

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14

**TITLE**

Fine-tuned large language models for answering questions about full-text biomedical research studies

**SHORT TITLE**

Fine-tuned LLM for answering questions about biomedical research studies

**AUTHORS**

Kaiming Tao<sup>1</sup>, Jinru Zhou<sup>1</sup>, Zachary A. Osman<sup>1</sup>, Vineet Ahluwalia<sup>1</sup>, Chiara Sabatti<sup>2</sup>, Robert W. Shafer<sup>1\*</sup>

**AFFILIATIONS**

<sup>1</sup>Division of Infectious Diseases, Dept. of Medicine, Stanford University, Stanford, CA, USA 94305; <sup>2</sup>Dept. of Biomedical Data Sciences, Stanford University, Stanford, CA, USA 94305.

\*Corresponding author (rshafer@stanford.edu).

15 **ABSTRACT**

16 **Background:** Few studies have explored the degree to which fine-tuning a large-language model (LLM)  
17 can improve its ability to answer a specific set of questions about a research study. **Methods:** We created  
18 an instruction set comprising 250 marked-down studies of HIV drug resistance, 16 questions per study,  
19 answers to each question, and explanations for each answer. The questions were broadly relevant to  
20 studies of pathogenic human viruses including whether a study reported viral genetic sequences and the  
21 demographics and antiviral treatments of the persons from whom sequences were obtained. We fine-  
22 tuned GPT-4o-mini (GPT-4o), Llama3.1-8B-Instruct (Llama3.1-8B), and Llama3.1-70B-Instruct (Llama3.1-  
23 70B) using a quantized low rank adapter (QLoRA). We assessed the accuracy, precision, and recall of each  
24 base and fine-tuned model in answering the same questions on a test set comprising 120 different  
25 studies. Paired t-tests and Wilcoxon signed-rank tests were used to compare base models to one  
26 another, fine-tuned models to their respective base model, and the fine-tuned models to one another.  
27 **Results:** Prior to fine-tuning, GPT-4o displayed significantly greater performance than both Llama3.1-70B  
28 and Llama3.1-8B due to its greater precision compared with Llama3.1-70B and greater precision and  
29 recall compared with Llama3.1-8B; there was no difference in performance between Llama3.1-70B and  
30 Llama3.1-8B. After fine-tuning, both GPT-4o and Llama3.1-70B, but not Llama3.1-8B, displayed  
31 significantly improved performance compared with their base models. The improved performance of  
32 GPT-4o resulted from a mean 6% increased precision and 9% increased recall; the improved  
33 performance of Llama3.1-70B resulted from a 15% increased precision. After fine-tuning, Llama3.1-70B  
34 significantly outperformed Llama3.1-8B but did not perform as well as the fine-tuned GPT-4o model  
35 which displayed superior recall. **Conclusion:** Fine-tuning GPT-4o and Llama3.1-70B, but not the smaller  
36 Llama3.1-8B, led to marked improvement in answering specific questions about research papers. The  
37 process we describe will be useful to researchers studying other medical domains.

38

39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54

### **AUTHOR SUMMARY**

Addressing key biomedical questions often requires systematically reviewing data from numerous studies—a process that demands time and expertise. Large language models (LLMs) have shown potential in screening papers and summarizing their content. However, few research groups have fine-tuned these models to enhance their performance in specialized biomedical domains. In this study, we fine-tuned three LLMs to answer questions about studies on the subject of HIV drug resistance including one proprietary LLM (GPT-4o-mini) and two open-source LLMs (Llama3.1-Instruct-70B and Llama 3.1-Instruct-8B). To fine-tune the models, we used an instruction set comprising 250 studies of HIV drug resistance and selected 16 questions covering whether studies included viral genetic sequences, patient demographics, and antiviral treatments. We then tested the models on 120 independent research studies. Our results showed that fine-tuning GPT-4o-mini and Llama3.1-Instruct-70B significantly improved their ability to answer domain-specific questions, while the smaller Llama3.1-Instruct-8B model was not improved. The process we described offers a roadmap for researchers in other fields and represents a step in our attempt towards developing an LLM capable of answering questions about research studies across a range of pathogenic human viruses.

55

56

## **INTRODUCTION**

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

## **MATERIALS AND METHODS**

74

### **Fine-tuning**

75

76

77

78

The systematic review of data from multiple research studies is often required to answer many biomedical questions. The use of automated software tools to assist in reviewing research papers is a topic of increasing interest (1–7). Several research studies have described the use of LLMs, primarily the GPT-4.0 API or ChatGPT, to screen papers for specific criteria and for summarizing their content (8–14).

We previously assessed the use of GPT-4 to answer questions about studies on HIV drug resistance (HIVDR) (15). In that study, we found that GPT-4 reproducibly answered a set of 60 questions with a precision of 87% and a recall of 73% without human feedback. However, its performance was not improved with a 2000-word instruction sheet. The lack of improvement with this form of prompt engineering, led us to assess the degree to which fine-tuning could improve the performance of an LLM to specific answer questions about published HIVDR research studies.

We selected questions designed to determine whether a study reported HIV sequences and whether the sequences and their associated data were made publicly available. A fine-tuned model capable of answering questions about viral sequences, their public availability, and the demographics and antiviral treatments of the persons from whom the sequenced viruses were obtained would be invaluable to virology researchers, journal editors, and funding agencies.

Figure 1 outlines the approach to fine-tuning, testing, and analysis used in this study. We worked with three LLMs: (1) GPT-4o mini-2024-07-18 (GPT-4o); (2) Meta-Llama 3.1-70B-Instruct (Llama3.1-70B); and (3) Meta-Llama 3.1-8B-Instruct (Llama3.1-8B). Llama3.1-70B and Llama3.1-8B have 70 billion and 8 billion parameters, respectively. The exact parameter count for GPT-4o has not been reported.

79           Research papers: We selected 250 curated research papers about HIV drug resistance from the  
80 Stanford HIV Drug Resistance Database (HIVDB) encompassing studies of (1) HIV sequences from  
81 infected persons who were either antiretroviral treatment (ART)-experienced or ART-naïve; (2) HIV  
82 isolates with known mutations undergoing *in vitro* susceptibility testing; (3) wildtype HIV isolates  
83 cultured in the presence of increasing concentrations of an antiretroviral drug (i.e., *in vitro* selection  
84 experiments); and (4) different approaches to HIV sequencing and cloning. The complete list of papers  
85 are in Supplementary Table 1.

