

Beyond the Turn-Based Game: Enabling Real-Time Conversations with Duplex Models

Anonymous ACL submission

Abstract

As large language models (LLMs) increasingly permeate daily lives, there is a growing demand for real-time interactions that mirror human conversations. Traditional turn-based chat systems driven by LLMs prevent users from verbally interacting with the system while it is generating responses. To overcome these limitations, we adapt existing LLMs to *duplex models* so that these LLMs can listen for users while generating output and dynamically adjust themselves to provide users with instant feedback. Specifically, we divide the queries and responses of conversations into several time slices and then adopt a time-division-multiplexing (TDM) encoding-decoding strategy to pseudo-simultaneously process these slices. Furthermore, to make LLMs proficient enough to handle real-time conversations, we build a fine-tuning dataset consisting of alternating time slices of queries and responses as well as covering typical feedback types in instantaneous interactions. Our experiments show that although the queries and responses of conversations are segmented into incomplete slices for processing, LLMs can preserve their original performance on standard benchmarks with a few fine-tuning steps on our dataset. Automatic and human evaluation indicate that duplex models make user-AI interactions more natural and human-like, and greatly improve user satisfaction compared to vanilla LLMs. Our duplex model and dataset will be released.

1 Introduction

Large language models (LLMs) have demonstrated impressive capabilities in various scenarios (OpenAI, 2023b; Achiam et al., 2023; Touvron et al., 2023; Team et al., 2023). These large models are deeply integrated with our daily lives and their extraordinary capabilities can satisfy users in many applications, such as coding assistants (Chen et al., 2021; GitHub, 2023b,a; Microsoft, 2024;

Rozière et al., 2023; Li et al., 2023b), task assistants (Wang et al., 2023b; Qian et al., 2023; OpenAI, 2024), virtual role play (Shao et al., 2023; Shanahan et al., 2023), and even emotional companions (Chaturvedi et al., 2023; Guingrich and Graziano, 2023; Pentina et al., 2023).

Despite ongoing advancements, interactions with LLMs often fail to provide users human-like interaction experience (Hill et al., 2015; Mou and Xu, 2017; Zhou et al., 2023). One reason is the turn-based nature of current chatbot implementations (Skantze, 2021), which is different from human conversations where there are many overlaps, interruptions, and silences (Zimmerman and West, 1996). Current human-LLM interactions necessitate that one participant remains entirely idle while the other generates responses. Interruptions are manually triggered with a “stop” button or by saying certain keywords, resulting in conspicuously artificial communication. In human conversations, participants simultaneously process incoming information and formulate responses, often in overlapping and interleaved contexts, thus allowing each other to interrupt or be interrupted.

To address this limitation, we introduce the concept of **duplex models**. Duplex models emulate human cognitive processes by synthesizing responses internally while simultaneously attending to incoming user inputs, akin to a person thinking while listening as well as speaking while observing. However, present autoregressive models face substantial challenges in adopting a duplex configuration, as they must process and encode a complete input message before generating any tokens, resulting in a turn-based conversation. Considering this, we propose a framework for quickly converting current LLMs into duplex models by processing queries and responses pseudo-simultaneously without significant alternations to their architectures.

Specifically, we propose a time-division-multiplexing (TDM) encoding-decoding strategy.

083 messages in dialogues are split into time slices and
084 the model processes time slices of input queries
085 incrementally and generates time slices of output
086 responses based on these partial input slices. When
087 a new input query arrives, the model immediately
088 halts its current generation process and starts a new
089 sequence that integrates the additional input, en-
090 abling swift responses. To adapt existing LLMs to
091 this format of time slices, we build a duplex dataset
092 for fine-tuning. The differences between our data
093 from the conventional supervised fine-tuning (SFT)
094 dataset are: (1) its input and output are time slices
095 and (2) it includes various interactive user interrup-
096 tions, such as generation termination, regeneration,
097 and dialogue reset.

098 To demonstrate the feasibility of duplex mod-
099 els, we train a prototype named MiniCPM-duplex,
100 based on MiniCPM—a robust and lightweight
101 LLM (Hu et al., 2024). Empirical results show
102 that MiniCPM-duplex has its original performance
103 on general benchmarks while enabling dynamic
104 responses to user queries. Additionally, we con-
105 duct a user study to compare the MiniCPM-duplex
106 with the original MiniCPM. The results indicate
107 that duplex models show significant improvements
108 in responsiveness, human-likeness, and user satis-
109 faction. Our contributions are fourfold:

110 (1) We introduce and define the concept of du-
111 plex models, which are designed to generate output
112 simultaneously as they receive input.

