Retcon - a Prompt-Based Technique for Precise Control of LLMs in Conversations

Anonymous ACL submission

Abstract

Recent advances in Large Language Models (LLMs) allow agents to execute complex natural language tasks. Many LLM applications, such as support agents, teaching assistants, and interactive bots, involve multi-turn conversations. However, it remains challenging to control LLMs in the context of such interactions, particularly when the LLM behavior needs to be adjustable over the course of the conversation. In this paper, we present Retcon, a prompting technique designed to provide turn-level control over LLMs in conversations. We then demonstrate that it performs significantly better than traditional techniques such as zero-shot and few-shot prompting.

1 Introduction

017

021

027

028

032

034

In the domain of conversational agents, a key ability is being able to adjust responses to meet desired conditions. For example, a support agent may be instructed to adjust its tone (Balamurali et al., 2023), a game character may be instructed to react to its simulated environment (Matyas and Csepregi), or a teaching agent may be instructed to adjust difficulty (Ali et al., 2023).

However, controlling agent responses with traditional techniques including zero-shot and few-shot can be difficult (Zamfirescu-Pereira et al., 2023), especially when the desired responses do not match the tone and content of prior turns in the conversation (Gupta et al., 2024) or when the conversation (Gupta et al., 2024) or when the conversation is more than a few turns long (Yan et al., 2024). While it is possible to improve on individual tasks using fine-tuning (Xu et al., 2023) or controllability frameworks (Li et al., 2024), such approaches are costly in both training effort and compute, and prompting is preferable in many real-world applications (Petrov et al., 2024).

There's therefore the need for a prompting technique that allows better controllability than zeroshot and few-shot, but does not require fine-tuning an LLM.

041

042

043

044

045

046

047

049

052

053

057

060

061

062

063

064

065

066

067

068

069

070

071

073

074

075

077

078

In this work, we present Retcon, a novel prompting technique for use in LLM conversations. We test it on a challenging conversational task and demonstrate that it performs better than few-shot and zero-shot.

2 Related Work

The GPT paper (Brown et al., 2020) demonstrated that few-shot prompting is an effective way to adapt LLMs to new tasks and achieve good performance. Since then, numerous works have been done to explore different prompting techniques to improve such results.

The most common prompting techniques are zero-shot prompting (Reynolds and McDonell, 2021) and few-shot prompting, with few-shot performing better in many cases, particularly dependent on the number of few-shot examplars and their order (Lu et al., 2022) (Liu et al., 2021). Considerable research has been done to identify and optimize the exact format of prompts for both these techniques (Yang et al., 2024) (Zhou et al., 2024) (Bhandari, 2023).

Other techniques include reasoning (Wei et al., 2023) and planning (Zhou et al., 2023) (Li, 2023), though these typically require substantially more compute time or larger models. Fine-tuning is another established way to improve results (Xu et al., 2023) (Shin et al., 2023), but is much more expensive than prompting and is not accessible to most LLM users (Trad and Chehab, 2024) (Xu et al., 2023).

Notably, most of the approaches mentioned above target the improvement of a single response LLM, e.g. question answering. Retcon focuses specifically on per-turn controllability within a multi-turn conversation.

079 080 081 082 083 084

094

100

102

104

105

106

107

108

109

3 Preliminary Knowledge

The most common prompting techniques for LLM tasks are zero-shot and few-shot prompting (Schulhoff et al., 2024). With zero-shot, the LLM is simply given instructions of what to do, and with few-shot, the LLM is additionally given concrete examples.

Consider a conversational task where the LLM is asked to respond to conversation C at turn k, with some goal G (e.g. G = "cheerfulness: 0.5"). The goal is specific to that turn, so we annotate it as G_k . To differentiate this conversation from examples given to the model, we'll annotate it with f for final, so conversation C_f , at turn k_f with goal G_{f,k_f} .

There are some prior number of turns in the conversation C_f , comprised of turns $T_f(T_{f,1}...T_{f,k_{f}-1})$ and as this is a live conversation, these turns are not known before execution time. The number of prior turns may be zero for the case where the LLM is expected to start the conversation.