86           Research questions: We designed 16 questions addressing key aspects of HIVDR including (1)  
87 Whether sequencing was performed on HIV isolates obtained from patients and whether the sequences  
88 were made publicly available (5 questions); (2) The demographics of patients whose viruses were  
89 sequenced (2 questions); (3) The treatment characteristics of patients whose viruses were sequenced (5  
90 questions); and (4) The technical aspects of sequencing (4 questions). For eight questions, the answer  
91 was a list of items; for seven questions the answer was yes or no; and for one question, the answer was  
92 a number. For studies in which sequencing was not performed, the answers to patient demographics,  
93 treatments, and technical aspects of sequencing were considered to be “not reported”. Table 2 presents  
94 the complete list of questions along with their frequencies of being classified as true (for Boolean  
95 questions), non-empty (for list questions), or non-zero (for the single numeric question) in both the 250  
96 study instruction set and the 120 study test set.

97           Instruction set: The instruction set contained 250 training samples. Each sample contained (1) a  
98 markdown version of one of the 250 papers containing its abstract, methods, results, discussion, and  
99 data sharing statement; (2) the 16 research questions; (3) the answers to the research questions; and (4)  
100 the explanation for each of the answers, including the text relevant to each answer. For questions not  
101 addressed by a study, the explanation indicated that the study did not address the question. The  
102 complete training set is in Supplementary Table 2.

103            Training hyperparameters: We used Hugging Face's parameter efficient fine-tuning (16) using  
104 Quantized Low-Rank Adaptation (QLoRA) (17,18). Because our dataset was complex, we used a rank of  
105 25, which is in the upper range of the recommended values of 4 to 32. We set our batch size to one  
106 because our sample sizes were large with the median number of tokens per sample being 9343 (range:  
107 4261-22085).

108            Implementation: Table 1 shows the GPU, VRAM, and time requirements associated with fine-  
109 tuning and testing each model used in this study. For GPT-4o, the GPU and VRAM requirements were not  
110 known because fine-tuning and testing were done using the OpenAI API (19).

111

## 112 Testing and analysis

113            We created a test set comprising 120 journal articles. One hundred studies were identified by  
114 querying PubMed for journal articles about HIV drug resistance published in 2023. Twenty additional  
115 studies were selected from HIVDB because they reported data on uncommon topics that were unlikely  
116 to be reported in the first 100 studies. Like the training set, these papers included studies of viral  
117 sequences from HIV-infected persons, *in vitro* susceptibility testing, *in vitro* selection experiments, and  
118 technical aspects of HIV sequencing. Supplementary Table 3 lists the 120 papers used for testing.

119            For each question, we evaluated the answers of the original and fine-tuned models. Model-  
120 generated answers were compared to the human-curated answers, which were considered to be the  
121 ground truth. For the seven Boolean questions, we calculated the number of true positives, true  
122 negatives, false positives, and false negatives, as well as the model's precision, recall, accuracy, and F1-  
123 score. For the eight list-based questions, we defined true positive when the model outputted a non-  
124 empty list exactly matching the human answer; true negative when both the human answer and the  
125 model output were empty lists; false positive when the human answer was an empty list while the  
126 model outputted a non-empty list; and false negative when the human answer was a non-empty list, but

127 the model output was either an empty list or a list that did not match the human list. A similar approach  
128 was applied to the sole numeric question in that a result of 0 was considered analogous to an empty list.  
129 Supplementary Table 4 lists the correct answers and the answers for the three base and fine-tuned  
130 models for 1920 questions (120 papers x 16 questions).

131 We used Fisher Exact Tests to compare the accuracy, recall, and precision of the base model and  
132 fine-tuned model on the individual questions. For these tests, a Benjamini-Hochberg adjustment was  
133 calculated for the 16 questions evaluated. We used parametric (paired T-tests) and nonparametric  
134 (Wilcoxon-Rank Sign Test) tests to compare the average accuracy, recall, precision, and F1-score across  
135 questions of (1) the base models to one another; (2) the fine-tuned models to one another; and (3) each  
136 fine-tuned model to its base model. For these nine comparison, a Benjamini-Hochberg adjustment was  
137 calculated. A summary of the statistical analysis is in Supplementary File 5.

138

139

## **RESULTS**

140 Figure 2 displays the accuracy, precision, recall, and F1-score of the base and fine-tuned models  
141 for each of the 16 questions, individually. Points to the upper left of the diagonal line indicate questions  
142 for which there was any improvement for the fine-tuned model compared with the base model. Table 3  
143 displays the accuracy, precision, recall, and F1-scores for the base and fine-tuned GPT-4o and Llama3.1-  
144 70B models for those questions for which there was a significant increase in either precision or recall for  
145 the fine-tuned model by Fisher Exact testing. For GPT-4o, the questions with improvements were Q2,  
146 Q6, Q9, Q11, Q14, Q15, and Q16. For Llama-70B, the questions with improvements were Q14, Q15, and  
147 Q16. The fine-tuning of Llama3.1-8B did not result in a significant increase in precision or recall for any  
148 question.

149 Figure 3 compares the overall mean accuracy, precision, recall, and F1-score for the 16 questions  
150 pooled over the 120 test set studies. Prior to fine-tuning, GPT-4o displayed significantly greater precision

151 and recall compared with Llama3.1-8B and significantly greater precision compared with Llama3.1-70B  
152 using both paired t-tests and Wilcoxon-ranked sign tests (Figure 2A). There were no statistically  
153 significant differences between Llama3.1-70B and Llama3.1-8B base models.

154 After fine-tuning, GPT-4o displayed 6% increased accuracy, 6% increased precision, 9% increased  
155 recall, and 8% increased F1-score (Figure 2B). Llama3.1-70B displayed 8% increased accuracy, 15%  
156 increased precision, 1% increased recall, and 8% increased F1-score. Llama3.1-8B did not display  
157 significantly improved performance after fine-tuning. The increased recall, accuracy, and F1- score for  
158 the fine-tuned GPT-4o model and the increased precision, accuracy, and F1-score for the Llama3.1-70B  
159 model were statistically significant using both paired t-tests and Wilcoxon-ranked sign tests.

