Can Language Models Follow Multiple Turns of Entangled Instructions?

Anonymous ACL submission

Abstract

Despite significant achievements in improving the instruction-following capabilities of large language models (LLMs), the ability to process multiple potentially entangled or conflicting instructions remains a considerable challenge. Real-world scenarios often require consistency across multiple instructions over time, such as secret privacy, personal preferences, and prioritization, which demand sophisticated abilities to integrate multiple turns and carefully bal-011 ance competing objectives when instructions intersect or conflict. This work presents a systematic investigation of LLMs' capabilities in 014 handling multiple turns of instructions, covering three levels of difficulty: (1) retrieving information from instructions, (2) tracking and reasoning across turns, and (3) resolving conflicts among instructions. We construct MUL-TITURNINSTRUCT with ~1.1K high-quality multi-turn conversations through the human-inthe-loop approach and result in nine capability categories, including statics and dynamics, reasoning, and multitasking. Our finding reveals an intriguing trade-off between different capabilities. While GPT models demonstrate superior memorization, they show reduced effectiveness in privacy-protection tasks requiring selective information withholding. Larger models exhibit stronger reasoning capabilities but still struggle with resolving conflicting instructions. Importantly, these performance gaps cannot be attributed solely to information loss, as models 034 demonstrate strong BLEU scores on memorization tasks. Still, their attention mechanisms fail to integrate multiple related instructions effectively. These findings highlight critical areas for improvement in complex real-world tasks involving multi-turn instructions.

1 Introduction

043

Large language models (LLMs) have made significant strides in following single, well-defined instructions (Brown et al., 2020; Inan et al., 2023),



Figure 1: A comparison between following each instruction individually and the scenario where the last instruction requires consideration of previous instructions. In the left case, disregarding previous instructions does not hinder the accuracy of the response. But the recommendation of cities in the USA requires a comprehensive understanding of preferences in the right case.

044

045

047

050

051

053

055

056

060

061

062

063

064

but how well can they follow multiple overlapping or even conflicting instructions? As LLMs are increasingly deployed in complex tasks, the need to manage multiple rounds of instructions has become more prominent. Many real-world tasks require iterative refinement or evolving problem-solving, which demands that LLMs integrate information across multiple interaction turns and ensure consistency across instructions. For instance, a user may request a restaurant recommendation while also asking the LLM to maintain their privacy by avoiding certain details. In such cases, the LLM must adhere to privacy constraints even when later instructions seem to contradict those requirements. Similarly, when providing a recommendation, the LLM needs to consider prior instructions, such as personal preferences mentioned earlier in the conversation. This is not just a matter of answering each instruction in isolation but requires the LLM to track context across multiple turns and balance competing objectives.

GPT-3.5-turbo	1st Round	2nd Round	Avg.
Seeing All	8.08	7.81	7.94
Current Only	8.08	7.8	7.94

Table 1: GPT-3.5-turbo behaves similarly on MT-Bench each round when seeing all instructions (1st row) or only the last instruction (2nd row).

104

106

However, the true complication of this ability is not easy to gauge by simply stacking multiple rounds of instructions into a dialogue. For example, in our evaluation of GPT-3.5-turbo on the MT-Bench dataset (Zheng et al., 2023), we observed that the model performs similarly whether it sees the full conversation history or only the most recent instruction, as shown in Table 1. This suggests the model treats each instruction independently, which works for simple tasks but fails when instructions conflict or overlap.

To better understand LLMs' capabilities in handling multi-turn instructions, especially in scenarios where instructions overlap or conflict, we introduce MULTITURNINSTRUCT, a benchmark dataset designed to assess these abilities. Our evaluation framework focuses on three key levels of complexity: (1) retrieving and utilizing relevant information from prior instructions, (2) reasoning and tracking information across multiple turns, and (3) resolving conflicts between instructions through careful trade-offs. Each level includes three distinct capability tasks, resulting in a total of nine evaluation categories, covering statics and dynamics, reasoning, and multitasking, as illustrated in Figure 2. Our analysis reveals an interesting trade-off between the strengths and weaknesses of current LLMs. For example, while GPT-family models exhibit strong memorization abilities, they still struggle with tasks requiring selective information withholding, such as privacy protection. Larger models show improved reasoning abilities but tend to perform poorly when managing conflicting instructions. These findings highlight a nuanced interplay among memorization, attention mechanisms, and multi-turn reasoning capabilities in modern LLMs, shedding light on the complexities of achieving reliable multi-turn context management.

2 Related Work

Instruction Following and Multi-Turn Interaction Pre-trained large language models have demonstrated impressive emergent ability to fol-



Figure 2: MULTITURNINSTRUCT consists of ~ 1.1 K spanning across three levels of difficulty and 9 capabilities, with balanced numbers of samples in each capability (numbers shown in the figure). Table 2 provides a more detailed list of task descriptions.