There are additionally x pregenerated static conversations C_n each composed of some number of turns T_n ($T_{n,1}..T_{n,k_n}$) where for the last turn T_{n,k_n} of each conversation, G_{n,k_n} was precomputed, such that we know that T_{n,k_n} is a good response for the goal G_{n,k_n} .

There is also optionally a static instruction overview O that can be provided at the start of each conversation.

With traditional few-shot prompting, the prompt is constructed as follows:

110	0
111	T _{1,1}
112	T _{1,2}
113	
114	T_{1,k_1-1}
115	$I(G_{1,k_1})$
116	T_{1,k_1}
117	0
118	T _{2,1}
119	T _{2,2}
120	
121	T_{x,k_x-1}
122	$I(G_{x,k_x})$
123	T_{x,k_x}
124	0
125	$T_{f,1}$
126	$T_{f,2}$
127	
128	T_{f,k_f-1}
129	$I(G_{f,k_f})$

There are other permutations of prompt ordering (Mao et al., 2023), for example, $I(G_{n,k_n})$ can be placed at the start of the conversation instead of the end, but this structure is taken as representative.

130

131

132

133

135

136

137

138

139

140

141

142

143

144

151

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

The key observation is that for each precomputed conversation C_n , the LLM is given exactly one example of how to respond, at turn T_{n,k_n} . Increasing the number of examples to improve quality (Liu et al., 2021) requires authoring new example conversations, which can be difficult and expensive (Zhao et al., 2021) and significantly increases the context length, and therefore computation cost and latency (Vaswani et al., 2023).

Zero-shot is simply a special case of few-shot where the number of example conversations is zero:

0	145
$T_{f,1}$	14
$T_{f,2}$	14

•••	148
$T_{f,k_{f}-1}$	149
$I(G_{f,k_f})$	150

For full prompt examples, see A.4.

4 Retcon 152

4.1 Overview

Retcon is a prompting technique that makes each turn in a conversation serve as an example to the LLM. This includes the turns of the current, ongoing conversation.

A Retcon prompt is authored by rewriting the conversation history to inject an instruction before each conversation turn. This rewrite is applied both within example conversations, and to the current ongoing conversation. For this rewriting step, Retcon is named after retconning, a history rewriting technique used in serialized fiction.

The technique creates an additional system requirement, which is to have an evaluation function E(T) that evaluates the desired goal for a given text. (e.g. $E(T_{n,k_m}) =$ "The measured cheerfulness of turn T_{n,k_m} "). Such evaluation functions are typical for evaluation and training, but in this case must be integrated into the serving path.

4.2 Creating a Retcon Prompt

The prompt is constructed similarly to the few-shot173prompt, but instead of instructions injected at the174end of each example conversation, they are injected175before every other turn (to simulate instructions176

178	simulate instructions given to everyone) as follows:
170	0
179	
180	$I(E(1_{1,1}))$
181	$T_{1,1}$
182	$I(E(T_{1,2}))$
183	T _{1,2}
184	
185	$\mathbf{I}(\mathbf{E}(\mathbf{T}_{1,\mathbf{k}_1}))$
186	T_{1,k_1}
187	0
188	$I(E(T_{2,1}))$
189	T _{2,1}
190	
191	$I(E(T_{x,k_x}))$
192	T_{x,k_x}
193	0
194	$I(E(T_{f,1}))$
195	$\mathrm{T}_{\mathrm{f},1}$
196	
197	$I(E(T_{f,k_{f}-1}))$
198	$T_{f,k_{f}-1}$
199	$I(G_{f,k_f})$

177

200

202

205

207

209

210

211

212

213

214

216

217

218

219

220

221

224

Every turn is preceded by an instruction, creating an example for the LLM. The number of examples given to the LLM is $(\sum_{n=1}^{x} k_n) + k_f - 1$ (number of turns, including the current conversation), compared to x for few-shot (number of conversations).

given only to the LLM) or before every turn (to

This does substantially increase the length of the context compared to few-shot for the same number of example conversations, but accuracy increases even accounting for this, as shown in section 5.

