.20	0.90
270	OpenAI	GPT-4.1	gpt-4.1-nano	4096	0.20	0.90
271	OpenAI	GPT-4o	gpt-4o	4096	0.20	0.90
272	OpenAI	GPT-4o	gpt-4o-mini	4096	0.20	0.90
273	Anthropic	Claude Opus	claude-opus-4-1-20250805	4096	n/a	n/a
274	Anthropic	Claude Opus	claude-opus-4-20250514	4096	0.20	0.90
275	Anthropic	Claude Sonnet	claude-sonnet-4-20250514	4096	0.20	0.90
276	Anthropic	Claude Sonnet	claude-3-7-sonnet-20250219	4096	0.20	0.90
277	Anthropic	Claude Haiku	claude-3-5-haiku-20241022	4096	0.20	0.90
278	Anthropic	Claude Haiku	claude-3-haiku-20240307	4096	0.20	0.90
279	<b>Open-source LMs</b>					
280	Cirrascale	OLMo 2	OLMo-2-0325-32B-Instruct	2048	0.20	0.90
281	TogetherAI	Marin	marin-8b-instruct	2048	0.20	0.90
282	Cirrascale	OLMo 2	OLMo-2-1124-13B-Instruct	2048	0.20	0.90
283	Cirrascale	OLMo 2	OLMo-2-1124-7B-Instruct	2048	0.20	0.90
284	Cirrascale	OLMo 2	OLMo-2-0425-1B-Instruct	2048	0.20	0.90
285	<b>Open-weight LMs</b>					
286	Together.ai	Meta Llama 4	Llama-Maverick-17B-128E-Instruct- FP8	4096	0.20	0.90
287						
288	Together.ai	Meta Llama 4	Llama-Scout-17B-16E-Instruct	4096	0.20	0.90
289	Together.ai	Meta Llama 3	Llama-3.3-70B-Instruct	4096	0.20	0.90
290	Together.ai	Meta Llama 3	Llama-3.1-405B-Instruct	4096	0.20	0.90
291	Together.ai	Meta Llama 3	Llama-3.3-70B-Instruct	4096	0.20	0.90
292	Together.ai	Deep Cogito	cogito-v2-preview-llama-70B	4096	0.20	0.90
293	Together.ai	Deep Cogito	cogito-v2-preview-llama-405B	4096	0.20	0.90
294	Together.ai	Qwen	QwQ-32B	4096	0.20	0.90
295	Together.ai	Qwen	Qwen3-235B-A22B-Instruct-2507	4096	0.20	0.90
296	Together.ai	Moonshot	Kimi-K2-Instruct	4096	0.20	0.90
297	Together.ai	Qwen	Qwen2.5-72B-Instruct	4096	0.20	0.90
298	Together.ai	DeepSeek	DeepSeek-R1-0528	4096	0.20	0.90
299	Together.ai	DeepSeek	DeepSeek-V3.1	4096	0.20	0.90
300	Together.ai	DeepSeek	DeepSeek-V3	4096	0.20	0.90
301	Together.ai	DeepSeek	DeepSeek-R1-Distill-Llama-70B	4096	0.20	0.90
302	Together.ai	GPT-OSS	gpt-oss-120b	4096	0.20	0.90
303	Together.ai	GPT-OSS	gpt-oss-20b	4096	0.20	0.90
304	Together.ai	Meta Llama 3	Llama-3-70b-chat-hf	4096	0.20	0.90
305	Together.ai	DeepSeek	DeepSeek-R1-Distill-Qwen-14B	4096	0.20	0.90
306	Together.ai	Mistral AI	Mistral-Small-24B-Instruct-2501	4096	0.20	0.90
307	Together.ai	Deep Cogito	cogito-v2-preview-llama-109B-MoE	4096	0.20	0.90
308	Together.ai	Meta Llama 3.1	Meta-Llama-3.1-8B-Instruct	4096	0.20	0.90
309	Together.ai	Google Gemma	gemma-3n-E4B-it	4096	0.20	0.90
310	Together.ai	Deep Cogito	cogito-v2-preview-deepseek-671b	4096	0.20	0.90
311	Together.ai	Meta Llama 3	Meta-Llama-3-8B-Instruct-Lite	4096	0.20	0.90
312	Together.ai	Qwen 2.5	Qwen2.5-7B-Instruct	4096	0.20	0.90
313	Together.ai	Mistral AI	Mistral-7B-Instruct-v0.3	4096	0.20	0.90
314	Together.ai	Meta Llama 3.2	Llama-3.2-3B-Instruct	4096	0.20	0.90
315	Together.ai	Mistral AI	Mistral-7B-Instruct-v0.2	4096	0.20	0.90
316	Together.ai	Mistral AI	Mistral-7B-Instruct-v0.1	4096	0.20	0.90

## 317 **B Appendix: Evaluation Framework**

### 318 **B.1 Architecture**

319 We designed the IFF framework to use a modular, layered architecture to support maintainability  
320 and reproducibility. The application layer exposes a user interface through command-line scripts.  
321 The primary entry points are `evaluation_bin.py` for running evaluations and `generate_responses.py`  
322 for producing model outputs. The business logic layer implements the core functionality. The  
323 `evaluation_lib.py` module orchestrates end-to-end evaluation, and the instruction modules define and  
324 resolve the logic for each verifiable constraint. The utility layer provides shared functions required  
325 by the business logic. The `instructions_util.py` module implements text processing and validation  
326 routines such as word counting, sentence splitting, and table detection that are reused by multiple  
327 checkers. The data layer holds the artifacts consumed and produced by the system, including input  
328 JSONL files with prompts and output JSONL files containing model responses and evaluation  
329 results.