113 (2) We propose a TDM encoding-decoding strat-
114 egy and a duplex-specific SFT dataset for imple-
115 menting duplex models.

116 (3) We confirm that segmenting time slices dur-
117 ing interactions does not compromise performance,
118 and notably enhances the responsiveness, human-
119 likeness, and overall satisfaction of conversations.

120 (4) We release the model and dataset and provide
121 a demo for users to experience firsthand.

122 2 Duplex Models

123 We define *duplex models* as models that can process
124 inputs and produce outputs simultaneously, and dy-
125 namically decide when to respond. It differs from
126 current LLMs-based chatbots where participants
127 must specify the end of inputs and only produce
128 outputs after processing the entire input. To convert
129 existing LLMs into duplex models, we split conver-
130 sation messages into time slices, and then propose
131 a TDM encoding-decoding mechanism to process
132 these slices. To enhance the processing of these

133 time slices, we further introduce duplex alignment
134 to adapt existing LLMs to duplex models.

135 2.1 Time-Division-Multiplexing 136 Encoding-Decoding

137 Current autoregressive language models struggle
138 to function as true duplex systems. During the
139 input phase, the LLM encodes the input into key-
140 value caches without generating any output. To
141 leverage autoregressive models in approximating
142 duplex models, we propose a TDM strategy. We
143 divide the conversation interaction into time slices
144 and process input slices immediately to produce
145 corresponding output slices.

146 Instead of requiring users to specify when the
147 model should respond, the duplex model infers re-
148 sponses after every k seconds, i.e., each time slice
149 spans k seconds. A special token (e.g., <idle>)
150 is used to indicate the model’s decision to remain
151 silent and wait for further inputs. If not used, the
152 generated slice is delivered to the user immediately.
153 This approach mimics human conversational pat-
154 terns more closely, as humans do not use special
155 tokens to signal the end of utterances and intuitively
156 determine the appropriate moments to respond to
157 inputs. Figure 1 illustrates the distinction between
158 duplex and conventional language models.

159 2.2 Time-Slicing Chunking

160 As shown in Figure 1, all the input queries and
161 output responses of conversations are in the slice
162 format. The size of slices has great implications
163 for the performance of a duplex model. Large slice
164 sizes result in greater response (or interruption)
165 latency, while smaller slice sizes may result in un-
166 necessarily long inputs (because some tokens are
167 added between the chunks). Our preliminary inves-
168 tigation and pilot experiments with our transformer-
169 based (Vaswani et al., 2017) models reveal that
170 time-slicing chunking at 2-second intervals bal-
171 ances response latency and user experience. As-
172 suming human beings usually speak 110-170 words
173 per minute¹, an appropriate size of time slices is
174 4-6 words.

175 2.3 Duplex Alignment

176 Normal LLMs are unable to handle time slices as
177 shown in Figure 2, so we need to fine-tune them
178 into duplex models. To achieve this, we construct
179 a duplex SFT duplex dataset.