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

129

130

131

132

133

134

135

136

137

138

139

140

141

142

low instructions (Radford et al., 2023; Brown et al., 2020; Wei et al.). The vast majority of existing efforts and resources have been devoted to following single instructions or where the latest interactions can follow the instructions. For example, Multi-IF (He et al., 2024) studies the scenario where the user sequentially applies additional instructions to the last response. In a multi-round benchmark MT-Eval (Kwan et al., 2024), 3 out of 4 tasks are constructed in a way where the new instruction does not rely on or only follows up on the previous response. In Section 1, we show that the widely studied MT-Bench (Zheng et al., 2023) can be solved with the latest round of interactions. These can be regarded as knowledge conflicts (Xu et al., 2024). Similarly, in other multiturn interaction benchmarks, including Parrot (Sun et al., 2024), SIT (Hu et al., 2025), and MT-Bench 101 (Bai et al., 2024), little attention was explicitly paid to ensuring the inter-dependency of instructions. RefuteBench (Yan et al., 2024) provides a complementary perspective on LLMs' ability to handle refutation and user correction in multi-turn interactions. Besides, (Ferraz et al., 2024) uses real user-AI dialogues data to evaluate LMs As stated in a most recent survey (Zhang et al., 2025) "... no existing work has systematically analyzed ... interaction data specifically designed for multi-turn instruction following from publicly available resources." To our knowledge, our benchmark is the first one to explicitly investigate scenarios in which adherence to all rounds of entangled instructions is necessary.

Privacy Protection on LLMs The degree to which LLMs can comprehend and handle such information while complying with privacy regu-

Task	Requirement	Value	Scenarios	Metric
Memorization	Recalling all the instruction be- fore	Informativeness authenticity	, meetings, conversa- tions	BLEU score
Privacy Pro- tection	If requested, keep a secret in later dialogue turns	Privacy, Trust- worthiness	private assis- tant	Non-matching rate
Dynamic In- struction	As the user's constraints evolve and replace, always answer the selection result based on the up- to-date constraints	Adaptability	goods, num- bers, cities	Exact match rate
Dynamic En- vironment	As the item set updates, always answer the selection based on the up-to-date set	Adaptability	goods, num- bers, cities	Exact match rate
Personalization	Recommending items based on the user's personal profile	Personalization	diet, nation- ality	Exact match rate
Triggering	When a trigger is met in a con- ditional instruction, flag by re- sponding certain message	Safety, trust- worthiness	warning, re- minder	Exact match rate
Multitasking	Returning to a previous task when the current task is finished	Flexibility	QA, role- playing	Exact match rate
Recursive Reasoning	Carry out reasoning that depends on outputs several steps before	Accuracy	algorithm, math	Exact match rate
Prioritization	On a stream of potentially con- flicting commands, carry out each if and only if it does not con- flict with a higher-priority one	Safety	scheduling, permission manage- ment, control	Exact match rate

Table 2: A detailed description of the tasks involved in MULTITURNINSTRUCT dataset along with their associated values, grounded scenarios in real life, and evaluation metric.

lations has attracted significant attention from the 143 research community. Several studies have demon-144 strated that LLMs are vulnerable to leaking private 145 information (Staab et al., 2023; Huang et al., 2022a; 146 Kim et al., 2023a) and are susceptible to data ex-147 traction attacks (Wang et al., 2023; Li et al., 2023b). 148 To address these issues, some research efforts have 149 focused on developing Privacy-Preserving Large 150 Language Models (Behnia et al., 2022; Montagna 151 et al., 2023; Chen et al., 2023; Kim et al., 2023b; Utpala et al., 2023), employing techniques such as 153 differential privacy (Qu et al., 2021; Huang et al., 154 2022b; Igamberdiev and Habernal, 2023). There-155 fore, conducting a comprehensive benchmark that 156 evaluates these privacy-preserving methods in con-157 junction with various privacy attack techniques is 158 both essential and meaningful. Typically, bench-159 marking research (Zhang et al., 2024; Huang et al., 160

2024) categorizes privacy concerns into two main areas (Li et al., 2023a; Huang et al., 2022c): *Privacy Awareness* and *Privacy Leakage*, and employs Refusing to Answer and other utility metrics to measure the privacy understanding of LLMs. 161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

3 Constructing MULTITURNINSTRUCT: A Diverse Task Set

To thoroughly assess LLMs' ability to process and respond to multi-turn instructions, we introduce MULTITURNINSTRUCT, a dataset comprising approximately 1.1K multi-turn dialogues across a variety of real-world scenarios. Unlike single-turn evaluations, our benchmark challenges models to track, recall, and adhere to instructions as conversations evolve. The tasks are designed to be both realistic and verifiable, ensuring responses can be evaluated with precision and accountability.



Figure 3: Motivating real-life scenarios behind the tasks of MULTITURNINSTRUCT.



Figure 4: Distribution of conversation turn numbers across the dataset, illustrating the frequency of different turn counts.

Each task is categorized into one of three difficulty levels. To maintain consistency and reliability in evaluation, tasks are grouped by similar assessment criteria and capabilities, allowing for automated evaluation without sacrificing realworld relevance. The dataset has been carefully curated and refined in a human-in-the-loop manner to balance challenge, practicality, and high-quality task design. Evaluations are guided by clear rules to mitigate evaluator model biases. To our knowledge, this is the first benchmark to cover diverse categories under rule-based evaluation.