### 330 **B.2 Core Components & Data Flow**

331 Benchmark construction begins with `build_input_jsonl`, which assembles test cases using  
332 `instructions_registry` to combine instruction types and parameters into financial prompts. The  
333 prompts are serialized to an input JSONL file. The `generate_responses` module reads this file and  
334 submits the prompts to external models through a multi-provider gateway. The evaluation phase is  
335 initiated by `evaluation_bin`, which uses `evaluation_lib` to load prompts and their corresponding re-  
336 sponses. For each prompt, `evaluation_lib` requests the appropriate checker from `instructions_registry`  
337 and applies the selected evaluation mode to the model output. The `instructions_registry` maintains  
338 the mapping from instruction identifiers to checker classes and returns the correct implementation  
339 at runtime. The `finance_instructions` module contains the concrete checkers, each inheriting from a  
340 common base and implementing the domain-specific rules required to validate its constraint. The  
341 `instructions_util` module supplies the supporting text-processing functions that are reused across  
342 checkers. The overall data flow is linear and reproducible: prompt generation produces the input  
343 JSONL, response generation produces the responses JSONL, and evaluation produces result files  
344 that record pass/fail outcomes and statistics.

### 345 **B.3 Design Patterns**

346 The registry pattern appears in `instructions_registry`, which centralizes the association between  
347 instruction identifiers and checker implementations so that the evaluation engine does not depend on  
348 concrete classes. The strategy pattern governs evaluation behavior, allowing Strict and Loose modes  
349 to be selected at runtime without changing the engine’s control flow. The factory pattern is used when  
350 instantiating checkers from identifiers in a prompt; `evaluation_lib` requests a class from the registry,  
351 constructs the checker with any provided parameters, and executes validation. This combination  
352 allows new instructions to be added by registering a checker class, without modifying orchestration  
353 code or other components.

## 354 **C Appendix: Evaluation Dataset**

355 This appendix provides detailed information on the prompts and instructions used in our experi-  
356 ments, along with samples of model outputs.

### 357 **C.1 System Prompt**

358 The LLMs are primed using the following system prompt to aid in aligning with our instruction  
359 following and improve performance.

360  
361 System prompt  
362



431 - HEADINGS: Use "#", "##", "###", etc. only if explicitly requested, with the  
432 wording provided.

#### 434 BLANK-LINE RULE

435 - Unless a specific blank-line pattern is requested (e.g., "one blank line between  
436 Block 1 and Block 2"), produce no blank lines. When a pattern is requested, follow  
437 it exactly and avoid any additional blank lines.

### 438 C.2 Sample Prompts

439 Representative examples of prompts are provided to illustrate the variety and structure of inputs used.

#### 440 C.2.1 Single Output Prompts

441

##### 442 Prompt

443 Rates quants want a 150-word explainer of how the SOFR-TONA basis feeds into  
444 cross-currency-swap valuation. TL;DR: the basis spread is added to one leg so both  
445 discount curves align before fair-value logic is applied. Close with the  
446 adjustment formula.

447 Instruction: Begin with exact italicised TL;DR sentence above and end with a  
448 fenced LaTeX formula block showing the basis-spread adjustment.

#### 449 C.2.2 Chained Output Prompts

450

##### 451 Prompt

452 The bank is scheduling a release of the core payments switch identified as release  
453 twenty twenty five point zero nine. The change window runs from Saturday twelve  
454 September two thousand twenty five at twenty two hundred hours until Sunday  
455 thirteen September two thousand twenty five at zero four hundred hours local time.  
456 The update affects instant bank transfers, automated clearing house files, and  
457 card authorisation routing. The contingency plan is to return to the prior  
458 production version within thirty minutes if specific failure indicators appear.  
459 Prepare a cutover runbook that coordinates technology, operations, and  
460 communications. Instruction: Begin with a one sentence goal statement, then set  
461 out the runbook in order: (1) provide a timeline table with columns time,  
462 activity, owner, and success signal from pre change to post change; (2) write  
463 three checklists titled pre cutover, cutover, and post cutover using Markdown  
464 checkboxes for items such as freezes, backups, message drains, smoke tests, and  
465 reconciliation; (3) add a short section titled rollback triggers that lists the  
466 exact indicators that require a return to the prior version; (4) conclude with a  
467 communications tree naming who is notified at start, at checkpoints, and at  
468 completion. Output in Markdown only.

### 469 C.3 Sample Responses

470 This section presents representative outputs of Llama 4 Maverick, including both correct and  
471 incorrect responses, to illustrate model performance for both strict and loose mode.