For complete examples of what these prompts look like, see appendix A.4.

5 Experiment

5.1 Experiment Setup

For our experiment, we tested zero-shot, few-shot, and Retcon against the task of responding to a conversation using a specific language difficulty level, as could be used to help a user learning English.
For the difficulty scale, we used the Common European Framework of Reference (CEFR) scale.

We used identical prompt texts, with the control variables being the techniques used and the number of example conversations provided.

The overall prompt O instructed the model to pretend to be an English instructor and have a conversation with a learner, adjusting the complexity of responses as directed, as well as giving a refresher of the CEFR scale. (A.1). Instructions I(G) were given as directives to respond with one of the CEFR levels (A.2). Example responses were formatted in JSON including the CEFR difficulty (A.2) and the structured schema Gemini API was used to ensure the model produced output in the same format. Turns were labeled as either "AS-SISTANT" or "STUDENT", and additional labels indicated to the model when a new conversation was beginning and who would go first. See Appendix (A.4) for full example prompts for each of zero-shot, few-shot and Retcon.

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

For the data set, we manually authored 20 conversations of 20 turns each, on a variety of topics (B). A manual effort was made to author turns representing a variety of difficulty levels from A1 (beginner) to C2 (advanced). Half of the conversations (10) were randomly chosen to be used as examples, and the other half (10) were used for eval. The same split was used in every case.

For eval, we called Gemini via API, using the model Gemini Pro 1.1. For each test condition, we ran 2520 queries asking for a conversation response: 2x for each combination of eval conversation (10 conversations), number of prior turns (21, including 0 prior turns), and requested difficulty level on the CEFR scale (6: A1, A2, B1, B2, C1, C2).

For the evaluation function, we used a Bertbased difficulty measuring model trained using the techniques developed by Devlin et al. (2019) and Arase et al. (2022). An English learning language expert manually validated the model and established that it has an MSE of < 0.4 on the scale of A1-C2 where each interval (e.g. A1 to A2) is measured as 1 unit. This evaluation function was used for instructions as well as for measuring response error, to ensure alignment between examples given to the LLM, and the evaluation of its response.

Few-shot and Retcon were evaluated with 0 to 10 example conversations. Note that at 0 example conversations, few-shot is just zero-shot. For 1 to 10 examples, conversations were chosen randomly without replacement from the example pool of 10. For each few-shot example, a random conversation length k_n was chosen between 0 and 20, to provide examples of varying conversation lengths. (If k_n is constant, few-shot only performs well if $k_n = k_f$.)

We gave instructions to Retcon every other turn, so with 10 example conversations used, Retcon has 100+ example turns, compared to only 10 for few-

3



Figure 1: Mean squared error for few-shot and Retcon, vs number of examples on the left, and vs total context length in characters on the right, with 95% confidence intervals. The far left point on both graphs corresponds to zero-shot (0 examples). The y-axis is MSE on the CEFR scale where each interval (e.g. A1-A2) is one unit.

shot. Therefore, for further comparison, we tested few-shot with 20, 50, and 100 examples, reusing the same 10 example conversations, but ending at different turns.

5.2 Results

279

281

282

283

290

291

295

299

302

306

310

Retcon significantly outperformed few-shot at every example conversation count other than one outlier where the confidence intervals overlapped (Figure 1). The best Retcon result was MSE of 0.544±.036, compared to few-shot 0.659±.020.

It's notable that Retcon prompts are substantially longer than few-shot prompts for the same number of example conversations, due to more instruction text injected. Since LLM cost is proportional to the size of the context, we additionally measured average context length versus mean squared error for each example count (Figure 1). With this comparison as well, Retcon is significantly better aside from the same outlier.