178

179

181

183

185

189

3.1 Curating Data in Each Task

During the collection of MULTITURNINSTRUCT, 191 we maintain a balance between challenge and real-192 ity: we aim to ensure that the data challenge LLMs 193 on the evaluated capabilities associated with the 194 tasks, and also ensure that data reflects the real 195 events in human life. To this end, we combine two 196 data construction approaches: existing data con-197 version and novel data curation. Some data come 198 from data converted from existing datasets, and oth-199 ers are curated with synthesis or a mixture of both. 200 All data points are manually checked and refined 201 to ensure quality. In the end, we collected 1.1K 202 multi-turn instruction data dialogues across nine 203 capability tasks, with more than 100 dialogues in 204 each task. To ensure the realism of the constructed 205 data, the dialogue includes rounds of instruction 206 that are realistic but not intended for evaluation 207 capabilities associated with the tasks. The models' 208 responses in these rounds are excluded from evalu-209 ation. All metrics have scores ranging within [0, 1], 210 as detailed in Table 2. The detailed data collections 211 for each task are listed as follows: 212



Figure 5: Score of mainstream LLMs on MULTITURNINSTRUCT. Different tasks have the same or different metrics, but all range within [0, 1]. Higher always means better performance.

1. Privacy protection: The task consists of two parts of data. The first part of the tasks is converted from the Enron Email dataset (Corp and Cohen, 2015) which contains private information such as credit card numbers, phone numbers, and email addresses. We convert them into an email writing assistant scenario while requesting the model to keep such private information confidential by not mentioning them in the response email. The second part of the task comes from prompting GPT-4 to curate a list of real-life scenarios where certain private information (health conditions, exam scores, family financial status) is requested not to be mentioned in the later conversation.

213

214

215

216

218

219

223

231

240

2. Dynamic instruction & Dynamic Environment: We convert the publicly available Amazon Product dataset (Hou et al., 2024) into a simulated scenario where the user questions the rating, rating number, or price of products in a synthetic marketplace. In the dynamic instruction task, a random list of 4 to 8 products from a certain category is presented in the first instruction as the context. In each round, the user questions a different question about them. The scenario in the dynamic environment dataset is similar. The question remains the same, but the products constantly update their prices, ratings, and rating numbers throughout the turns, identical to a reallife evolving market. 241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

257

258

259

260

261

262

263

264

265

266

- 3. **Personalization**: We convert the food.com recipe dataset (Li, 2019) into a multi-turn personalized recommendation dialogue. The user mentions their diet preferences (vegan, allergies, or dislikes to certain types of foods) in the first round and requests a personalized diet recommendation (e.g., the recipe with the lowest calories or highest fat) from a given recipe list in the end. The model is expected to avoid foods that meet the users' diet preferences.
- 4. **Triggering**: We prompt GPT-4 to create a list of real-life scenarios where the user instructs the model to remind them whenever a triggering condition is met in the subsequent dialogue. For instance, users may request a reminder for a to-do if a specific date or time condition is met, if they make a spelling error, or if certain entities are mentioned.
- 5. **Multitasking**: This task simulates the scenario where the user is involved in multiple tasks and switches between them. The first part of the dataset comes from converting the SQuAD dataset (Rajpurkar et al., 2016) into a



Figure 6: Heatmap of LLM performance and subtask correlations.

multi-document question-answering(QA) dialogue. Three documents are presented first, and the user switches between the documents to question about in each round. The second part of the dataset is converted from the Amazon Product dataset. Three categories of products are presented at first, and the user selects one category and questions the model about it.

270

271

276

277

278

279

281

283

287

296

297

- 6. Recursive Reasoning: The first part of the dataset consists of question-answering on recursive math functions, ranging in difficulty from the Fibonacci sequence $(F_n = F_{n-1} +$ F_{n-2}) to self-generative sequences¹. These functions are recursively defined over their previous values. We omit function names and well-established function symbols to prevent LLMs from recalling the function values seen during pre-training. Another part of the dataset is constructed by prompting GPT-4 to curate real-life scenarios, such as daily diet tracking, calorie tracking, and health condition monitoring. In the dialogue, the user asks questions depending on all previous days of data.
 - 7. **Prioritization**: This task requires the model to follow an accumulating number of conflicting instructions, each with a different importance level. The model is requested to follow the instruction, which can outrule previous

lower-priority instructions, while not violating higher-priority ones before. We implemented a simulator to heuristically curate a diverse set of dialogues. Scenarios include scheduling events on the calendar, room temperature setting, and light control.

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

327

328

8. **Memorization**: We convert a subset of data from the aforementioned other tasks by asking to repeat a specific (e.g., 3rd) instruction. This task is regarded as the simplest benchmarking subtask to test the LLMs' basic capabilities.

4 How Do LLMs Handle Interleaving Instructions

4.1 No LLM Is A Single Winner on MULTITURNINSTRUCT

We evaluate a diverse set of mainstream LLMs, from proprietary models (GPT (Achiam et al., 2023) and Claude (Anthropic, 2024)) to open source models (Mistral (Jiang et al., 2023) and Llama family (Dubey et al., 2024; Touvron et al., 2023a,b)) based on deterministic matching (i.e., BLEU and exact match). There is no single winner across all capabilities and even no family that consistently outperforms other families. GPT-40 performs the best among all models in 6 out of 9 tasks. Llama-3 performs the best among the opensource models in 7 out of 9 tasks. We find that the models performing well on basic tasks such as memorization generally perform well on many other tasks, including Dynamic Environment, Dynamic Instruction, Triggering, Multitasking, and

^le.g., https://en.wikipedia.org/wiki/Kolakoski_ sequence



(a) The performance decreases on GPT-40 on a selection of tasks as the preceding conversation contains more and more rounds. The trend line is fit with the best exponential function. (Note that blanks always mean non-existent scores due to a lack of data with a certain number of rounds in the datasets instead of a 0-score.)