¹<https://debatix.com/en/speech-calculator/>

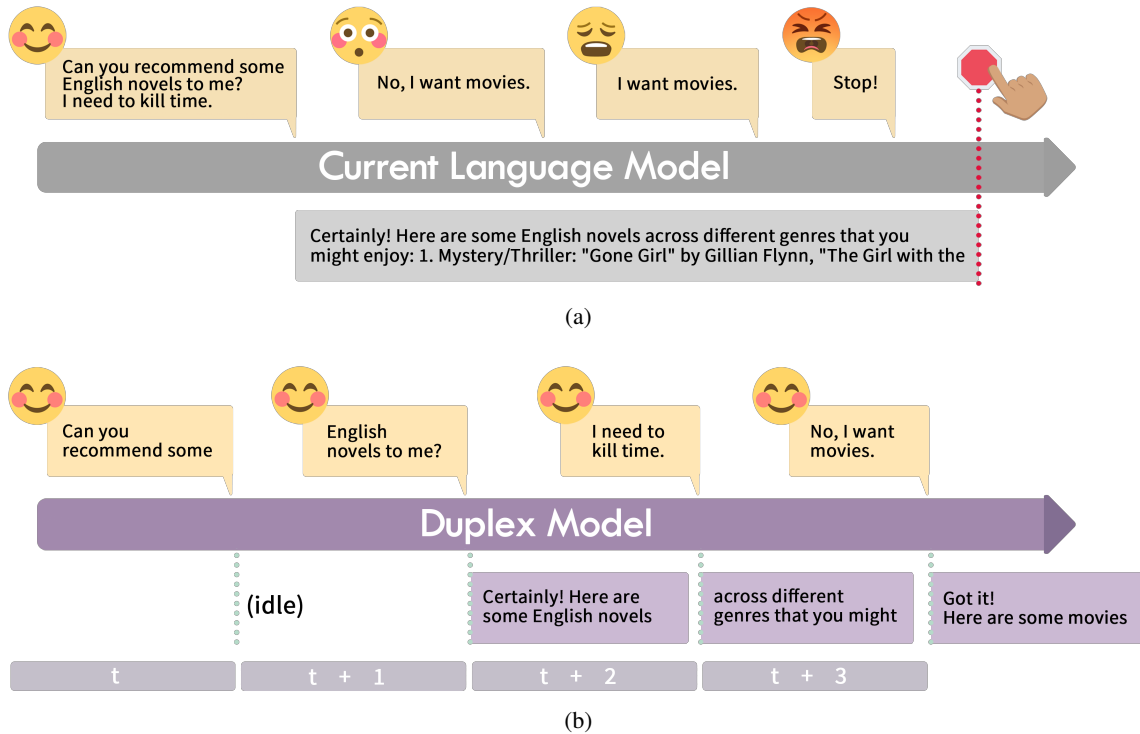


Figure 1: Illustration of the input/output processing scheme of traditional models (1a) and duplex models (1b). Traditional models receive the complete input from the user before generating the response. In contrast, duplex models process the input and generate the output in an online manner.

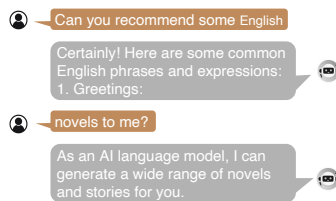


Figure 2: Responses of MiniCPM when inputs are time slices.

3 Supervised Fine-Tuning Duplex Dataset

We create **Duplex-UltraChat** for tuning current LLMs into duplex models. Different from existing dialogue datasets, Duplex-UltraChat has no special tokens or keywords to indicate the beginning or end of messages. Messages are split into time slices. A slice is either the actual message of an individual or a special “idle” token to indicate silence. Each individual may interrupt by generating a response before the other party’s message is completed.

Duplex-UltraChat is derived from UltraChat (Ding et al., 2023) to reduce annotation costs. We heuristically inject appropriate random interruptions to simulate realistic scenarios. Powerful LLMs rewrite the interruptions to ensure diversity and naturalness. Each user message is

randomly split into 4-6 words. Assistant messages are split into 10-token slices.

During the construction of the dataset, we abide by the following two design choices: user behavior is unpredictable and the assistant should be polite. Examples in the dataset can be categorized as uninterrupted dialogues and dialogues with interruptions. As shown in Table 1, there are six categories of duplex data consisting of over 4.8M dialogues. Each piece of data has an average length of 2,570.2 tokens encoded by the tokenizer of MiniCPM-duplex and 170.4 slice pairs.

3.1 Uninterrupted Dialogue

Basic Ordinary uninterrupted dialogue data is obtained by splitting existing dialogue messages into slices. When the user input is unfinished, the output of the duplex model should be <idle>. Meanwhile, when the duplex model is generating output, the user is set to quiet and its input is <idle>. Figure 3 shows an example of basic duplex data.

Topic Interweaving People may discuss several topics interweavably ignoring coherence. To mimic such behavior, we interlace sentences of 3-5 dialogues while keeping their orders, and split each sentence into time slices as the basic type does.

Example Type	# Dialogues	Avg. # Slice Pairs	Avg. # Tokens
Basic	1,458,353	153.9	2,342.2
Topic Interweaving	489,065	427.7	6819.6
Generation Termination	1,468,141	89.3	1,318.0
Regeneration	806,687	171.2	2,590.4
Dialogue Reset	300,318	194.7	2,906.5
Back on Topic	327,286	199.1	2495.6
Total	4,849,850	170.4	2,570.2

Table 1: The statistics of Duplex-UltraChat. The tokens are produced by the tokenizer of our MiniCPM-duplex.