472 The model produces outputs that fully adhere to the prompt, meeting all specified instructions and  
473 correctly following the intended structure and formatting.

474

##### 475 ##### Prompt

476

477 Write about 250 words on Meta's Q3 setup using only the facts below. Include upside  
478 levers (Reels RPM, tighter opex) and downside watch items (foreign exchange drag,  
479 ads macro). Do not invent data.

480 Facts: revenue \$36.8 billion; EBIT margin 39%.

481 Begin with one sentence in bold stating your overall call. Add the heading "Upside  
482 Watch" followed by exactly three bullets, then the heading "Downside Watch" followed  
483 by exactly three bullets; each bullet must begin with one of these tokens : Increase,  
484 Tighten, Improve, Monitor, Reduce, Shift, Manage. If your opening bold sentence  
485 contains the word "probably". Insert a single bold line that reads "Confidence: Low."

486 immediately before the final line. Finish with one italic line "Downside risk: FX  
487 drag; ads macro."  
488  
489 ##### Response (PASSED)  
490  
491 \*\*Meta is probably well-positioned for Q3 given its historical performance and  
492 current trends.\*\*  
493     Upside Watch  
494     - Improve Reels RPM as the format continues to gain traction.  
495     - Tighten opex by optimizing operational efficiencies.  
496     - Increase ad load without compromising user experience.  
497     Downside Watch  
498     - Monitor foreign exchange fluctuations that could negatively impact revenue.  
499     - Manage ads macro environment that may affect advertiser spending.  
500     - Reduce dependence on a single revenue stream to mitigate potential risks.  
501     \*\*Confidence: Low.\*\*  
502     \*Downside risk: FX drag; ads macro.\*

503 The example below illustrates a case in which the model does not fully adhere to the prompt,  
504 producing outputs that deviate from the specified instructions and formatting requirements. In this  
505 instance, the model fails to follow multiple instructions provided in the prompt.

506 ##### Prompt  
507  
508 Summarize the three deficiencies identified by the model-risk audit in mortgage  
509 prepayment model version 5.2. The deficiencies are: an outdated seasoning curve,  
510 a missing macroeconomic linkage for refinancing incentive, and an inadequate back-  
511 testing window. Present them clearly and concisely, without adding or inventing any  
512 further issues.  
513 List the three deficiencies as bullets and end with the remediation deadline <span  
514 style="color:red">2025-10-31</span> in bold red text.  
515  
516 ##### Response (FAILED)  
517  
518 - The mortgage prepayment model version 5.2 has an outdated seasoning curve.  
519 - The model lacks a macroeconomic linkage for refinancing incentive.  
520 - The back-testing window used is inadequate.  
521 Remediation deadline: <font color="red">\*\*2025-10-31\*\*</font>

523 The example below illustrates another case in which the model does adhere to the prompt, again  
524 producing outputs that match the specified instructions and formatting.

525 ##### Prompt  
526  
527 Create a detailed checklist for the loan-closing team, covering the settlement  
528 of a seven hundred fifty million-dollar leveraged-loan allocation scheduled for  
529 tomorrow on ClearPar. The checklist must be grouped into three sections: Booking,  
530 Know Your Customer (KYC), and Cash Settlement. Each step within the checklist should  
531 be preceded by an empty tick box written as [ ]. After the checklist, add the sign-  
532 off line:  
533 Ops Lead: \_\_\_\_\_  
534 Follow the group headings exactly and append the sign-off line.  
535  
536 ##### Response (Passed)  
537  
538 # Booking  
539  
540 [ ] Verify loan allocation amount of \$750,000,000 in ClearPar system  
541 [ ] Confirm settlement date matches tomorrow's business date  
542 [ ] Validate borrower entity details and loan facility information  
543 [ ] Check trade confirmation details against original allocation notice  
544 [ ] Ensure proper booking codes and GL accounts are assigned  
545 [ ] Verify interest rate, spread, and payment frequency parameters  
546 [ ] Confirm loan maturity date and amortization schedule  
547 [ ] Review and approve any fee structures (commitment, facility, agent fees)  
548 [ ] Validate collateral details and security interest documentation  
549

```

550 [ ] Ensure proper risk weighting and regulatory capital allocation
551 [ ] Confirm booking system integration with downstream systems
552 [ ] Generate and review preliminary booking entries
553
554 # Know Your Customer (KYC)
555
556 [ ] Verify borrower KYC documentation is current and complete
557 [ ] Confirm beneficial ownership information is up to date
558 [ ] Check sanctions screening results for all related parties
559 [ ] Validate anti-money laundering (AML) compliance status
560 [ ] Review politically exposed person (PEP) screening results
561 [ ] Confirm customer due diligence (CDD) documentation
562 [ ] Verify enhanced due diligence (EDD) requirements if applicable
563 [ ] Check regulatory watch list screening results
564 [ ] Validate customer risk rating and approval status
565 [ ] Confirm compliance with know your customer's customer (KYCC) requirements
566 [ ] Review any adverse media or negative news screening results
567 [ ] Ensure all KYC exceptions have proper approvals
568
569 # Cash Settlement
570
571 [ ] Confirm settlement instructions in ClearPar match internal records
572 [ ] Verify correspondent banking relationships for settlement
573 [ ] Check available liquidity and funding arrangements
574 [ ] Validate wire transfer instructions and beneficiary details
575 [ ] Confirm settlement currency and any FX requirements
576 [ ] Review nostro account balances for settlement currency
577 [ ] Verify cut-off times for same-day settlement processing
578 [ ] Check for any settlement holds or restrictions
579 [ ] Confirm backup settlement procedures are in place
580 [ ] Validate settlement reporting requirements
581 [ ] Ensure proper cash management and liquidity monitoring
582 [ ] Prepare settlement confirmation and reporting documentation
583
584 Ops Lead: _____