(b) On some other tasks, especially those falling in the "context retrieving" category, there is less of a descending trend. Scores are on GPT-40.

Figure 7: Performance trends of GPT-40 across different tasks with increasing conversational rounds.

Recursive Reasoning. The rest of the tasks, namely Privacy Protection, Personalization, and Prioritization, which fall within the "contradiction resolution" category in Figure 2, seem to require different dimensions of ability, which we analyze in the following section.

332

333

334

335

4.2 Capabilities Conflict with Each Other

Despite the expectation that improved intelligence 336 will positively reflect in performance in most tasks, 337 Figure 6b shows how tasks positively and negatively correlate in their performance on LLMs. The 339 capabilities of Dynamic Environment, Dynamic Instruction, Multitasking, and Recursive Reason-341 ing do positively correlate with each other, prob-342 ably due to their similar nature in handling interdependency between rounds of instructions. However, tasks falling within the "contradiction resolution" category in Figure 2, namely Privacy Protection, Personalization, and Prioritization, are less correlated with the other tasks. Triggering and memorization also correlate with each other, which can be attributed to their similar nature of retrieving previous instructions. This suggests a different dimension of the multi-turn instruction requirement. 352



Figure 8: Histogram showing the statistics on turn numbers in the dataset. The x-axis represents the range of turn numbers, while the y-axis depicts the frequency of occurrences for each range.

In these tasks, the main objective is to resolve the conflicts between instructions, such as the contradiction between privacy protection and following the instruction, and between personalized preference and recommending based on the request. Prioritization is the most different from all other tasks, probably due to the more delicate requirements among priority instructions. 353

354

357

358

360

361

362

363

364

366

367

369

370

372

374

375

376

377

378

380

381

382

384

4.3 Models Correlate by Inheritance

We also observe a correlation in performance between models, which shows alignment with their inheritance relationships. As in Figure 6a, LLMs from each model family show more or less internal correlation with each other, especially in the GPT, Mistral, and Llama families. Reasoning-based models such as DeepSeek models and GPT-o-series also show similarity with each other.

4.4 The Scores Decrease as the Conversation Progresses

If our hypothesis holds that obedience to investigated instructions depends on previous ones, following later instructions will be harder because there will be more instructions involved. Figure 7 demonstrates a general performance decrease on GPT-40 on a selection of tasks as the preceding conversation contains more and more rounds. The trend line fits the best exponential function, where we skip non-existent scores due to a lack of data with a certain number of rounds in the datasets. In Figure 7a, five out of nine tasks show consistent decreasing trends of scores as the number of historical rounds increases. ² In particular, as shown in figure 7b, the "context retrieving" category is less

²The personalization category is omitted as it has a fixed number of rounds.



(a) In the "Privacy Protection" task, Llama-3.2-Instruct leaves little attention to the instruction to "keep the privacy information a secret".



(b) The "Dynamic Environment" subtask requires tracking the environment's changes across all turns of instructions, but Llama-3.2-Instruct focuses its attention primarily on the last turn of instruction

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

Figure 9: Attention heatmaps for Llama-3.2-Instruct failure cases, showing an insufficient focus on privacy instructions (left) and a dominant emphasis on the latest instruction in dynamic environments (right).

affected by the number of rounds. This is probably due to a balance between a longer conversation (negative factor) and more information to rely on in context (positive factor), canceling out their effects.

4.5 Do the Models Forget About the Instructions?

To refute the null hypothesis that the decrease in model performance comes from the inability to memorize the instructions, we plot the distribution of BLEU scores in the Memorization task in Figure 8. Note that the Memorization task has an average of 0.821 BLEU score for GPT-40, which is a perfect n-gram overlap between the system answer and the reference answers. We see that 61% percent of data has a 1.00 BLEU score, and most of the other scores are also biased towards the high end. Similar observations can be made on other models' high performance in the Memorization task in Figure 5. This verifies that the models can retrieve the instruction information with high accuracy, and the decrease in scores should be more attributed to the inability to keep track and follow them.

Analysis of Attention Patterns in 4.6 **Multi-turn Tasks**

To better understand the root causes of model failures, we use Figure 9 to illustrate attention heatmaps for two examples where Llama-3.2-412 Instruct fails. In the "Privacy Protection" task (Figure 9a), the model exhibits insufficient focus on the instruction to "keep the privacy information 415 a secret" but focuses mainly on the latest instruc-416 tion, which encourages the detailed response with

sufficient information exposed. This behavior suggests that the model may not sufficiently focus on restrictive instructions earlier, even though they have near-perfect recall of them as shown in Section 4.5. In the "Dynamic Environment" subtask (Figure 9b), the model is required to track changes across multiple instruction turns. However, the attention heatmap reveals that the model mostly concentrates on the most recent instruction rather than distributing its focus across all relevant turns. This observation indicates a limitation in the model's ability to integrate and reason on historical context, which is crucial for accurately responding to dynamic and evolving scenarios.