160 Figure 2C shows that the fine-tuned GPT-4o model displayed significantly greater accuracy,  
161 recall, and F1-score compared with both the fine-tuned Llama3.1-8B and Llama3.1-70B models using  
162 both paired t-tests and Wilcoxon-ranked sign tests. Llama3.1-70B displayed significantly greater accuracy,  
163 recall, and F1-score compared with Llama3.1-8B using both paired t-tests and Wilcoxon-ranked sign  
164 tests.

165

166

## **DISCUSSION**

167 Fine-tuning a foundation model provides significant advantages for handling domain-specific  
168 tasks. By training a model on targeted data, it becomes more reliable and effective in delivering accurate  
169 results without requiring complex prompts. Using a pre-trained model for fine-tuning lowers  
170 computational costs making it a highly efficient approach for specialized use cases. This study  
171 demonstrates that fine-tuning GPT-4o and Llama3.1-70B significantly improved their ability to answer  
172 questions about research studies in a specialized medical field, specifically those questions included in  
173 their training.



174 Our findings can be distilled into four main observations: (1) Prior to fine-tuning, GPT-4o  
175 displayed significantly greater performance than both Llama3.1-70B and Llama3.1-8B as a result of  
176 increased precision compared with Llama3.1-70B and increased precision and recall compared with  
177 Llama3.1-8B, while no difference in performance between Llama3.1-70B and Llama3.1-8B was observed;  
178 (2) After fine-tuning both GPT-4o and Llama3.1-70B, but not Llama3.1-8B, displayed significantly  
179 improved performance compared with their base models; (3) The improved performance of GPT-4o was  
180 a result of its 6% increased precision and 9% increased recall while the improved performance of  
181 Llama3.1-70B resulted from its 15% increased precision; (4) After fine-tuning, Llama3.1-70B  
182 outperformed Llama3.1-8B primarily as a result of its improved precision, but still did not perform as  
183 well as the fine-tuned GPT-4o model which displayed superior recall.

184 Most studies evaluating the potential use of LLMs for answering questions about research  
185 studies have evaluated the use of foundational models to determine whether the titles and abstracts of  
186 a study were likely to meet the inclusion criteria for a systematic review (8–14,20). Few have evaluated  
187 fine-tuned foundational models for their ability to answer questions about full-text research papers  
188 (21,22).

189 We examined the effects of fine-tuning on three models: GPT-4o, selected for its top-tier  
190 performance and ease of fine-tuning with just an API and Python script (23); Llama3.1-70B, chosen for  
191 its long context length and high ranking among open-source models (24); and Llama3.1-8B, to assess  
192 fine-tuning's impact on a smaller model. We intended to fine-tune the even larger Llama3.1-405B model,  
193 but the most widely available GPUs lacked the memory capacity to train this model, despite several  
194 optimization attempts. Testing this larger model would have required even more memory, and renting  
195 the necessary GPUs was cost-prohibitive at significantly more than \$10,000 for both fine-tuning and  
196 testing.

197           LoRA and QLoRA are widely used approaches for fine-tuning foundational models, as they adjust  
198 a small subset of parameters, reducing both memory usage and computational costs compared to full  
199 fine-tuning (17). LoRA adapters can be reused to fine-tune multiple foundational models and enable low-  
200 cost re-tuning when these models are updated. Moreover, LoRA adapters can be easily swapped or  
201 combined, facilitating modular specialization (17). QLoRA further introduces innovations that optimize  
202 memory usage while maintaining performance (18).

203           We selected a narrow topic to determine whether a foundation model could be fine-tuned to  
204 answer questions about full-text research studies. Given the specificity of our topic, we chose not to  
205 expand our training set or further optimize our model. However, the questions able to be answered by  
206 the fine-tuned models target key data types broadly relevant to studies of pathogenic human viruses,  
207 including those with available antiviral treatments, as well as those with pandemic potential. Therefore,  
208 the success that we have described is a step towards accomplishing the more ambitious goal of  
209 developing a fine-tuned model capable of answering questions broadly applicable to all pathogenic  
210 human viruses.

211

212 Financial disclosure statement

213 This work was funded by a grant from the National Institutes of Health: 2R24AI13661806. The  
214 funder played no role in this review.

215

216 Competing interests

217 RWS has received honoraria for participation in advisory boards from Gilead Sciences and  
218 GlaxoSmithKline and speaking honoraria from Gilead Sciences and ViiV Healthcare.

219

220 Data availability statement

221 All data generated or analyzed during this study are included in this published article and its  
222 supplementary information files. The Llama3.1-70B and Llama3.1-8B LoRA adapters developed  
223 for this study was shared on the Hugging Face platform (<https://huggingface.co/kmtao/llama3.1-8B-HIVDB-adapter>, <https://huggingface.co/kmtao/llama3.1-70B-HIVDB-adapter>).

225

226 Author contributions

227 Conceptualization: K.T., V.A., and R.W.S; Data Curation: K.T, J.Z and Z.A.O; Formal analysis:  
228 K.T, J.Z, C.S and R.W.S; Methodology: K.T, J.Z, J.A and R.W.S; Writing – Original Draft  
229 Preparation: K.T, J.Z and R.W.S; Writing – Review & Editing: K.T, J.Z, V.A, C.S., and R.W.S