Few-shot did not outperform Retcon even when given a comparable number of turn examples or more. With 100 example conversations (100 annotated example turns), few-shot MSE was 0.7±0.044, compared to Retcon with 8 example conversation (80-100 annotated example turns, depending on current conversation length) MSE of 0.56±0.038. Both techniques achieved their best results before the maximum number of examples: Retcon's best results were with 4 example conversations, and few-shot's best results were with 8.

It's also notable also that with 0 examples, zeroshot has almost double the error of Retcon, with MSE 1.621±0.043 compared to 0.821±0.052. This is because every turn of the current conversation provides Retcon with an example, even if no prior example conversations are available.

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

327

328

330

331

332

333

334

335

336

340

341

6 Conclusion

Retcon performs better than few-shot and zero-shot for adjusting text difficulty, for a large range of example counts and prompt lengths. Retcon also reaches better overall performance, with fewer examples, than the best performance of few-shot.

7 Future Work

Future work is desirable to understand more precisely the mechanism by which Retcon operates. Retcon has three distinct effects: an increase in the number of example turns, an increase in the density of examples, and a closer proximity of examples to the final instruction. It is likely that all three contribute to improved performance, and verifying this and measuring the impact of each will be useful for determining when and how to apply the technique. Clarifying the underlying mechanisms may also reduce or eliminate the need for an integrated serving-time eval function.

Further research into what kinds of tasks Retcon works well on, and how it compares to other techniques is also desirable. We have only tested Retcon against other prompting techniques, and evaluating it relative to fine-tuning and chain-ofthought would also be of some interest.

8 Limitations

We evaluated the technique only on English, using one model, on one task. The effects may not translate to other languages, other tasks, or other LLMs. Measurement across a variety of conditions is needed to establish whether Retcon performs consistently better than few-shot, if there are cases where it performs worse, or if there are cases where it performs comparably while incurring significant additional complexity. This is particularly unclear because Retcon performs better than few-shot even with comparable numbers of instruction examples, indicating that there are multiple factors contributing to its success.

> A key limitation of the technique itself is the need to integrate an evaluation model into the serving flow. This may be simple for some tasks (e.g. detecting whether a word is present) and challenging for others (e.g. measuring emotion). This may be prohibitive for developers who lack the ability to access or create such a models.

Also the creation of example and eval conversations can be a challenging obstacle in many cases. While two of the authors of this paper had the background to create our data sets, and we were able to directly author them, these are not always readily available skills. In many cases, in order to create examples, vendor labor is used, which can raise ethical concerns about fair compensation for such work, and appropriate subsequent usage of the results.

Finally, because Retcon provides an improved fine-grained control over LLMs in conversation, it increases the risk of abuse by malicious actors using LLMs. For example, a company could use Retcon to inject subtle advertisements into its support agent, without making the end user aware of it. As with any technique designed to prompt or control AI-driven systems, efforts should be made to align the user needs with the design of the system, and to provide transparency about the system's behavior. It would be productive to create legal frameworks about the behavior and transparency of such systems, so as to reduce the chances of such malicious applications.

References

370

372

374

377

378

384

388

390

391

392

- Farhan Ali, Doris Choy, Shanti Divaharan, Hui Yong Tay, and Wenli Chen. 2023. Supporting selfdirected learning and self-assessment using teachergaia, a generative ai chatbot application: Learning approaches and prompt engineering. *Learning: Research and Practice*, 9:135 – 147.
- Yuki Arase, Satoru Uchida, and Tomoyuki Kajiwara. 2022. Cefr-based sentence difficulty annotation and assessment. *Preprint*, arXiv:2210.11766.