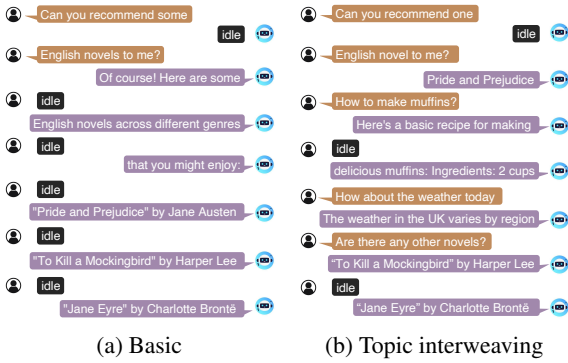


Figure 3: An example of uninterrupted dialogue in Duplex-UltraChat.

3.2 Dialogues with Interruptions

In realistic human conversions, the individuals may start speaking before the other part is done with their message. Therefore, to simulate such scenarios, we inject four interruptions into the data as shown in Figure 4.

Generation Termination Forced interruptions are when users directly speak out their next sentence regardless of the status of the assistant. To generate such data, we randomly choose a location in an assistant message, discard the remaining part of the message, and insert a new user input at that location. We prefix the user input with one of the 11 pre-defined transitional sentences (see Appendix A.1). This input is rewritten by ChatGPT to ensure a natural and varied transition. The target output is idle tokens because the assistant is expected to terminate its current response.

Generation termination requires the assistant to learn to stop speaking when the user is forcibly interrupting it and be robust to incomplete messages in the chat history. Since this interruption may be regarded as impolite, our dataset does not contain situations where the user is interrupted.

Regeneration Another scenario where the user interrupts the assistant is when the user is dissatisfied with the current response. In conventional LLM-based chatbots, the user must first stop the generation with a button, and prompt the model with the updated prompt. In contrast, duplex models allow the user to directly interrupt and reinput the new prompt while generating outputs. To create such data, we randomly sample a user message and repeat it with one of 15 pre-defined transition sentences (given in Appendix A.2). ChatGPT rewrites this repetition message for better coherence. Then, the chat history and repetition message are fed to ChatGPT to generate the annotation.

Dialogue Reset Here, we consider situations where the user wants to chat abruptly on an entirely different topic while the assistant is generating output. To create such data, we randomly sample five dialogues and truncate the first four dialogues at random locations before concatenation. We define 18 kinds of transitional sentences in Appendix A.3, including one empty string. We randomly choose a transitional sentence, and prefix it with the first sentence of the new dialogue. Each message is then rewritten by ChatGPT. If the selected transitional sentence is the empty string, we do not rewrite the input, which simulates certain users who wish to start a new dialogue as fast as possible.

Back on Topic When the user only interrupts a question without attempting to stop the assistant or change the topic, the assistant should answer the question and then continue the unfinished statement. To construct this type of data, we randomly select a within a message from the assistant, and annotate a question about a statement by the assistant. GPT-4 (Achiam et al., 2023) is used to generate the answer to the user’s question and continue the interrupted message with coherence.

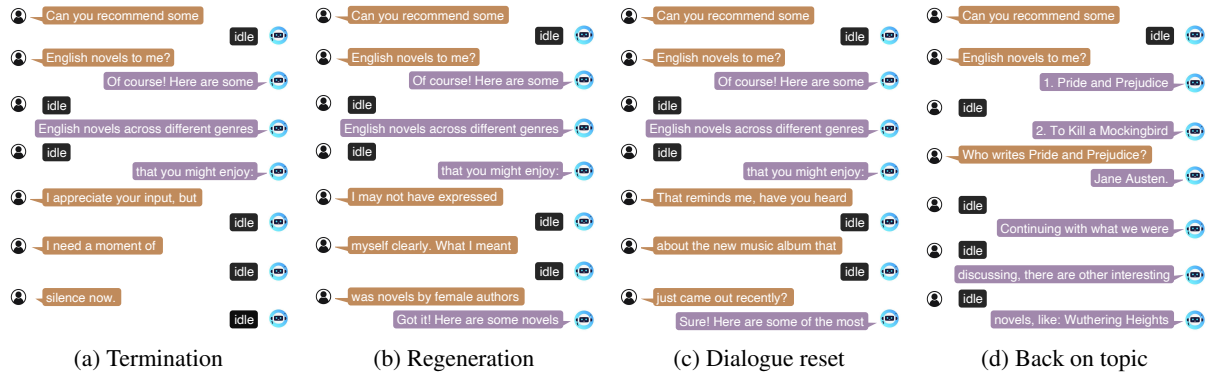


Figure 4: Some examples from Duplex-UltraChat.