230

231

## REFERENCES

- 232 1. van de Schoot R, de Bruin J, Schram R, Zahedi P, de Boer J, Weijdema F, et al. An open source  
233 machine learning framework for efficient and transparent systematic reviews. *Nat Mach Intell*. 2021  
234 Feb;3(2):125–33.
- 235 2. Cierco Jimenez R, Lee T, Rosillo N, Cordova R, Cree IA, Gonzalez A, et al. Machine learning  
236 computational tools to assist the performance of systematic reviews: A mapping review. *BMC*  
237 *Medical Research Methodology*. 2022 Dec 16;22(1):322.
- 238 3. Blaizot A, Veettil SK, Saidoung P, Moreno-Garcia CF, Wiratunga N, Aceves-Martins M, et al. Using  
239 artificial intelligence methods for systematic review in health sciences: A systematic review.  
240 *Research Synthesis Methods*. 2022;13(3):353–62.
- 241 4. Dijk SHB van, Brusse-Keizer MGJ, Bucsán CC, Palen J van der, Doggen CJM, Lenferink A. Artificial  
242 intelligence in systematic reviews: promising when appropriately used. *BMJ Open*. 2023 Jul  
243 1;13(7):e072254.
- 244 5. Jin Q, Kim W, Chen Q, Comeau DC, Yeganova L, Wilbur WJ, et al. MedCPT: Contrastive Pre-trained  
245 Transformers with large-scale PubMed search logs for zero-shot biomedical information retrieval.  
246 *Bioinformatics*. 2023 Nov 1;39(11):btad651.
- 247 6. Santos ÁO dos, da Silva ES, Couto LM, Reis GVL, Belo VS. The use of artificial intelligence for  
248 automating or semi-automating biomedical literature analyses: A scoping review. *Journal of*  
249 *Biomedical Informatics*. 2023 Jun 1;142:104389.
- 250 7. Kebede MM, Le Cornet C, Fortner RT. In-depth evaluation of machine learning methods for semi-  
251 automating article screening in a systematic review of mechanistic literature. *Research Synthesis*  
252 *Methods*. 2023;14(2):156–72.
- 253 8. Alshami A, Elsayed M, Ali E, Eltoukhy AEE, Zayed T. Harnessing the Power of ChatGPT for Automating  
254 Systematic Review Process: Methodology, Case Study, Limitations, and Future Directions. *Systems*.  
255 2023 Jul;11(7):351.
- 256 9. Schopow N, Osterhoff G, Baur D. Applications of the Natural Language Processing Tool ChatGPT in  
257 Clinical Practice: Comparative Study and Augmented Systematic Review. *JMIR Medical Informatics*.  
258 2023 Nov 28;11(1):e48933.
- 259 10. Syriani E, David I, Kumar G. Assessing the Ability of ChatGPT to Screen Articles for Systematic  
260 Reviews [Internet]. *arXiv*; 2023 [cited 2023 Nov 14]. Available from:  
261 <https://arxiv.org/abs/2307.06464>
- 262 11. Khraisha Q, Put S, Kappenberg J, Warraitch A, Hadfield K. Can large language models replace humans  
263 in systematic reviews? Evaluating GPT-4's efficacy in screening and extracting data from peer-  
264 reviewed and grey literature in multiple languages. *Research Synthesis Methods* [Internet]. 2024  
265 [cited 2024 Mar 17];n/a(n/a). Available from:  
266 <https://onlinelibrary.wiley.com/doi/abs/10.1002/jrsm.1715>

- 267 12. Guo E, Gupta M, Deng J, Park YJ, Paget M, Naugler C. Automated Paper Screening for Clinical  
268 Reviews Using Large Language Models: Data Analysis Study. *Journal of Medical Internet Research*.  
269 2024 Jan 12;26(1):e48996.
- 270 13. Polak MP, Morgan D. Extracting accurate materials data from research papers with conversational  
271 language models and prompt engineering. *Nat Commun*. 2024 Feb 21;15(1):1569.
- 272 14. Dennstädt F, Zink J, Putora PM, Hastings J, Cihoric N. Title and abstract screening for literature  
273 reviews using large language models: an exploratory study in the biomedical domain. *Syst Rev*. 2024  
274 Jun 15;13(1):158.
- 275 15. Tao K, Osman ZA, Tzou PL, Rhee SY, Ahluwalia V, Shafer RW. GPT-4 performance on querying  
276 scientific publications: reproducibility, accuracy, and impact of an instruction sheet. *BMC Medical  
277 Research Methodology*. 2024 Jun 25;24(1):139.
- 278 16. Huggingface. PEFT (PEFT) [Internet]. 2023 [cited 2024 Oct 2]. Available from:  
279 <https://huggingface.co/PEFT>
- 280 17. Hu EJ, Shen Y, Wallis P, Allen-Zhu Z, Li Y, Wang S, et al. LoRA: Low-Rank Adaptation of Large Language  
281 Models [Internet]. arXiv; 2021 [cited 2024 Jul 5]. Available from: <http://arxiv.org/abs/2106.09685>
- 282 18. Dettmers T, Pagnoni A, Holtzman A, Zettlemoyer L. QLoRA: Efficient Finetuning of Quantized LLMs  
283 [Internet]. arXiv; 2023 [cited 2024 Aug 30]. Available from: <http://arxiv.org/abs/2305.14314>
- 284 19. OpenAI Platform [Internet]. [cited 2024 Aug 30]. Available from: <https://platform.openai.com>
- 285 20. Zhang G, Jin Q, Zhou Y, Wang S, Idnay BR, Luo Y, et al. Closing the gap between open-source and  
286 commercial large language models for medical evidence summarization [Internet]. arXiv; 2024 [cited  
287 2024 Sep 12]. Available from: <http://arxiv.org/abs/2408.00588>
- 288 21. Shah A, Mehendale S, Kanthi S. Efficacy of Large Language Models in Systematic Reviews [Internet].  
289 arXiv; 2024 [cited 2024 Sep 12]. Available from: <http://arxiv.org/abs/2408.04646>
- 290 22. Susnjak T, Hwang P, Reyes NH, Barczak ALC, McIntosh TR, Ranathunga S. Automating Research  
291 Synthesis with Domain-Specific Large Language Model Fine-Tuning [Internet]. arXiv; 2024 [cited  
292 2024 Oct 11]. Available from: <http://arxiv.org/abs/2404.08680>
- 293 23. API Reference - OpenAI API [Internet]. [cited 2024 Sep 13]. Available from:  
294 <https://platform.openai.com/docs/api-reference/fine-tuning>
- 295 24. LYSYS.org. Chatbot Arena Leaderboard - a Hugging Face Space by lmsys [Internet]. [cited 2024 Sep  
296 13]. Available from: <https://huggingface.co/spaces/lmsys/chatbot-arena-leaderboard>
- 297  
298

299 **SUPPORT INFORMATION CAPTIONS**

300 Supplementary Table 1 (S1Table.xlsx): The list of research studies used for the instruction set.

301 Supplementary Table 2 (S2Table.xlsx): The list of training examples included in the instruction set.

302 Supplementary Table 3 (S3Table.xlsx): The list of research studies used for testing.

303 Supplementary Table 4 (S4Table.xlsx): The correct answers and the answers of each of the six models

304 (base and fine-tuned for GPT-4o, Llama3.1-70B, and Llama3.1-8B) for 1920 questions (120 test studies x

305 16 questions).