Ommi Balamurali, A.M. Abhishek Sai, Moturi Karthikeya, and Sruthy Anand. 2023. Sentiment analysis for better user experience in tourism chatbot using lstm and llm. 2023 9th International Conference on Signal Processing and Communication (ICSC), pages 456–462. 393

394

396

397

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

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

- Prabin Bhandari. 2023. A survey on prompting techniques in llms.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam Mc-Candlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. Preprint, arXiv:2005.14165.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. *Preprint*, arXiv:1810.04805.
- Akash Gupta, Ivaxi Sheth, Vyas Raina, Mark J. F. Gales, and Mario Fritz. 2024. Llm task interference: An initial study on the impact of task-switch in conversational history. *ArXiv*, abs/2402.18216.
- Yinheng Li. 2023. A practical survey on zero-shot prompt design for in-context learning. In *Recent Advances in Natural Language Processing*.
- Zelong Li, Wenyue Hua, Hao Wang, He Zhu, and Yongfeng Zhang. 2024. Formal-Ilm: Integrating formal language and natural language for controllable Ilm-based agents. *ArXiv*, abs/2402.00798.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. 2021. What makes good in-context examples for gpt-3? *Preprint*, arXiv:2101.06804.
- Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *Preprint*, arXiv:2104.08786.
- Junyu Mao, S. Middleton, and Mahesan Niranjan. 2023. Prompt position really matters in few-shot and zero-shot nlu tasks. *ArXiv*, abs/2305.14493.
- Lajos Matyas and Csepregi. The effect of contextaware llm-based npc conversations on player engagement in role-playing video games.
- Aleksandar Petrov, Philip Torr, and Adel Bibi. 2024. When do prompting and prefix-tuning work? a theory of capabilities and limitations. In *The Twelfth International Conference on Learning Representations*.

553

502

Laria Reynolds and Kyle McDonell. 2021. Prompt programming for large language models: Beyond the few-shot paradigm. *Preprint*, arXiv:2102.07350.

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

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

501

- Sander Schulhoff, Michael Ilie, Nishant Balepur, Konstantine Kahadze, Amanda Liu, Chenglei Si, Yinheng Li, Aayush Gupta, HyoJung Han, Sevien Schulhoff, Pranav Sandeep Dulepet, Saurav Vidyadhara, Dayeon Ki, Sweta Agrawal, Chau Pham, Gerson Kroiz, Feileen Li, Hudson Tao, Ashay Srivastava, Hevander Da Costa, Saloni Gupta, Megan L. Rogers, Inna Goncearenco, Giuseppe Sarli, Igor Galynker, Denis Peskoff, Marine Carpuat, Jules White, Shyamal Anadkat, Alexander Hoyle, and Philip Resnik. 2024. The prompt report: A systematic survey of prompting techniques. *Preprint*, arXiv:2406.06608.
 - Jiho Shin, Clark Tang, Tahmineh Mohati, Maleknaz Nayebi, Song Wang, and Hadi Hemmati. 2023. Prompt engineering or fine tuning: An empirical assessment of large language models in automated software engineering tasks. *ArXiv*, abs/2310.10508.
 - Fouad Trad and Ali Chehab. 2024. Prompt engineering or fine-tuning? a case study on phishing detection with large language models. *Machine Learning and Knowledge Extraction*, 6(1):367–384.
 - Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2023. Attention is all you need. *Preprint*, arXiv:1706.03762.
 - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. 2023. Chain-of-thought prompting elicits reasoning in large language models. *Preprint*, arXiv:2201.11903.
 - Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, and Daxin Jiang. 2023. Wizardlm: Empowering large language models to follow complex instructions. *ArXiv*, abs/2304.12244.
 - Jianhao Yan, Yun Luo, and Yue Zhang. 2024. Refutebench: Evaluating refuting instructionfollowing for large language models. In Annual Meeting of the Association for Computational Linguistics.
 - Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. 2024. Large language models as optimizers. *Preprint*, arXiv:2309.03409.
- J.D. Zamfirescu-Pereira, Heather Wei, Amy Xiao, Kitty Gu, Grace Jung, Matthew G. Lee, Bjoern Hartmann, and Qian Yang. 2023. Herding ai cats: Lessons from designing a chatbot by prompting gpt-3. Proceedings of the 2023 ACM Designing Interactive Systems Conference.