```

585 Selection prompts are effective because they require the model to generate outputs that fit into a  
586 rigid structure, leaving little room for ambiguity. This ensures we can clearly evaluate whether the  
587 model followed the exact formatting rules and constraints.

```

588
589 ##### Prompt

```

```

590
591 FinOps review for Q2-25 cloud spend: AWS Compute $650 k, Storage $280 k, Network
592 $90 k, Other $120 k (Total $1.14 m); Azure Compute $540 k, Storage $180 k, Network
593 $60 k, Other $70 k (Total $0.85 m); GCP Compute $720 k, Storage $210 k, Network $70
594 k, Other $50 k (Total $1.05 m). Build a table by provider and line item, add a Grand
595 Total row, and flag every quarterly cost above $1 million with Δ.
596 Show the Grand Total row and mark each cost > $1 m with Δ.

```

```

597
598 ##### Response (PASSED)

```

```

599
600 | Provider | Compute | Storage | Network | Other | Total |
601 |-----|-----|-----|-----|-----|-----|
602 | AWS      | $650k  | $280k  | $90k   | $120k  | $1.14m Δ |
603 | Azure   | $540k  | $180k  | $60k   | $70k   | $0.85m |
604 | GCP     | $720k  | $210k  | $70k   | $50k   | $1.05m Δ |
605 | **Grand Total** | **$1.91m Δ** | **$670k** | **$220k** | **$240k** | **$3.04m
606 Δ** |
607

```

## 608 C.4 Instruction Registry

609 This section provides a comprehensive registry of all instructions used.

- 610 1. fin:equities\_bold\_intro\_italic\_risk
- 611 2. fin:credit\_table\_spread\_vs\_carry

612 3. fin:fx\_calc\_codeblock\_limit  
613 4. fin:compliance\_rule10b5\_numbered  
614 5. fin:ops\_settlement\_checklist  
615 6. fin:ir\_six\_bullets\_verb\_buyback  
616 7. fin:treasury\_liquidity\_risk\_section  
617 8. fin:deriv\_black76\_latex\_sigma  
618 9. fin:risk\_var\_numbered\_boldusd  
619 10. fin:pe\_subheaders\_dashes  
620 11. fin:quant\_pseudocode\_comments  
621 12. fin:crypto\_recap\_percent\_italic  
622 13. fin:abs\_table\_then\_comments  
623 14. fin:reit\_underline\_wordlimit  
624 15. fin:structured\_protect\_terms  
625 16. fin:ecb\_timestamp\_boldrates  
626 17. fin:ratings\_three\_numbered  
627 18. fin:pension\_table\_footnote  
628 19. fin:margin\_im\_alert  
629 20. fin:etf\_timed\_checklist  
630 21. fin:fintech\_swot\_licence\_italic  
631 22. fin:cet1\_formula\_block\_end  
632 23. fin:ma\_synergies\_arrows  
633 24. fin:esg\_scopes\_bold  
634 25. fin:tax\_carried\_interest  
635 26. fin:ifrs17\_callout\_two  
636 27. fin:retail\_liquidity\_underline  
637 28. fin:embond\_yield\_boldprice  
638 29. fin:rrp\_drivers\_subs\_bold  
639 30. fin:climate\_grid\_critical  
640 31. fin:syndicate\_table\_boldcover  
641 32. fin:aml\_numbered\_high  
642 33. fin:xccy\_basis\_tldr\_latex  
643 34. fin:clearpar\_groups\_signoff  
644 35. fin:modelrisk\_deadline  
645 36. fin:carbon\_divergences\_two  
646 37. fin:clo\_table\_italic\_shortest  
647 38. fin:vc\_proscons\_underline\_terms  
648 39. fin:ndf\_numbered\_bold\_vals  
649 40. fin:lch\_margin\_timestamp\_code  
650 41. fin:marketing\_irr\_threehashtags  
651 42. fin:project\_keyparties\_boldusd  
652 43. fin:finops\_table\_grandtotal\_warn  
653 44. fin:payments\_iso20022\_qa\_underline  
654 45. fin:munis\_tey\_footnote\_bold  
655 46. fin:index\_rebal\_subject\_bullets  
656 47. fin:correlations\_table\_red  
657 48. fin:scf\_redflags\_callout\_limit  
658 49. fin:sdr\_csv\_italic\_note  
659 50. fin:eqderivs\_gamma\_table\_bold  
660 51. fin:muni\_budget\_gap\_snapshot  
661 52. fin:esg\_csr\_checklist  
662 53. fin:settlement\_runbook\_notes  
663 54. fin:ma\_board\_slide\_apextech  
664 55. fin:rmbs\_poolcard\_stratification  
665 56. fin:fxhedge\_calendar\_novamed  
666 57. fin:airline\_fuel\_exposureemap  
667 58. fin:aml\_case\_northbridge  
668 59. fin:treasury\_monthend\_runsheet  
669 60. fin:rm\_securitization\_factsheet