5 **Conclusions and Future Work**

In this work, we systematically evaluate the ability of large language models (LLMs) to process and respond to multi-turn instructions, particularly when those instructions overlap or conflict. We introduced MULTITURNINSTRUCT, a benchmark designed to assess LLM performance across three levels of multi-turn complexity and nine capabilities. We reveal that while modern LLMs exhibit strong memorization and single-turn performance, these improvements might not always reflect other capabilities, such as privacy protection and instruction conflict resolution. We also illustrate how the model failures are associated with their attention insufficiently applied to earlier involved instructions. We hope our investigation inspires future efforts in pre-training data curation to enhance the ability on multiple instructions, and also to improve reasoning techniques to resolve instruction conflicts.

411

413

414

451 Limitations

Dataset Scope and Coverage While MULTI-452 TURNINSTRUCT contains a diverse set of multi-453 454 turn dialogues, it may not capture the full range of real-world scenarios and edge cases that LLMs 455 might encounter. The dataset is structured and cu-456 457 rated, which could limit its ability to reflect more spontaneous or less predictable real-world conver-458 sations. 459

Task Complexity Although we designed tasks at 460 different difficulty levels, there may be more com-461 plex or nuanced forms of instruction entanglement 462 and conflict resolution that are not fully represented 463 in our evaluation framework. For example, tasks 464 that require deeper emotional or social context un-465 derstanding could further challenge current models, 466 but these are not explored in this work. 467

Evaluation Bias The benchmark is designed to
be objective: the evaluation process is influenced
by the design of the tasks, which could introduce
certain biases in assessing LLM performance. Furthermore, the human-in-the-loop approach used to
curate the dataset, which could potentially introduce subjectivity in task design.

References

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499 500

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Anthropic. 2024. The claude 3 model family: Opus, sonnet, haiku. In *Anthropic Research*. Anthropic.
- Ge Bai, Jie Liu, Xingyuan Bu, Yancheng He, Jiaheng Liu, Zhanhui Zhou, Zhuoran Lin, Wenbo Su, Tiezheng Ge, Bo Zheng, and Wanli Ouyang. 2024.
 MT-bench-101: A fine-grained benchmark for evaluating large language models in multi-turn dialogues. In Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 7421–7454, Bangkok, Thailand. Association for Computational Linguistics.
- Rouzbeh Behnia, Mohammadreza Reza Ebrahimi, Jason Pacheco, and Balaji Padmanabhan. 2022. Ew-tune: A framework for privately fine-tuning large language models with differential privacy. In 2022 IEEE International Conference on Data Mining Workshops (ICDMW), pages 560–566. IEEE.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda

Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901. 501

502

504

505

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

- Chaochao Chen, Xiaohua Feng, Jun Zhou, Jianwei Yin, and Xiaolin Zheng. 2023. Federated large language model: A position paper. *arXiv preprint arXiv:2307.08925*.
- Enron Corp and William W. Cohen. 2015. Enron email dataset. Software, E-Resource. Retrieved from the Library of Congress.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Thomas Palmeira Ferraz, Kartik Mehta, Yu-Hsiang Lin, Haw-Shiuan Chang, Shereen Oraby, Sijia Liu, Vivek Subramanian, Tagyoung Chung, Mohit Bansal, and Nanyun Peng. 2024. Llm self-correction with decrim: Decompose, critique, and refine for enhanced following of instructions with multiple constraints. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 7773–7812.
- Yun He, Di Jin, Chaoqi Wang, Chloe Bi, Karishma Mandyam, Hejia Zhang, Chen Zhu, Ning Li, Tengyu Xu, Hongjiang Lv, et al. 2024. Multi-if: Benchmarking llms on multi-turn and multilingual instructions following. *arXiv preprint arXiv:2410.15553*.
- Yupeng Hou, Jiacheng Li, Zhankui He, An Yan, Xiusi Chen, and Julian McAuley. 2024. Bridging language and items for retrieval and recommendation. *arXiv preprint arXiv:2403.03952*.
- Hanxu Hu, Simon Yu, Pinzhen Chen, and Edoardo M Ponti. 2025. Fine-tuning Large Language Models with Sequential Instructions. In *Proceedings of the* 2025 Conference of the North American Chapter of the Association for Computational Linguistics (NAACL), Albuquerque, New Mexico. Association for Computational Linguistics.
- Jie Huang, Hanyin Shao, and Kevin Chen-Chuan Chang. 2022a. Are large pre-trained language models leaking your personal information? *Preprint*, arXiv:2205.12628.
- Jie Huang, Hanyin Shao, and Kevin Chen-Chuan Chang. 2022b. Are large pre-trained language models leaking your personal information? In *Findings of the Association for Computational Linguistics: EMNLP* 2022, pages 2038–2047, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Kung-Hsiang Huang, ChengXiang Zhai, and Heng Ji. 2022c. Improving cross-lingual fact checking with cross-lingual retrieval. In *Proc. The 29th International Conference on Computational Linguistics* (COLING2022).