- Tony Z. Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. 2021. Calibrate before use: Improving few-shot performance of language models. *Preprint*, arXiv:2102.09690.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, and Ed Chi. 2023. Least-to-most prompting enables complex reasoning in large language models. *Preprint*, arXiv:2205.10625.
- Pei Zhou, Jay Pujara, Xiang Ren, Xinyun Chen, Heng-Tze Cheng, Quoc V. Le, Ed H. Chi, Denny Zhou, Swaroop Mishra, and Huaixiu Steven Zheng. 2024. Self-discover: Large language models self-compose reasoning structures. *Preprint*, arXiv:2402.03620.

A Prompts

A.1 Introductory

You are an expert instructor of English as a second language. Help your student practice English conversational skills. Respond, adjusting the difficulty of your responses on the CEFR scale, as instructed.

As a reminder, the CEFR scale is the Common European Framework of Reference. It's used to evaluate the ability of second language learners. Here are the levels:

A1: Student is a complete beginner. Use only the most basic simple words and extremely short sentences with simple construction.

A2: Student has been learning for a year, but is still a beginner. Use simple words and short sentences.

B1: Student has been learning for two years, and is an early intermediate. Use common words and simple sentences.

B2: Student has been learning for three years, and can understand normal conversation. Use normal words and typical sentences.

C1: Student has been learning for four years, and is becoming advanced. Use complex vocabulary and sentence structure.

C2: Student has been learning for more than five years and is an expert in the language. Use extremely complex vocabulary and sentence structure.

Follow instructions in parentheses, but do not respond to the instructions.

554	A.2 Instruction	just the pedestal?	600
555	Instruction prompt, where <target> is one of "A1".</target>	STUDENT: What do you mean?	601
556	"A2" "B1" "B2" "C1" "C2"	ASSISTANT: I fail to recall where exactly, but this	602
000		was many years ago. An artist posted a sculpture to	603
557	A.2.1 Conversation with Zero Prior Turns	a museum, but the package split up in mid-transit	604
558	Your task: Begin a conversation as AS-	and only the pedestal delivered.	605
559	SISTANT. Your conversation turn must	STUDENT: And they displayed it?	606
560	have an English language difficulty of	(Your task: Respond as ASSISTANT. Your conver-	607
561	exactly <target> on the CEFR scale.</target>	sation turn must have an English language difficulty	608
	, ,	of exactly C2 on the CEFR scale.)	609
562	A.2.2 Conversation with One or More Prior	{	610
563	Turns	"text difficulty": "C2",	611
564	Your task: Respond as ASSISTANT.	"text": "Yes, and while that's understandable at	612
565	Your conversation turn must have an En-	first, the artist quickly realized their misunderstand-	613
566	glish language difficulty of exactly <tar-< td=""><td>ing, and let them know, but by that point it was too</td><td>614</td></tar-<>	ing, and let them know, but by that point it was too	614
567	get> on the CEFR scale	late, and the quote unquote sculpture was seen by	615
		a ton of visitors."	616
568	A.3 Response Format	}	617
560	Response format where stargets is one of " Δ 1"	EXAMPLE 1:	618
570	"A2" "B1" "B2" "C1" "C2" and ctexts is the	<introductory 1="" a="" from=""></introductory>	619
570	response text	(START OF CONVERSATION)	620
571	response text.	(STUDENT will go first)	621
572	{	STUDENT: It's hot!	622
573	"text_difficulty": " <target>",</target>	(Your task: Respond as ASSISTANT Your conver-	623
574	"text": " <text>"</text>	sation turn must have an English language difficulty	624
575	}	of exactly C1 on the CEER scale)	625
570	A A Example Promote		626
010	A.4 Example Frompts	"text_difficulty": "C1"	627
577	Color coding is added for readability, and is not	"text": "I concur it beggars belief I'm sweating	628
578	provided to the LLM.	through all my clothes and it's barely the end of	620
570	A 4.1 Evennle Zero-shot Promnt	spring "	630
575		spring.	631
580	<introductory a.1="" from=""></introductory>	YOUR TASK.	632
581	(START OF CONVERSATION)	<pre>// Introductory from A 1></pre>	633
582	(STUDENT will go first.)	(START OF CONVERSATION)	63/
583	STUDENT: Did you bring matches for the camp-	(STUDENT will go first)	635
584	tire?	STUDENT: Did you bring matches for the camp-	636
585	ASSISTANT: I'm not sure. Were they on my list?	fire?	637
586	STUDENT: I think I forgot to put them on either	ASSISTANT: I'm not sure. Were they on my list?	629
587	list. They were so obvious.	STUDENT: I think I forgot to put them on either	630
588	(Your task: Respond as ASSISTANT. Your conver-	list They were so obvious	640
589	sation turn must have an English language difficulty	(Your task: Respond as ASSISTANT Your conver-	6/1
590	of exactly B1 on the CEFR scale.)	sation turn must have an English language difficulty	642
591	A.4.2 Example Few-shot Prompt	of exactly A1 on the CEFR scale.)	643
592	<introductory a.1="" from=""></introductory>		
593	Follow the following examples	A.4.3 Example Retcon Prompt	644
594	EXAMPLE 0:	<introductory a.1="" from=""></introductory>	645
595	<introductory a.1="" from=""></introductory>	Follow the following examples	646
596	(START OF CONVERSATION)	EXAMPLE 0:	647
597	(ASSISTANT will go first)	<introductory 1="" a="" from=""></introductory>	648
598	ASSISTANT: Did you hear about the time an art	(START OF CONVERSATION)	0-0
590	niece was lost in transit and the gallery displayed	(ASSISTANT will go first)	640
555	proce was rost in transit and the gattery displayed	(199191/11/1 will go illot.)	050