4 Experimental Details

4.1 Training

We start from the public checkpoint of MiniCPM-2.4B (Hu et al., 2024)² and fine-tune it on Duplex-UltraChat as well as the SFT data that MiniCPM uses to obtain MiniCPM-duplex.

We make the following modifications to MiniCPM: (1) we append a special end-of-sentence token (i.e., `<eos>`) to each response of the duplex model, and (2) we add a special token `<idle>` to represent empty input or output.

The training of MiniCPM-duplex uses the following hyperparameters: 10^{-3} maximum learning rate, warmupstableexp (Hu et al., 2024) learning rate scheduler, a batch size of 800, and a maximum length of 4,096. We train for 10,000 steps on 40 NVIDIA A100 GPUs for 36 hours.

4.2 Baseline

Since our MiniCPM-duplex and MiniCPM are derived from the same checkpoint, we verify the effectiveness of our method by comparing it against the vanilla MiniCPM.

4.3 Evaluation

We evaluate the duplex model with three kinds of metrics: automatic metrics, GPT-4, and human. Automatic metrics, like accuracy and pass rate, are widely used for convenience and low cost.

4.3.1 GPT-4 Evaluation

To evaluate the multi-turn dialogue ability of MiniCPM-duplex, we benchmark it on MT-Bench (Zheng et al., 2024) with GPT-4 as the judge.

²<https://huggingface.co/openbmb/MiniCPM-2B-sft-bf16>, denoted MiniCPM.

To mimic real-time scenarios, we chunk each instruction in MT-Bench into multiple 4-6 word slices and feed one slice at a time. Then we concatenate all output segments from the duplex model to form the final output. For the traditional model, we directly feed the entire prompt to the model.

Both models use the same decoding parameters: random sampling, a temperature of 0.8, a top- p value of 0.8, and a top- k value of 0. The maximum length is set to 4,096. For the duplex model, we set the maximum token generated per chunk to 10.

4.3.2 Human Evaluation

When using humans as evaluators, we consider the following four aspects.

Responsiveness This metric measures whether a model will respond to a user query and the latency if it responds, which is a perceived latency. Many factors may contribute to greater response latency, including the speech-to-text and text-to-speech conversion time, model inference time, network latency, and the interaction strategy that the model utilizes. There is no obvious difference between the actual inference latency of MiniCPM-duplex and MiniCPM.

Human-Likeness Inspired by the Turing test, we wish to develop a language model that chats in a way indistinguishable from humans. Therefore, we define human-likeness as a metric that measures the degree of the similarity of a model to humans.

Faithfulness Faithfulness is a widely used metric in the evaluation of LLMs (Arras et al., 2017; Serrano and Smith, 2019; Jain and Wallace, 2019; DeYoung et al., 2020; Adlakha et al., 2023; Chen et al., 2023b). Here, we use it to reflect the degree how the model follows a user’s instruction, which

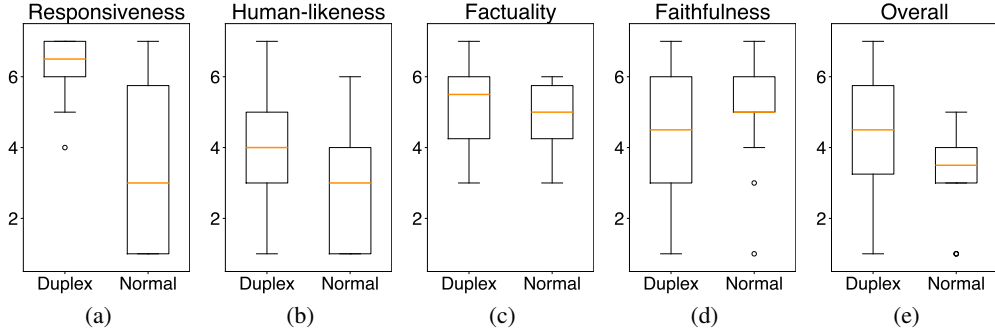


Figure 5: The human evaluation score distributions for MiniCPM and MiniCPM-duplex regarding responsiveness, human-likeness, factuality, faithfulness, and overall satisfaction.