670 61. fin:pe\_quarterly\_letter\_skeleton  
671 62. fin:claims\_reserve\_rollforward  
672 63. fin:merchant\_chargeback\_pack  
673 64. fin:jetfuel\_hedgeplan\_outline  
674 65. fin:bank\_gap\_snapshot  
675 66. fin:lockup\_monitoring\_plan  
676 67. fin:branchclosure\_oakview\_plan  
677 68. fin:merchant\_onboarding\_arcadia  
678 69. fin:realestate\_valuation\_runbook  
679 70. fin:greenbond\_allocation\_riverton  
680 71. fin:procurement\_cards\_auditpack  
681 72. fin:research\_blueharbor\_compliance  
682 73. fin:payments\_cutover\_runbook  
683 74. fin:venturedebt\_covenant\_monitor  
684 75. fin:payments\_switch\_cutover\_runbook  
685 76. fin:restricted\_list\_adds\_table  
686 77. fin:fund\_profile\_fees\_contact  
687 78. fin:treasury\_accounts\_rationalization  
688 79. fin:sec\_filing\_pack\_zephyr  
689 80. fin:esg\_extract\_orion\_2024  
690 81. fin:credit\_limit\_changes\_summary  
691 82. fin:saas\_amendment\_change\_json  
692 83. fin:pci\_scope\_extract\_2025  
693 84. fin:ap\_payment\_run\_2025\_10\_15  
694 85. fin:tax\_1099\_candidates\_2025  
695 86. fin:saas\_quanta\_commercial\_controls  
696 87. fin:revops\_deals\_hygiene\_q4\_2025  
697 88. fin:ir\_roadshow\_blueharbor\_oct2025

## 698 **NeurIPS Paper Checklist**

### 699 **1. Claims**

700 Question: Do the main claims made in the abstract and introduction accurately reflect the paper’s  
701 contributions and scope?

702 Answer: [YES]

703 Justification: The abstract and introduction claim to introduce a new benchmark, IFF, for  
704 financial instruction following, present a large-scale evaluation of modern LMs, and release the  
705 associated code and data. These contributions are described in detail in the Methodology section  
706 and the results are presented in the Results section. The appendix provides further details on the  
707 models, framework, and dataset.

708 Guidelines:

- 709 • The answer NA means that the abstract and introduction do not include the claims made in  
710 the paper.
- 711 • The abstract and/or introduction should clearly state the claims made, including the contri-  
712 butions made in the paper and important assumptions and limitations. A No or NA answer  
713 to this question will not be perceived well by the reviewers.
- 714 • The claims made should match theoretical and experimental results, and reflect how much  
715 the results can be expected to generalize to other settings.
- 716 • It is fine to include aspirational goals as motivation as long as it is clear that these goals are  
717 not attained by the paper.

### 718 **2. Limitations**

719 Question: Does the paper discuss the limitations of the work performed by the authors?

720 Answer: [YES]

721 Justification: The paper includes a dedicated “Limitations” section. We discuss that our evalu-  
722 ation relies on a single response sample per model, which may not fully account for stochastic  
723 model outputs, and note that this was due to the high API costs of generating multiple responses.

724 Guidelines:

- 725 • The answer NA means that the paper has no limitation while the answer No means that the  
726 paper has limitations, but those are not discussed in the paper.
- 727 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 728 • The paper should point out any strong assumptions and how robust the results are to  
729 violations of these assumptions (e.g., independence assumptions, noiseless settings, model  
730 well-specification, asymptotic approximations only holding locally). The authors should  
731 reflect on how these assumptions might be violated in practice and what the implications  
732 would be.
- 733 • The authors should reflect on the scope of the claims made, e.g., if the approach was only  
734 tested on a few datasets or with a few runs. In general, empirical results often depend on  
735 implicit assumptions, which should be articulated.
- 736 • The authors should reflect on the factors that influence the performance of the approach.  
737 For example, a facial recognition algorithm may perform poorly when image resolution  
738 is low or images are taken in low lighting. Or a speech-to-text system might not be used  
739 reliably to provide closed captions for online lectures because it fails to handle technical  
740 jargon.
- 741 • The authors should discuss the computational efficiency of the proposed algorithms and  
742 how they scale with dataset size.

- 743 • If applicable, the authors should discuss possible limitations of their approach to address  
744 problems of privacy and fairness.
- 745 • While the authors might fear that complete honesty about limitations might be used by  
746 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
747 limitations that aren't acknowledged in the paper. The authors should use their best  
748 judgment and recognize that individual actions in favor of transparency play an important  
749 role in developing norms that preserve the integrity of the community. Reviewers will be  
750 specifically instructed to not penalize honesty concerning limitations.