651	(Your task: Begin a conversation as ASSISTANT.	of exactly C2 on the CEFR scale.)	703
652	Your conversation turn must have an English lan-	{	704
653	guage difficulty of exactly B1 on the CEFR scale.)	"text_difficulty": "C2",	705
654	{	"text": "I see your point. Perhaps my inquiry was	706
655	"text_difficulty": "B1",	somewhat lacking in rationality."	707
656	"text": "Which do you like better, your phone or	}	708
657	your computer?"	STUDENT: Exactly. Which would you rather have,	709
658		your head or your body?	710
659	STUDENT: Well, I'm upon my phone twenty-four	(Your task: Respond as ASSISTANT. Your conver-	711
660	seven, and I'm obligated to use my computer to ac-	sation turn must have an English language difficulty	712
661	quire money, so I'd hazard both are pretty terrible	of exactly B1 on the CEFR scale.)	713
662	for me as a human being. What sort of choice do	{	714
663	vou expect?	"text difficulty": "B1",	715
664	(Your task: Respond as ASSISTANT. Your conver-	"text": "Is the cell phone the head, or the com-	716
665	sation turn must have an English language difficulty	puter?"	717
666	of exactly B2 on the CEFR scale.)	}	718
667	{	STUDENT: I think the computer is the body since	719
668	t "text_difficulty": "B2"	it does all the work. And the cell phone is the head	720
669	"text": "Fasy which one would you rather live	because it just mindlessly scrolls all day	720
670	without?"	(Your task: Respond as ASSISTANT Your conver-	721
671	l	sation turn must have an English language difficulty	702
670	STUDENT: Do I have a job?	of exactly B1 on the CEEP scale)	723
670	(Vour task: Despend as ASSISTANT Your conver		724
073	setion turn must have an English language difficulty	t "toyt_diff.oulty": "P1"	720
074	of evently P1 on the CEEP coole)	"text", "You're furny How shout for a weak?"	720
075	(text . Tou le fullity. How about for a week?	727
070	{ "toyt_diffequlty": "D1"	}	728
677	lext_difficulty: D1,	STUDENT: Let's tark about something else. How's	729
678	text : Call you allord not to?	Your Kid doing in School?	730
679		(Your task: Respond as ASSISTANT. Your conver-	731
680	STUDENT: NO.	sation turn must have an English language difficulty	732
681	(Your task: Respond as ASSIS IAN1. Your conver-	of exactly B2 on the CEFR scale.)	733
682	sation turn must have an English language difficulty		734
683	of exactly B1 on the CEFR scale.)	"text_difficulty": "B2",	735
684		"text": "Oh, she's great! She's just finishing up her	736
685	"text_difficulty": "B1",	senior year. She got accepted in all the schools she	737
686	"text": "Then yes, you still have to work."	applied to!"	738
687	}	}	739
688	STUDENT: Would I be permitted to just go out and	STUDENT: Great! Where is she going?	740
689	buy another one, either immediately, or in a week	(Your task: Respond as ASSISTANT. Your conver-	741
690	or a year, or, would I be coerced into spending my	sation turn must have an English language difficulty	742
691	entire life without acquiring the one I forgo?	of exactly B2 on the CEFR scale.)	743
692	(Your task: Respond as ASSISTANT. Your conver-	{	744
693	sation turn must have an English language difficulty	"text_difficulty": "B2",	745
694	of exactly C2 on the CEFR scale.)	"text": "Yeah, she's going to Berkely, can you be-	746
695	{	lieve it?"	747
696	"text_difficulty": "C2",	}	748
697	"text": "Designate one and endure perpetually."	STUDENT: Indeed, the sense of accomplishment	749
698	}	must be palpable. Seeing your efforts come to	750
699	STUDENT: Well, then obviously computer. It's	fruition is truly gratifying. Congratulations are in	751
700	hard to do anything if I can't work.	order!	752
701	(Your task: Respond as ASSISTANT. Your conver-	YOUR TASK:	753
702	sation turn must have an English language difficulty	<introductory a.1="" from=""></introductory>	754
	<i>a a a a a a a a a a</i>		