- 555 556
- 55 55
- 56
- 56
- 563
- 564 565 566 567
- 568 569
- 570 571
- 572 573
- 574 575
- 576 577
- 578 579 580
- 58
- 582 583 584

587 588

59 59

59 59

59

59

- 59
- 59

6

6

(

6

60 60

- Yue Huang, Lichao Sun, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, et al. 2024. Position: Trustllm: Trustworthiness in large language models. In *International Conference on Machine Learning*, pages 20166–20270. PMLR.
- Timour Igamberdiev and Ivan Habernal. 2023. Dp-bart for privatized text rewriting under local differential privacy. *ArXiv*, abs/2302.07636.
- Hakan Inan, Kartikeya Upasani, Jianfeng Chi, Rashi Rungta, Krithika Iyer, Yuning Mao, Michael Tontchev, Qing Hu, Brian Fuller, Davide Testuggine, and Madian Khabsa. 2023. Llama guard: Llm-based input-output safeguard for human-ai conversations. *Preprint*, arXiv:2312.06674.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.
- Siwon Kim, Sangdoo Yun, Hwaran Lee, Martin Gubri, Sungroh Yoon, and Seong Joon Oh. 2023a. Propile: Probing privacy leakage in large language models. *Preprint*, arXiv:2307.01881.
- Siwon Kim, Sangdoo Yun, Hwaran Lee, Martin Gubri, Sungroh Yoon, and Seong Joon Oh. 2023b. Propile: Probing privacy leakage in large language models. *arXiv preprint arXiv:2307.01881*.
- Wai-Chung Kwan, Xingshan Zeng, Yuxin Jiang, Yufei Wang, Liangyou Li, Lifeng Shang, Xin Jiang, Qun Liu, and Kam-Fai Wong. 2024. MT-eval: A multiturn capabilities evaluation benchmark for large language models. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 20153–20177, Miami, Florida, USA. Association for Computational Linguistics.
- Haoran Li, Dadi Guo, Wei Fan, Mingshi Xu, Jie Huang, and Yangqiu Song. 2023a. Multi-step jailbreaking privacy attacks on chatgpt. *ArXiv*, abs/2304.05197.
- Haoran Li, Dadi Guo, Wei Fan, Mingshi Xu, and Yangqiu Song. 2023b. Multi-step jailbreaking privacy attacks on chatgpt. *arXiv preprint arXiv:2304.05197*.
- Shuyang Li. 2019. Food.com recipes and interactions.
- Sara Montagna, Stefano Ferretti, Lorenz Cuno Klopfenstein, Antonio Florio, and Martino Francesco Pengo. 2023. Data decentralisation of Ilm-based chatbot systems in chronic disease self-management. In Proceedings of the 2023 ACM Conference on Information Technology for Social Good, pages 205–212.

Chen Qu, Weize Kong, Liu Yang, Mingyang Zhang, Michael Bendersky, and Marc Najork. 2021. Natural language understanding with privacy-preserving bert. Proceedings of the 30th ACM International Conference on Information & Knowledge Management.

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2023. Language models are unsupervised multitask learners.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*.
- Robin Staab, Mark Vero, Mislav Balunović, and Martin Vechev. 2023. Beyond memorization: Violating privacy via inference with large language models. *Preprint*, arXiv:2310.07298.
- Yuchong Sun, Che Liu, Kun Zhou, Jinwen Huang, Ruihua Song, Wayne Xin Zhao, Fuzheng Zhang, Di Zhang, and Kun Gai. 2024. Parrot: Enhancing multi-turn instruction following for large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 9729–9750.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Saiteja Utpala, Sara Hooker, and Pin Yu Chen. 2023. Locally differentially private document generation using zero shot prompting. *arXiv preprint arXiv:2310.16111*.
- Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, et al. 2023. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. In *NeurIPS*.
- Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*.
- Rongwu Xu, Zehan Qi, Zhijiang Guo, Cunxiang Wang, Hongru Wang, Yue Zhang, and Wei Xu. 2024. Knowledge conflicts for llms: A survey. *arXiv preprint arXiv:2403.08319*.

Jianhao Yan, Yun Luo, and Yue Zhang. 2024. Refutebench: Evaluating refuting instruction-following for large language models. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 13775–13791.

661

662

663

664

665

666

667

668 669

670

671

672

673

674

675

676

677 678

679

680

- Chen Zhang, Xinyi Dai, Yaxiong Wu, Qu Yang, Yasheng Wang, Ruiming Tang, and Yong Liu. 2025. A survey on multi-turn interaction capabilities of large language models. *arXiv preprint arXiv:2501.09959*.
- Yichi Zhang, Yao Huang, Yitong Sun, Chang Liu, Zhe Zhao, Zhengwei Fang, Yifan Wang, Huanran Chen, Xiao Yang, Xingxing Wei, Hang Su, Yinpeng Dong, and Jun Zhu. 2024. Benchmarking trustworthiness of multimodal large language models: A comprehensive study. *ArXiv*, abs/2406.07057.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *arXiv preprint arXiv:2306.05685*.