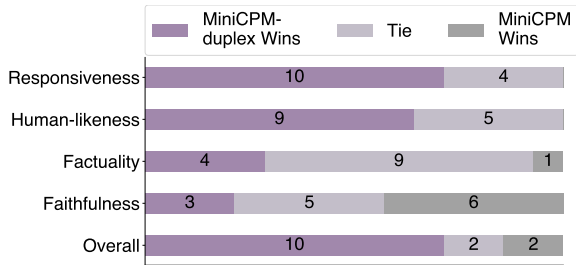


Figure 6: Win rates between MiniCPM and MiniCPM-duplex on responsiveness, human-likeness, factuality, faithfulness, and overall satisfaction.

is similar to (Adlakha et al., 2023).

Factuality This metric measures the degree of hallucination of a LLM (Rudinger et al., 2018; Tian et al., 2023; Chen et al., 2023a; Wang et al., 2023a; Nakano et al., 2021).

4.4 Interactive Demo

We implement an interactive demo with a user interface such that human evaluators can make evaluations based on actual interaction experience. In the demo, users chat with an assistant using voice. The assistant is either implemented with the vanilla MiniCPM or our MiniCPM-duplex. The conversion between speech and text is implemented with Google’s cloud-based ASR and TTS API³.

This demo supports both vanilla MiniCPM and MiniCPM-duplex. For the vanilla MiniCPM, the program automatically detects pauses in the user’s voice. On each pause, the speech is converted to text, which is then sent to the model. MiniCPM performs regular text generation, and each output token is passed to the ASR module, before being re-

³Speech-to-text API: <https://cloud.google.com/speech-to-text/docs/reference/rest>. Text-to-speech API: <https://cloud.google.com/text-to-speech/docs/reference/rest>.

turned to the user. Meanwhile, the user has to wait until the speech response is done before the next query. When interacting with MiniCPM-duplex, the user’s speech is processed every 2 seconds. When the MiniCPM-duplex does not generate the idle token, the text generation will be transcribed into audio and played out. The user’s voice will be captured, transcribed, and fed to the model regardless of whether the assistant speaks.

Benchmark	MiniCPM	MiniCPM-duplex
C-Eval	50.52	50.06
CMMLU	51.30	48.53
MMLU	53.45	53.76
BBH	37.25	36.35
HumanEval	50.00	49.39
MBPP	38.09	38.30
GSM8K	42.30	46.10
MATH	10.56	9.32
ARC-e	84.60	85.19
ARC-c	69.80	70.05
HellaSwag	61.40	60.79

Table 2: Performances of MiniCPM and MiniCPM-duplex on standard benchmarks.

Metric	MiniCPM	MiniCPM-duplex
Responsiveness	3.43	6.21
Human-Likeness	2.79	4.00
Factuality	4.93	5.21
Faithfulness	5.14	4.50
Overall	3.29	4.36

Table 3: Average human evaluation scores on responsiveness, human-likeness, factuality, faithfulness, and overall satisfaction. Higher is better.

Score	MiniCPM	MiniCPM-duplex
First turn	7.17	5.83
Second turn	5.85	4.84
Avg.	6.51	5.33

Table 4: MT-bench results of MiniCPM and MiniCPM-duplex. Higher is better.

4.5 User Study

Specifically, we recruit 14 participants consisting of 5 males and 9 females from 18 to 35 years old. Each participant holds a Bachelor’s or Master’s degree. Details on employment, payment, and ethical review are in Appendix C.

During the experiment, we rename MiniCPM-duplex as Model A, and MiniCPM as Model B to ensure anonymity. Participants are unaware of the difference between the two models beforehand. We specify the odd-numbered participants interact with Model A first, and the even-numbered ones first chat with Model B to eliminate the influence of chatting order. When finishing chatting with a model, the participant should score it and continue interacting with the other one. After the experiment, participants could modify and confirm scores for both models. Each participant is assigned at least 5 sessions of multi-turn dialogues with each model. The first sentence of sessions should be the same for both models. To help the participants come up with topics to chat about, we provide them with a reference note containing sample instructions from AlpacaEval (Li et al., 2023c).

Questionnaire Design The questionnaire consists of six questions. The first five questions prompt the user to rate the model based on responsiveness, human-likeness, faithfulness, factuality, and overall experience. The answer choices for these questions are scores from 1 to 7, where 1 represents disappointment, 4 represents indifference, and 7 represents excellence. The final question is open to suggestions on improving our duplex model. The actual questions are listed in Appendix B.2.