### 751 3. Theory Assumptions and Proofs

752 Question: For each theoretical result, does the paper provide the full set of assumptions and a  
753 complete (and correct) proof?

754 Answer: [NA]

755 Justification: This paper introduces a new benchmark and reports empirical findings. It does not  
756 include any theoretical results, theorems, or proofs.

757 Guidelines:

- 758 • The answer NA means that the paper does not include theoretical results.
- 759 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
760 referenced.
- 761 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 762 • The proofs can either appear in the main paper or the supplemental material, but if they  
763 appear in the supplemental material, the authors are encouraged to provide a short proof  
764 sketch to provide intuition.
- 765 • Inversely, any informal proof provided in the core of the paper should be complemented by  
766 formal proofs provided in appendix or supplemental material.
- 767 • Theorems and Lemmas that the proof relies upon should be properly referenced.

### 768 4. Experimental Result Reproducibility

769 Question: Does the paper fully disclose all the information needed to reproduce the main  
770 experimental results of the paper to the extent that it affects the main claims and/or conclusions  
771 of the paper (regardless of whether the code and data are provided or not)?

772 Answer: [YES]

773 Justification: We provide detailed descriptions of our evaluation methodology, including the  
774 “strict” and “loose” accuracy metrics. The appendix contains the full list of models evaluated and  
775 their hyperparameters, a description of the evaluation framework, and representative examples  
776 from our benchmark dataset to ensure our results can be reproduced.

777 Guidelines:

- 778 • The answer NA means that the paper does not include experiments.
- 779 • If the paper includes experiments, a No answer to this question will not be perceived well  
780 by the reviewers: Making the paper reproducible is important, regardless of whether the  
781 code and data are provided or not.
- 782 • If the contribution is a dataset and/or model, the authors should describe the steps taken to  
783 make their results reproducible or verifiable.
- 784 • Depending on the contribution, reproducibility can be accomplished in various ways. For  
785 example, if the contribution is a novel architecture, describing the architecture fully might  
786 suffice, or if the contribution is a specific model and empirical evaluation, it may be  
787 necessary to either make it possible for others to replicate the model with the same dataset,  
788 or provide access to the model. In general, releasing code and data is often one good way

789 to accomplish this, but reproducibility can also be provided via detailed instructions for  
790 how to replicate the results, access to a hosted model (e.g., in the case of a large language  
791 model), releasing of a model checkpoint, or other means that are appropriate to the research  
792 performed.

793 • While NeurIPS does not require releasing code, the conference does require all submissions  
794 to provide some reasonable avenue for reproducibility, which may depend on the nature of  
795 the contribution. For example

796 (a) If the contribution is primarily a new algorithm, the paper should make it clear how to  
797 reproduce that algorithm.

798 (b) If the contribution is primarily a new model architecture, the paper should describe the  
799 architecture clearly and fully.

800 (c) If the contribution is a new model (e.g., a large language model), then there should  
801 either be a way to access this model for reproducing the results or a way to reproduce  
802 the model (e.g., with an open-source dataset or instructions for how to construct the  
803 dataset).

804 (d) We recognize that reproducibility may be tricky in some cases, in which case authors  
805 are welcome to describe the particular way they provide for reproducibility. In the case  
806 of closed-source models, it may be that access to the model is limited in some way  
807 (e.g., to registered users), but it should be possible for other researchers to have some  
808 path to reproducing or verifying the results.

## 809 5. Open access to data and code

810 Question: Does the paper provide open access to the data and code, with sufficient instructions  
811 to faithfully reproduce the main experimental results, as described in supplemental material?

812 Answer: [YES]

813 Justification: Yes, we are releasing the complete benchmark, including all prompts, instructions,  
814 and evaluation code, under an open-source license. For the review process, we provide an  
815 anonymized link to the repository. The public URL will be included in the final camera-ready  
816 version of the paper.

817 Guidelines:

818 • The answer NA means that paper does not include experiments requiring code.

819 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/public/  
820 guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.

821 • While we encourage the release of code and data, we understand that this might not be  
822 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not includ-  
823 ing code, unless this is central to the contribution (e.g., for a new open-source benchmark).

824 • The instructions should contain the exact command and environment needed to run to  
825 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
826 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.

827 • The authors should provide instructions on data access and preparation, including how to  
828 access the raw data, preprocessed data, intermediate data, and generated data, etc.

829 • The authors should provide scripts to reproduce all experimental results for the new pro-  
830 posed method and baselines. If only a subset of experiments are reproducible, they should  
831 state which ones are omitted from the script and why.

832 • At submission time, to preserve anonymity, the authors should release anonymized versions  
833 (if applicable).

834 • Providing as much information as possible in supplemental material (appended to the paper)  
835 is recommended, but including URLs to data and code is permitted.

## 836 6. Experimental Setting/Details

837 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpara-  
838 meters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

839 Answer: [YES]

840 Justification: Our study evaluates existing models in a zero-shot setting, so no training was  
841 performed. We provide all relevant testing details in the appendix, which includes a table of  
842 the models used and the specific inference hyperparameters (e.g., temperature, max tokens) for  
843 each one.