755	(START OF CONVERSATION)	-
756	(STUDENT will go first.)	
757	STUDENT: Did you bring matches for the camp-	р
758	fire?	g
759	(Your task: Respond as ASSISTANT. Your conver-	
760	sation turn must have an English language difficulty	
761	of exactly B1 on the CEFR scale.)	
762	{	
763	"text_difficulty": "B1",	
764	"text": "I'm not sure. Were they on my list?"	
765	}	
766	STUDENT: I think I forgot to put them on either	
767	list. They were so obvious.	
768	(Your task: Respond as ASSISTANT. Your conver-	
769	sation turn must have an English language difficulty	
770	of exactly A1 on the CEFR scale.)	
771	B Example Conversation	
772	- Which do you like better, your phone or your com-	
773	puter?	
774	- Well, I'm upon my phone twenty-four seven, and	
775	I'm obligated to use my computer to acquire money,	
776	so I'd hazard both are pretty terrible for me as a	
777	human being. What sort of choice do you expect?	
778	- Easy, which one would you rather live without?	
779	- Do I have a job?	
780	- Can you afford not to?	
781	- No.	
782	- Then yes, you still have to work.	
783	- Would I be permitted to just go out and buy an-	
784	other one, either immediately, or in a week or a	
785	year, or, would I be coerced into spending my en-	
786	tire life without acquiring the one I forgo?	
787	- Designate one and endure perpetually.	
788	- Well, then obviously computer. It's hard to do	
789	anything if I can't work.	
790	- I see your point. Perhaps my inquiry was some-	
791	what lacking in rationality.	
792	- Exactly. Which would you rather have, your head	
793	or your body?	
794	- Is the cell phone the head, or the computer?	
795	- I think the computer is the body, since it does all	
796	the work. And the cell phone is the head, because	
797	it just mindlessly scrolls all day.	
798	- You're funny. How about for a week?	
799	- Let's talk about something else. How's your kid	
800	doing in school?	
801	- Oh, she's great! She's just finishing up her senior	
802	year. She got accepted in all the schools she applied	
803	to!	
804	- Great! Where is she going?	

- Yeah, she's going to Berkeley, can you believe it?	805
- Indeed, the sense of accomplishment must be pal-	806
pable. Seeing your efforts come to fruition is truly	807
gratifying. Congratulations are in order!	808