5 Results

Standard Benchmarks MiniCPM-duplex is benchmarked on several standard benchmarks, including multitask (C-Eval (Huang et al., 2024), CMMLU (Li et al., 2023a), MMLU (Hendrycks et al., 2020), BBH (Suzgun et al., 2023)), code

(HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021)), math (GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021)), and reasoning (ARC-e, ARC-c (Clark et al., 2018), HelLaSwag (Zellers et al., 2019)) with the LLM evaluation platform, UltraEval (He et al., 2024). Table 2 indicates that adapting to duplex models does not significantly harm its performance on general benchmarks.

GPT-4 Evaluation Table 4 shows the GPT-4 evaluation results on MT-Bench. MiniCPM-duplex is slightly inferior to MiniCPM mainly due to that MiniCPM-duplex tends to generate shorter responses. GPT-4 favors longer responses, whereas users prefer chat models that give concise answers.

Human Evaluation We have received 14 questionnaire. Table 3 lists the average scores of both models on five metrics. The duplex model surpasses the normal model by 81.05%, 43.37%, and 32.52% on responsiveness, human-likeness, and overall experience respectively.

Apart from absolute scores, we compare the ratings of the two models and count the number of evaluators that rate one model higher. The comparison results are shown in Figure 6. MiniCPM is more faithful than the duplex model mainly because it uses more diverse SFT data. Whereas the duplex model wins in other aspects, with an exceptionally large margin on responsiveness and human-likeness.

From these results, we conclude that duplex models can provide a better user experience in acting as the backbone model in AI assistants compared to ordinary language models.

6 Analysis & Discussion

6.1 Analysis

The superior performance of the duplex model is mainly due to its underlying receive/generate mechanism. Rather than strictly turn-based dialogue where users must explicitly signal the beginning and end of messages, duplex models behave more like human beings. Besides, the duplex model has learned when to speak at the fine-tuning stage on the Duplex-UltraChat, which makes it more human-like. Such ability is essential in passing a non-turn-based version of the Turing test, which is a more realistic test for whether a machine can be indistinguishable from humans.

468	6.2 Discussions	
469	We highlight some important open problems associated with duplex models below.	
470		
471	High-quality duplex data is urgently needed	
472	Existing dialogue datasets are inherently turn-based, which does not represent realistic and complex human conversations. Despite some success in empirical results with our synthetically generated duplex dataset, it still lags behind the practical demands. Two out of the 14 participants pointed out that they preferred concise responses rather than tedious answers.	
473		
474		
475		
476		
477		
478		
479		
480	We manually inspect 10 out of 90 chat sessions and find that the duplex model fails to remain silent once and interrupts the user unexpectedly once, showing that there is room for improvement. Thus, high-quality duplex datasets are in urgent need.	
481		
482		
483		
484		
485	A new decoding strategy is needed to improve the chat experience	
486	There are failed cases where the duplex model interrupted users unexpectedly. Balancing response speed and user experience is an open problem. Besides, to be more human-like, the duplex model should learn to start dialogues or topics actively.	
487		
488		
489		
490		
491		
492	A custom TTS system is needed to smooth the output voice	
493	The duplex model generates output chunk by chunk, which causes the output voice to be chunked. This results in incoherent intonation and volume, harming the user experience because existing TTS software does not support transcribing sequentially provided text chunks into a contiguous smooth voice. Overcoming this problem will improve the user experience considerably.	
494		
495		
496		
497		
498		
499		
500		
501	7 Related Work	
502		
503	7.1 Dialogue Dataset	
504	Dialogue data can be divided into two categories: single-turn and multi-turn.	
505		
506	Single-Turn Self-instruct (Wang et al., 2023c) is a synthetic instruction-following dataset of over 82K instances generated by GPT-3.5. Taori et al. (2023) adopt the data construction pipeline from Wang et al. (2023c) and construct Alpaca, a dataset with 52K instances. GPT-4-LLM (Peng et al., 2023) improves the Alpaca by replacing the data generator with GPT-4. It also adopts a Chinese version of Alpaca and Unnatural Instructions (Honovich et al., 2023). Besides, there are several high-	
507		
508		
509		
510		
511		
512		
513		
514		
	quality datasets, such as BELLE (Ji et al., 2023) and GPT-4ALL (Anand et al., 2023), among others.	515
		516
	Multi-Turn DailyDialog (Li et al., 2017) consists of over 13K dialogues annotated by humans, covering diverse daily conversation scenarios. Baize (Xu et al., 2023) generates multi-turn dialogues with ChatGPT by a prompting framework called self-chat where seed questions are from Quora and Stack Overflow, two popular question-answering websites. SODA (Kim et al., 2022) contains dialogues involving social commonsense. UltraChat (Ding et al., 2023) focuses on 30 meta-concepts and 20 types of materials and consists of over 1.4M dialogues.	517
		518
		519
		520
		521
		522
		523
		524
		525
		526
		527
		528
	7.2 Dialogue Models	529
	Chat-based models have gained widespread popularity since the release of ChatGPT. Some notable chat-based LLMs include the Claude series (Anthropic, 2023, 2024), Qwen series (Qwen, 2024), the Mistral series (Jiang et al., 2023) and LLaMa series (Touvron et al., 2023), among others. Most of these models, especially open-sourced ones, are purely text-based.	530
		531
		532
		533
		534
		535
		536
		537
	To enhance user experience, several applications support voice interaction. One instance is ChatGPT, where users press a button before speaking and indicate the end of speech with a button or pausing (OpenAI, 2023a). Then ChatGPT processes the received signal and produces a response until it finishes or users interrupt it by pressing a button. Such an implementation is unrealistic because it requires the user to specify the beginning and end of inputs. Whereas, our MiniCPM-duplex may improve this interactive experience by teaching the model to learn when to speak and when to be silent.	538
		539
		540
		541
		542
		543
		544
		545
		546
		547
		548
		549
	8 Conclusion	550
	We have introduced the concept of duplex models and provided one implementation. To this end, we also constructed the first non-turn-based dialogue dataset, Duplex-UltraChat, by injecting diverse kinds of interruptions into existing dialogue datasets. Our model, MiniCPM-duplex, is competitive with traditional models when evaluated on ordinary benchmarks while outperforming them in terms of responsiveness, human-likeness, and overall satisfaction. We believe that this work represents an essential step toward building machines that behave more human-like beyond current turn-based conversations.	551
		552
		553
		554
		555
		556
		557
		558
		559
		560
		561
		562
		563