306 Supplementary File 5 (S5Appendix.xlsx): Tab 1 contains the raw data and results of Fisher Exact Tests for

307 each of the 16 questions for each of the six models (base and fine-tuned for GPT-4o, Llama3.1-70B, and

308 Llama3.1-8B). Tab 2 contains the raw data and nine comparisons between the models. Specifically,

309 paired t-tests and Wilcoxon signed-rank tests were used to compare the base models to one another, the

310 fine-tuned models to their respective base model, and the fine-tuned models to one another. Tab 3

311 illustrates how the Benjamini-Hochberg adjustment for the nine model comparisons was performed.

312

313

## **FIGURE LEGENDS**

### **Figure 1**

315 Approach to fine-tuning (A), testing (B), and analyses (C) performed in this study. Fine-tuning was  
316 performed using an instruction set comprising 250 marked-down research studies, 16 questions about  
317 each study, answers to each question, and explanations for each answer. GPT-4o was fine-tuned using  
318 the OpenAI API; Llama3.1-70B and Llama3.1-8B were fine-tuned using QLoRA (A). The accuracy,  
319 precision, recall, and F1-score of each base and fine-tuned model was assessed using a test set  
320 comprising 16 questions applied to 120 different published research studies on HIV drug resistance (B).  
321 Parametric (paired t-tests) and nonparametric (Wilcoxon signed-rank tests) methods were used to  
322 compare the performance of the base models to one another, the fine-tuned models to their respective  
323 base model, and the fine-tuned models to one another (C).

324

### **Figure 2**

326 Comparison of base and fine-tuned models for each of the 16 questions applied to 120 published test  
327 studies. The accuracy (A), precision (B), recall (C), and F1-score (D) of the GPT4o, Llama3.1-70B, and  
328 Llama3.1-8B models are shown with the metrics for the base model indicated on the X-axis and for the  
329 fine-tuned model indicated on the Y-axis. Points to the left of the diagonal line indicate those questions  
330 for which there was an improvement for the fine-tuned model compared with the base model.

331

### **Figure 3**

333 Comparisons of the performance of the base models to one another (A), the fine-tuned models to their  
334 respective base model (B), and the fine-tuned models to one another (C). The histograms in figures 3A  
335 and 3C represent the performance of the base and fine-tuned models, respectively. The histograms in  
336 figure 3B represent the difference in performance between the fine-tuned and base model. The error

337 bars in figure 3B represent the standard error of the mean of the paired differences between the fine-  
338 tuned and base models. The mean differences in accuracy, precision, recall, and F1-score between  
339 models are indicated above the relevant histograms when statistically significant using both parametric  
340 (paired t-test) and nonparametric (Wilcoxon signed-rank test) methods. The p values shown are for the  
341 paired t-test performed on the aggregate data for each of the 16 questions. After adjustment for nine  
342 comparisons, the Benjamini-Hochberg false discovery rate was  $\leq 0.05$  for each of the p values shown.  
343  
344



**Table 1. GPU, VRAM, and Time Requirements Associated with Fine-Tuning and Testing**

<u>Model</u>	<u>Fine-Tuning</u>			<u>Testing the Base Model</u>			<u>Testing the Fine-Tuned Model</u>		
	<u>GPU</u>	<u>VRAM</u>	<u>Time</u>	<u>GPU</u>	<u>VRAM</u>	<u>Time</u>	<u>GPU</u>	<u>VRAM</u>	<u>Time</u>
GPT-4o-mini	NA	NA	2h	NA	NA	1h	NA	NA	1h
Llama3.1 8B	1 A100	80G	1h	1 A100	80G	2h	2 A100	160G	13h
Llama3.1 70B	3 A100	240G	7h	3 A100	240G	5h	4 A100	320G	21h

Footnote: GPU (graphical processing unit); VRAM (video random access memory); A100 (Nvidia A100 tensor core GPU); VRAM is indicated as gigabytes.

**Table 2. Complete List of Questions with their Frequencies of True, Non-Empty or Non-Zero in Both Instruction Set and Test Set**

	Question	Subject	Type	Instruction set (%)	Test set (%)
Q1	Does the paper report HIV sequences from patient samples?	Data availability	Boolean	85.6	66.7
Q2	Does the paper report in vitro drug susceptibility data?	Data availability	Boolean	20	20.8
Q3	Were sequences from the paper made publicly available?	Data availability	Boolean	56.4	17.5
Q4	What were the GenBank accession numbers for sequenced HIV isolates?	Data availability	List	54.4	12.5
Q5	How many individuals had samples obtained for HIV sequencing?	Data availability	Number	82.8	64.2
Q6	From which countries were the sequenced samples obtained?	Demographics	List	76.8	56.7
Q7	From what years were the sequenced samples obtained?	Demographics	List	64	51.7
Q8	Were samples cloned prior to sequencing?	Technical	Boolean	2.8	2.5
Q9	Which HIV genes were reported to have been sequenced?	Technical	List	91.2	75.0
Q10	What method was used for sequencing?	Technical	List	64.8	45.8
Q11	What type of samples were sequenced?	Technical	List	79.2	52.5
Q12	Were any sequences obtained from individuals with virological failure on a treatment regimen?	Treatment	Boolean	36.4	30.8
Q13	Were the patients in the study in a clinical trial?	Treatment	Boolean	14.4	15.8
Q14	Does the paper report HIV sequences from individuals who had previously received ARV drugs?	Treatment	Boolean	46.4	46.7
Q15	Which drug classes were received by individuals in the study before sample sequencing?	Treatment	List	36.8	38.3
Q16	Which drugs were received by individuals in the study before sample sequencing?	Treatment	List	34.4	32.5