844 Guidelines:

- 845 • The answer NA means that the paper does not include experiments.
- 846 • The experimental setting should be presented in the core of the paper to a level of detail  
847 that is necessary to appreciate the results and make sense of them.
- 848 • The full details can be provided either with the code, in appendix, or as supplemental  
849 material.

## 850 7. Experiment Statistical Significance

851 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
852 information about the statistical significance of the experiments?

853 Answer: [NO]

854 Justification: We report the exact accuracy scores for each model on our benchmark based on  
855 a single run. We do not report error bars or conduct statistical significance tests, as our current  
856 evaluation is based on a single response generated for each prompt. This is noted as a limitation  
857 in our paper.

858 Guidelines:

- 859 • The answer NA means that the paper does not include experiments.
- 860 • The authors should answer “Yes” if the results are accompanied by error bars, confidence  
861 intervals, or statistical significance tests, at least for the experiments that support the main  
862 claims of the paper.
- 863 • The factors of variability that the error bars are capturing should be clearly stated (for  
864 example, train/test split, initialization, random drawing of some parameter, or overall run  
865 with given experimental conditions).
- 866 • The method for calculating the error bars should be explained (closed form formula, call to  
867 a library function, bootstrap, etc.)
- 868 • The assumptions made should be given (e.g., Normally distributed errors).
- 869 • It should be clear whether the error bar is the standard deviation or the standard error of  
870 the mean.
- 871 • It is OK to report 1-sigma error bars, but one should state it. The authors should preferably  
872 report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality  
873 of errors is not verified.
- 874 • For asymmetric distributions, the authors should be careful not to show in tables or figures  
875 symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- 876 • If error bars are reported in tables or plots, The authors should explain in the text how they  
877 were calculated and reference the corresponding figures or tables in the text.

## 878 8. Experiments Compute Resources

879 Question: For each experiment, does the paper provide sufficient information on the computer  
880 resources (type of compute workers, memory, time of execution) needed to reproduce the  
881 experiments?

882 Answer: [NO]

883 Justification: We do not provide a detailed breakdown of the compute resources used to run our  
884 evaluation framework. The majority of the models evaluated were accessed via third-party APIs,  
885 meaning the computational load for model inference was handled externally. For the self-hosted  
886 open-source models, the required resources are specified in their respective documentation.

887 Guidelines:

- 888 • The answer NA means that the paper does not include experiments.
- 889 • The paper should indicate the type of compute workers CPU or GPU, internal cluster, or  
890 cloud provider, including relevant memory and storage.
- 891 • The paper should provide the amount of compute required for each of the individual  
892 experimental runs as well as estimate the total compute.
- 893 • The paper should disclose whether the full research project required more compute than the  
894 experiments reported in the paper (e.g., preliminary or failed experiments that didn't make  
895 it into the paper).

## 896 9. Code Of Ethics

897 Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS  
898 Code of Ethics <https://neurips.cc/public/EthicsGuidelines>

899 Answer: [YES]

900 Justification: We have reviewed the NeurIPS Code of Ethics and confirm our research is in full  
901 compliance. The benchmark dataset was authored by the research team and does not contain any  
902 real-world, sensitive, or personally identifiable information. Our work does not involve human  
903 subjects.

904 Guidelines:

- 905 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 906 • If the authors answer No, they should explain the special circumstances that require a  
907 deviation from the Code of Ethics.
- 908 • The authors should make sure to preserve anonymity (e.g., if there is a special consideration  
909 due to laws or regulations in their jurisdiction).

## 910 10. Broader Impacts

911 Question: Does the paper discuss both potential positive societal impacts and negative societal  
912 impacts of the work performed?

913 Answer: [NO]

914 Justification: Our paper focuses on the positive impact of improving model reliability in finance.  
915 While we acknowledge the risks of instruction-following failures, we do not include a detailed  
916 discussion of potential negative societal impacts or malicious uses of more capable models. Our  
917 work is foundational, intended to improve evaluation rather than deploy a new application.

918 Guidelines:

- 919 • The answer NA means that there is no societal impact of the work performed.
- 920 • If the authors answer NA or No, they should explain why their work has no societal impact  
921 or why the paper does not address societal impact.
- 922 • Examples of negative societal impacts include potential malicious or unintended uses  
923 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g.,  
924 deployment of technologies that could make decisions that unfairly impact specific groups),  
925 privacy considerations, and security considerations.

- 926 • The conference expects that many papers will be foundational research and not tied to  
927 particular applications, let alone deployments. However, if there is a direct path to any  
928 negative applications, the authors should point it out. For example, it is legitimate to point  
929 out that an improvement in the quality of generative models could be used to generate  
930 deepfakes for disinformation. On the other hand, it is not needed to point out that a generic  
931 algorithm for optimizing neural networks could enable people to train models that generate  
932 Deepfakes faster.
- 933 • The authors should consider possible harms that could arise when the technology is being  
934 used as intended and functioning correctly, harms that could arise when the technology is  
935 being used as intended but gives incorrect results, and harms following from (intentional  
936 or unintentional) misuse of the technology.
- 937 • If there are negative societal impacts, the authors could also discuss possible mitigation  
938 strategies (e.g., gated release of models, providing defenses in addition to attacks, mecha-  
939 nisms for monitoring misuse, mechanisms to monitor how a system learns from feedback  
940 over time, improving the efficiency and accessibility of ML).