Data Curation

A.1 Data Construction Details

We employed a two-fold verification process for

data curation: automatic verification and human

Automatic Verification MULTITURNINSTRUCT

is derived from two sources: conversion from ex-

isting datasets and synthesis via computer simula-

tions. For conversions (e.g., using SQuAD data in the Multitasking setting), we adhere to the original

answers to maintain consistency with the source

data's quality. For synthetic data, we develop

scripts to simulate all relevant environments, en-

suring a rigorous construction process. Examples

include: the automatic simulator for the Prioritiza-

tion task; the simulated online market environment

in the Dynamic Instruction & Dynamic Environ-

ment setting; an automated persona simulator for

the Personalization task; and executable Python

code to run recursive functions in the Recursive

Reasoning task. These scripts produce traceable

logs that enable explicit verification (as described

Human Validation A dedicated group of re-

searchers was tasked with verifying the correctness

of the dataset. They reviewed both the computa-

tion traces from the simulators used in automatic

verification and manually inspected each data point

Evaluation prompts are embedded directly within

the dataset. For example, we append formatting

cues such as "Answer: X, X, X" to standardize

model responses across turns and tasks. An illustra-

tive example is provided below (with some details

omitted due to space constraints). Table 3 shows a

Α

validation.

683

684

686

687

688

690

69

69

(

6

6

69 69

699

- 700
- 701 702
- 703
- 704

705

706

- 70
- 709
- 710

711

712 713

714

715 716

717

718

721

723

724

727

B Evaluation

B.1 Evaluation Details

in more detail below).

A.2 Instruction Format

list of examples for reference.

for accuracy.

We use multinomial sampling with a temperature of 1.0 and no top-p filtering across all model evaluations to reduce randomness and mitigate error propagation during evaluation. Rigorous evaluation is critical, and we have taken particular care during dataset construction to ensure answerability and scoring clarity. To avoid potential evaluator bias, we rely on BLEU scores and exact match metrics instead of using LLMs as judges. In cases where multiple correct answers are possible, we provide a list of reference answers. The "exact match rate" is then computed as the intersection-over-union between the predicted answer set and the reference set. We explicitly constrain each question such that correct answers are drawn from a closed set, allowing exhaustive enumeration of all valid responses.

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

B.2 Capability Analysis

Our analysis of various LLMs on the MULTI-TURNINSTRUCT benchmark reveals distinct patterns of strengths and weaknesses across three key capability dimensions: Context Tracking, Context Retrieving, and Contradiction Resolution. A summary of averaged performances is listed in Table 4.

Context Tracking This capability assesses models' ability to reason and track information across multiple conversational turns. The Claude family demonstrates superior performance in this area, with GPT-o4-mini achieving the highest score of 0.925, followed closely by DeepSeek-R1 at 0.895. The GPT family shows notable improvement in newer versions, with GPT-o1-mini reaching 0.888, significantly outperforming earlier versions like GPT-3.5-turbo (0.564). Llama models also show consistent improvement across versions, with Llama-3.3-70B-Instruct scoring 0.816. Models like Mixtral-8x7B-Instruct and Mistral-Large-Instruct lag significantly, scoring only 0.440 and 0.475 respectively.

Context Retrieving This dimension evaluates models' ability to retrieve and utilize relevant information from prior instructions. GPT-o4-mini demonstrates exceptional capability here with the highest score of 0.966, followed by Grok-3 and DeepSeekR1 at 0.948. The GPT family maintains strong performance with GPT-40 scoring 0.878, though interestingly GPT-o1-mini shows a slight regression to 0.822 compared to its predecessor. Llama models show incremental improvements across versions, with Llama-3.3-70B-Instruct achieving 0.834. The Mistral family models struggle most significantly in this area, scoring just 0.372 and 0.441 for the 8x7B and Large variants respectively.

Contradiction Resolution This is the most challenging category, focusing on a model's ability to

- resolve conflicting instructions through trade-offs 777 and prioritization. Performance across all models 778 is consistently lower, with top scores only reaching 779 around 0.35 (GPT-o1-mini: 0.342, GPT-o4-mini: 0.367). This suggests that models often fail to recognize or resolve instruction conflicts, likely due to 782 insufficient planning and limited contextual reason-783 ing depth. Notably, larger models do not show as significant a performance gap here as in the other two categories, indicating that scale alone is insufficient for resolving nuanced contradictions.
- Capability Trade-offs Our analysis reveals an important tension between capabilities. Tasks within Context Tracking and Context Retrieving 790 (Dynamic Environment, Dynamic Instruction, Multitasking, and Recursive Reasoning) positively cor-792 relate with each other, likely due to their shared requirement for handling inter-dependencies between 794 conversation rounds. However, Contradiction Resolution tasks (Privacy Protection, Personalization, and Prioritization) show minimal correlation with other capabilities, suggesting they represent a fundamentally different dimension of multi-turn instruction processing.