**Table 3. Accuracy, Precision, Recall and F1 Score for the Research Questions for which an Improvement was Observed After Fine-Tuning (FT) of Either GPT-4o or Llama3.1-70B**

	<u>Accuracy</u>		<u>Precision</u>		<u>Recall</u>		<u>F1-Score</u>	
	B	FT	B	FT	B	FT	B	FT
<b><i>GPT-4o</i></b>								
Q2. Does the paper report in vitro drug susceptibility data?	85.0	95.8 <sup>**†</sup>	58.1	91.7 <sup>**†</sup>	100.0	88.0	73.5	89.8
Q6. From which countries were the sequenced samples obtained?	80.0	89.2	100.0	93.7	64.7	86.8 <sup>**</sup>	78.6	90.1
Q9. Which HIV genes were reported to have been sequenced?	61.7	73.3	97.8	91.4	50.0	71.1 <sup>**</sup>	66.2	80.0
Q11. What type of samples were sequenced?	80.0	86.7	95.4	88.5	65.1	85.7 <sup>*</sup>	77.4	87.1
Q14. Does the paper report HIV sequences from individuals who had previously received ARV drugs?	84.2	91.7	100.0	91.1	66.1	91.1 <sup>**</sup>	79.6	91.1
Q15. Which drug classes were received by individuals in the study before sample sequencing?	57.5	77.5 <sup>**</sup>	44.9	77.1 <sup>**†</sup>	47.8	58.7	46.3	66.7
Q16. Which drugs were received by individuals in the study before sample sequencing?	57.5	74.2 <sup>**†</sup>	33.3	66.7 <sup>*(§)</sup>	30.8	41.0	32.0	50.8
<b><i>Llama3.1-70B</i></b>								
Q14. Does the paper report HIV sequences from individuals who had previously received ARV drugs?	74.2	84.2	73.6	100.0 <sup>***</sup>	69.6	66.1	71.6	79.6
Q15. Which drug classes were received by individuals in the study before sample sequencing?	34.2	70.8 <sup>***</sup>	29.6	76.2 <sup>***</sup>	52.2	34.8	37.8	47.8
Q16. Which drugs were received by individuals in the study before sample sequencing?	27.5	71.7 <sup>***</sup>	17.6	69.2 <sup>***</sup>	33.3	23.1	23.0	34.6

---

Footnote: OR: Odds ratio of Fisher's Exact Test. \*\*\* : unadjusted  $p < 0.001$ ; \*\* : unadjusted  $p < 0.01$ ; \* : unadjusted  $p < 0.05$ . After adjustment for 16 comparisons, the Benjamini-Hochberg false discovery rate was  $\leq 0.05$  for each significant comparison except for those indicated by <sup>(†)</sup> for which it was 0.06 and <sup>(§)</sup> for which it was 0.09.