## 941 11. Safeguards

942 Question: Does the paper describe safeguards that have been put in place for responsible release  
943 of data or models that have a high risk for misuse (e.g., pretrained language models, image  
944 generators, or scraped datasets)?

945 Answer: [NA]

946 Justification: This question is not applicable. Our paper introduces a benchmark dataset of  
947 synthetic prompts and evaluation code, not a new generative model or a scraped dataset. The  
948 prompts are authored by the research team and do not contain sensitive or private information,  
949 posing a low risk of misuse.

950 Guidelines:

- 951 • The answer NA means that the paper poses no such risks.
- 952 • Released models that have a high risk for misuse or dual-use should be released with  
953 necessary safeguards to allow for controlled use of the model, for example by requiring  
954 that users adhere to usage guidelines or restrictions to access the model or implementing  
955 safety filters.
- 956 • Datasets that have been scraped from the Internet could pose safety risks. The authors  
957 should describe how they avoided releasing unsafe images.
- 958 • We recognize that providing effective safeguards is challenging, and many papers do not  
959 require this, but we encourage authors to take this into account and make a best faith effort.

## 960 12. Licenses for existing assets

961 Question: Are the creators or original owners of assets (e.g., code, data, models), used in the  
962 paper, properly credited and are the license and terms of use explicitly mentioned and properly  
963 respected?

964 Answer: [NO]

965 Justification: We credit the creators of the models we evaluate by citing their original research  
966 papers. However, we do not explicitly state the license or terms of use for each model in the  
967 paper itself. The models used are all well-known and their licensing information is publicly  
968 available.

969 Guidelines:

- 970 • The answer NA means that the paper does not use existing assets.
- 971 • The authors should cite the original paper that produced the code package or dataset.
- 972 • The authors should state which version of the asset is used and, if possible, include a URL.

- 973 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 974 • For scraped data from a particular source (e.g., website), the copyright and terms of service  
975 of that source should be provided.
- 976 • If assets are released, the license, copyright information, and terms of use in the package  
977 should be provided. For popular datasets, <https://paperswithcode.com/datasets> has  
978 curated licenses for some datasets. Their licensing guide can help determine the license of  
979 a dataset.
- 980 • For existing datasets that are re-packaged, both the original license and the license of the  
981 derived asset (if it has changed) should be provided.
- 982 • If this information is not available online, the authors are encouraged to reach out to the  
983 asset’s creators.

### 984 13. **New Assets**

985 Question: Are new assets introduced in the paper well documented and is the documentation  
986 provided alongside the assets?

987 Answer: [YES]

988 Justification: Yes, the new asset (IFF benchmark) is documented throughout the paper. We detail  
989 its design principles and content in the Methodology section. The appendix provides further  
990 details on the framework architecture, and we include a README file with usage instructions  
991 in the public code repository.

992 Guidelines:

- 993 • The answer NA means that the paper does not release new assets.
- 994 • Researchers should communicate the details of the dataset/code/model as part of their  
995 submissions via structured templates. This includes details about training, license, limita-  
996 tions, etc.
- 997 • The paper should discuss whether and how consent was obtained from people whose asset  
998 is used.
- 999 • At submission time, remember to anonymize your assets (if applicable). You can either  
1000 create an anonymized URL or include an anonymized zip file.

### 1001 14. **Crowdsourcing and Research with Human Subjects**

1002 Question: For crowdsourcing experiments and research with human subjects, does the paper  
1003 include the full text of instructions given to participants and screenshots, if applicable, as well  
1004 as details about compensation (if any)?

1005 Answer: [NA]

1006 Justification: Not applicable. Our paper does not involve crowdsourcing or research with human  
1007 subjects. The benchmark prompts were created by the authors of the paper.

1008 Guidelines:

- 1009 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1010 human subjects.
- 1011 • Including this information in the supplemental material is fine, but if the main contribution  
1012 of the paper involves human subjects, then as much detail as possible should be included  
1013 in the main paper.
- 1014 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or  
1015 other labor should be paid at least the minimum wage in the country of the data collector.

### 1016 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 1017 Subjects**

1018 Question: Does the paper describe potential risks incurred by study participants, whether such  
1019 risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals  
1020 (or an equivalent approval/review based on the requirements of your country or institution) were  
1021 obtained?

1022 Answer: [NA]

1023 Justification: Not applicable. As our research did not involve human subjects, Institutional  
1024 Review Board (IRB) approval was not required.

1025 Guidelines:

- 1026 • The answer NA means that the paper does not involve crowdsourcing nor research with  
1027 human subjects.
- 1028 • Depending on the country in which research is conducted, IRB approval (or equivalent)  
1029 may be required for any human subjects research. If you obtained IRB approval, you should  
1030 clearly state this in the paper.
- 1031 • We recognize that the procedures for this may vary significantly between institutions and  
1032 locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines  
1033 for their institution.
- 1034 • For initial submissions, do not include any information that would break anonymity (if  
1035 applicable), such as the institution conducting the review.