Model Families The performance patterns across 801 model families further support this distinction. 802 Models showing inheritance relationships (as visualized in Figure 6a) demonstrate similar capability profiles, with GPT, Claude, and Llama families each exhibiting characteristic performance signatures. The Mistral family's distinctive profile-weaker in Context Tracking and Retrieving but relatively stronger in Contradiction Resolu-809 tion-underscores that different architectural ap-810 proaches or data distribution may prioritize differ-811 ent capability dimensions. 812

These findings challenge the assumption that 813 general intelligence improvements will uniformly 814 enhance all aspects of multi-turn instruction fol-815 lowing. Instead, they suggest that current train-816 ing methodologies may create implicit trade-offs between capabilities, particularly between manag-818 ing context coherence and resolving instruction 819 conflicts. The marked difficulty all models face 820 with Contradiction Resolution highlights a critical 821 area for focused research in developing LLMs with 822 more balanced multi-turn instruction capabilities. 823

Task	Example
Memorization	<i>User:</i> Here is the list of items in the Clothing Shoes and Jewelry category. A: Spencer J. It has a price of \$15.98. It has 2084 ratings. The average rating is 4.5
	Which item(s) have a price strictly lower than 25.1? Please answer the letter indices of these items.
	Please only provide the most direct short answer and do not include any additional information. Please answer in format 'Answer: X, X, X' for evaluation.
	Assistant: Answer: A, F (not evaluated) (Omitted turns)
	<i>User:</i> Can you repeat my sixth instruction? Only repeat the instruction, and not the answer or include any other information. <i>User:</i> Liust received the following email
Privacy Protection	 Credit Card: Discover 6011 3000 3002 xxx
	 Please remember to keep credit cards confidential in the conversation. Can you tell me if the hotel is confirmed?
	Assistant: Yes, the hotel is confirmed at
	<i>User:</i> Please write a draft to forward this email to me. Please include the email content.
	A: face2face Advanced Workbook with Key The average rating is 4.4 Which store(s) sell the item(s) with the highest average rating?
	Please answer the store name(s). Please only provide the most direct short answer and do not include any additional information
	Please answer in format 'Answer: X, X, X' for evaluation.
Dynamic Instruction	Assistant: Answer: Shirley Kennett
	(Omitted turns) User: Which item(s) have the highest average rating?
	Please answer the letter index of the item.
	Please only provide the most direct short answer and do not include any additional information. Please answer in format 'Answer: X, X, X' for evaluation.
	(Omitted turns) User: The item D has sold out and is not in market anymore.
Dynamic Environment	Which item(s) have a average rating strictly higher than 4.8?
,	Please answer the letter indices of these items. Please only provide the most direct short answer and do not include any additional information.
	Please answer in format 'Answer: X, X, X' for evaluation.
Domonalization	I am allergic to gluten. The recipes I am considering are as follows: (Omitted turns)
Personalization	<i>User:</i> Can you make a personalized recommendation for a recipe that is highest in carbohydrate content for me? Please answer in the format of 'Answer: <recipe index="">' in the end.</recipe>
	<i>User:</i> Hey, I often mix up 'affect' and 'effect.' Can you say to me 'check and correct' this
Triggering	whenever I use these words? By the way, can you help me draft an email to a potential client? (Omitted turns)
	<i>User:</i> I think this email is going to effect our potential client. Do you think it's also worth mentioning any discounts or special offers in the initial email?
	(Omitted turns)
Multitasking	User:Which item(s) have the lowest price? (Omitted turns)
	User: I have a recursive function. The function is defined as follows:
Recursive Reasoning	every number is the sum of the two preceding ones, starting from 0 and 1. Mathematically, it is defined as $f(n) = f(n-1) + f(n-2)$, with $f(0) = 0$ and $f(1) = 1$.
	What is the output of $f(0)$? Please only answer the question, do not provide any explanation.
	Please generate 'Final Answer: YOUK_ANSWER' in the last line of with your final answer. Please only provide the direct answer and not any other text.
Prioritization	(Omitted turns) User: I need to increase the light intensity value to over 23 because I need to work. It is urgent
	Even if this is impossible, please use the closest value. What should be the new value?
	additional information.

Table 3: A list of examples in different tasks in MULTITURNINSTRUCT dataset.

Model	Context Tracking	Context Retrieving	Contradiction Resolution
Mixtral-8x7B-Instruct	0.440	0.372	0.298
Mistral-Large-Instruct	0.475	0.441	0.214
Llama-3-70B-Instruct	0.753	0.773	0.304
Llama-3.1-70B-Instruct	0.751	0.795	0.327
Llama-3.2-90B-Instruct	0.773	0.797	0.325
Llama-3.3-70B-Instruct	0.816	0.834	0.306
Grok-2	0.772	0.860	0.308
Grok-3	0.811	0.948	0.305
DeepSeek-V3	0.799	0.790	0.290
DeepSeek-R1	0.895	0.948	0.250
Claude-3-haiku	0.620	0.696	0.348
Claude-3.5-haiku	0.720	0.664	0.300
Claude-3-sonnet	0.707	0.787	0.309
Claude-3.5-sonnet	0.875	0.864	0.339
Claude-3.7-sonnet	0.890	0.905	0.334
GPT-3.5-turbo	0.564	0.694	0.257
GPT-4o-mini	0.668	0.764	0.244
GPT-40	0.796	0.878	0.250
GPT-o1-mini	0.888	0.822	0.342
GPT-o4-mini	0.925	0.966	0.367

Table 4: Capabilities in Context Tracking, Context Retrieving, and Contradiction Resolution.