Figure 1

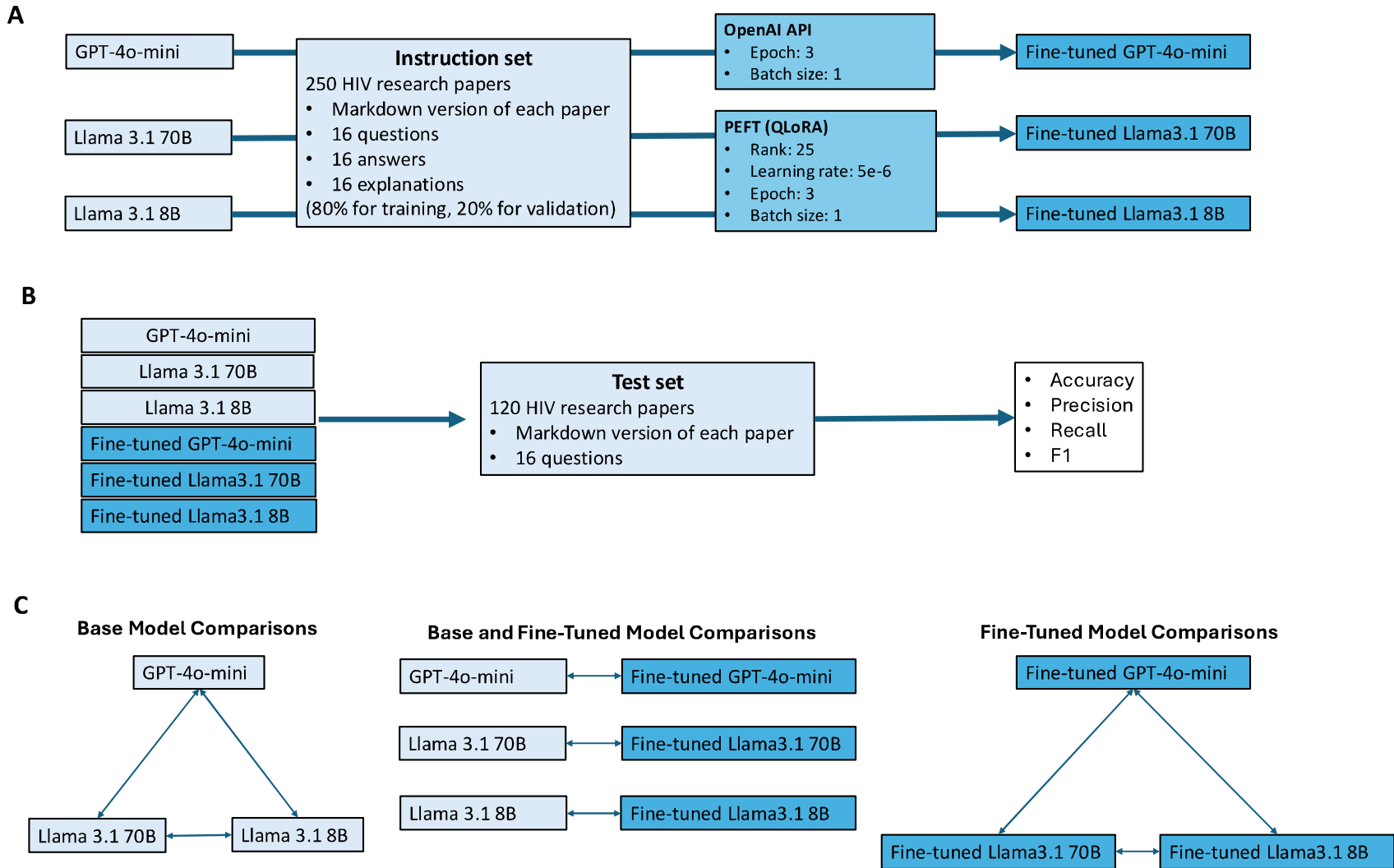


Figure 2

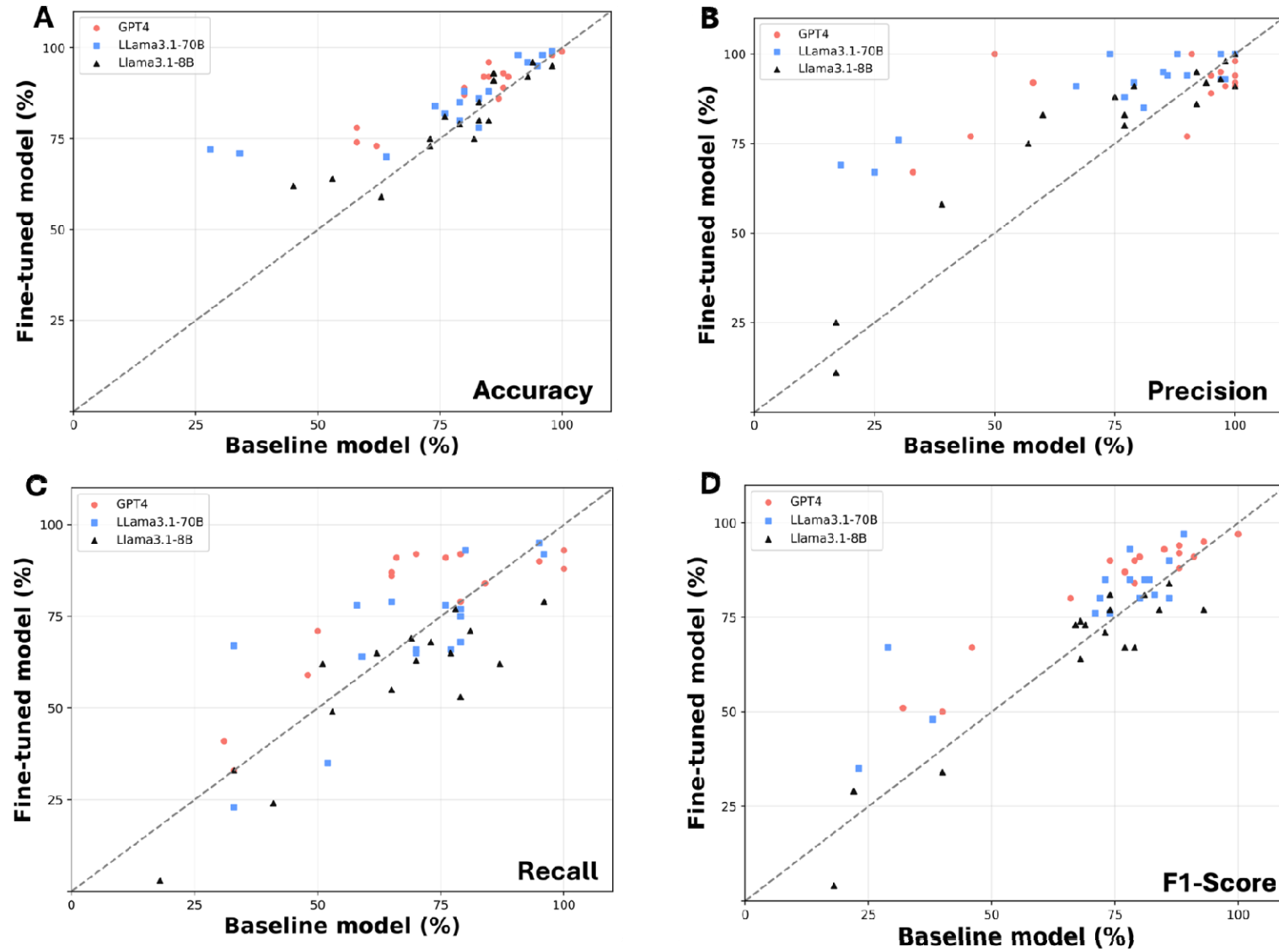
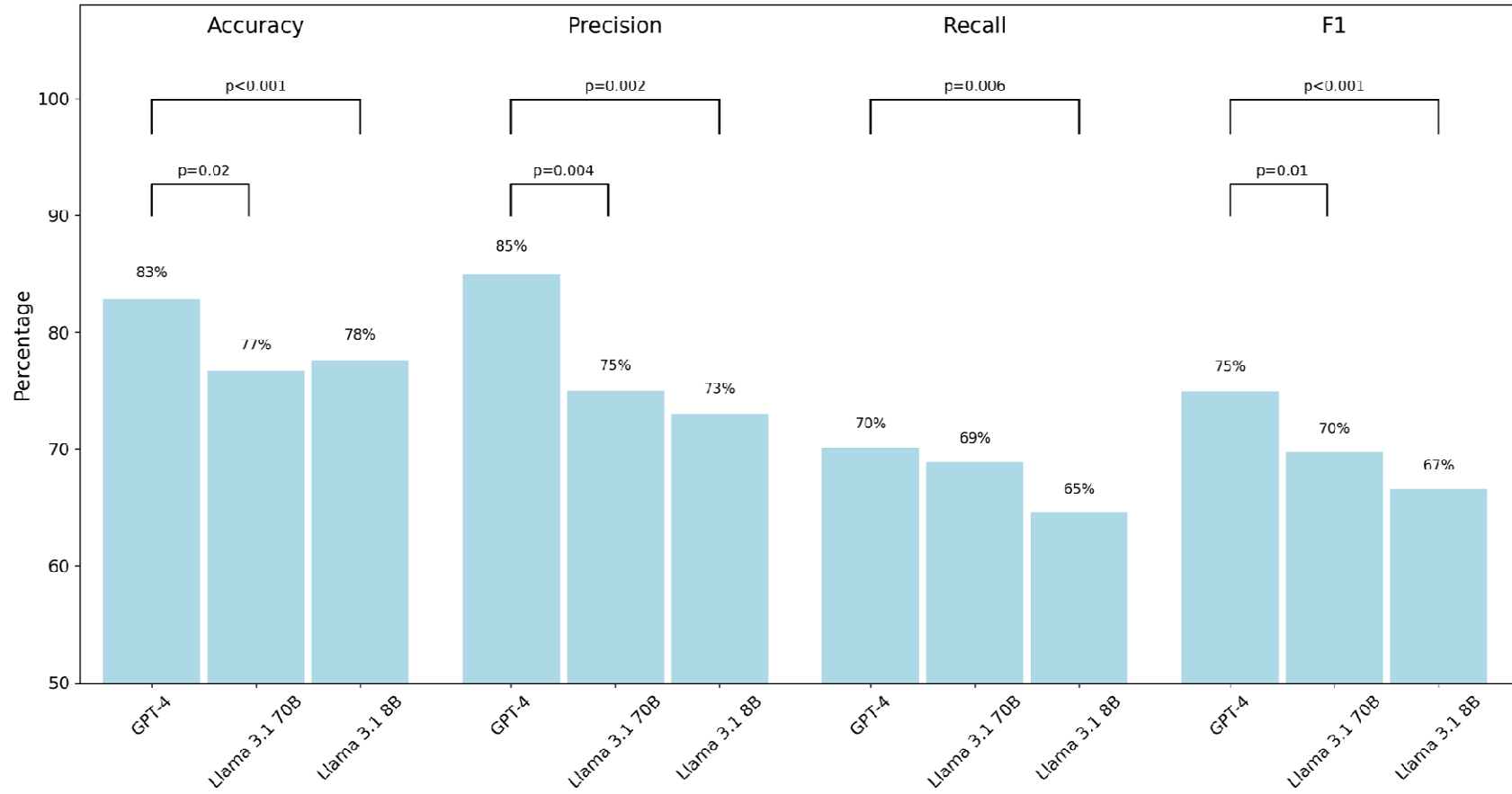
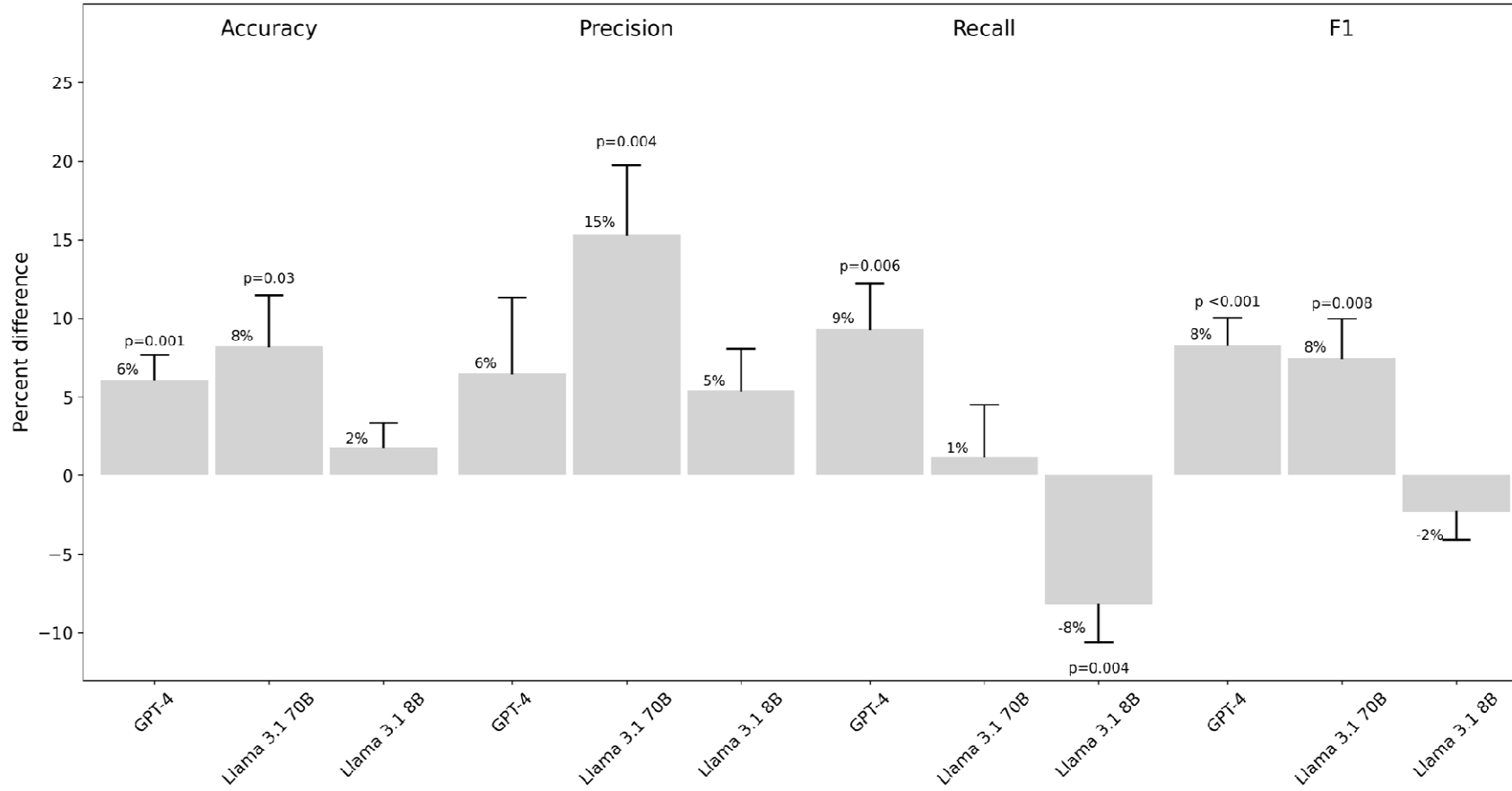


Figure 3

## A. Comparison of the Base Models



## B. Comparison of the Base and Fine-Tuned Models





## C. Comparison of the Fine-Tuned Models