564
565
566
567
568
569
570
571
572
573
574

Limitations

In this paper, we propose and verify the viability of duplex models. However, our implementation, MiniCPM-duplex, is a pseudo-duplex model, since it cannot perform encoding and decoding simultaneously. Consequently, our fixed-interval decoding strategy introduces a new hyperparameter that compromises responsiveness and context length (as discussed in Section 2.2). These limitations call for a new architecture that better supports the input-output scheme of duplex models.

References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*. 576
577
578
579
580

Vaibhav Adlakha, Parishad BehnamGhader, Xing Han Lu, Nicholas Meade, and Siva Reddy. 2023. Evaluating correctness and faithfulness of instruction-following models for question answering. *arXiv preprint arXiv:2307.16877*. 581
582
583
584
585

Yuvanesh Anand, Zach Nussbaum, Brandon Duderstadt, Benjamin Schmidt, and Andriy Mulyar. 2023. Gpt4all: Training an assistant-style chatbot with large scale data distillation from gpt-3.5-turbo. *GitHub*. 586
587
588
589

Anthropic. 2023. Introducing claude 2.1. <https://www.anthropic.com/news/claude-2-1>. 590
591

Anthropic. 2024. Introducing the next generation of claude. <https://www.anthropic.com/news/claude-3-family>. 592
593
594

Leila Arras, Franziska Horn, Grégoire Montavon, Klaus Robert Müller, and Wojciech Samek. 2017. " what is relevant in a text document?": An interpretable machine learning approach. *PLoS one*, 12(8):e0181142. 595
596
597
598

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*. 599
600
601
602
603

Rijul Chaturvedi, Sanjeev Verma, Ronnie Das, and Yogesh K. Dwivedi. 2023. [Social companionship with artificial intelligence: Recent trends and future avenues](#). *Technological Forecasting and Social Change*, 193:122634. 604
605
606
607
608

Liang Chen, Yang Deng, Yatao Bian, Zeyu Qin, Bingzhe Wu, Tat-Seng Chua, and Kam-Fai Wong. 2023a. Beyond factuality: A comprehensive evaluation of large language models as knowledge generators. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6325–6341. 609
610
611
612
613
614
615

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*. 616
617
618
619
620
621

Yijie Chen, Yijin Liu, Fandong Meng, Yufeng Chen, Jinan Xu, and Jie Zhou. 2023b. Improving translation faithfulness of large language models via augmenting instructions. *arXiv preprint arXiv:2308.12674*. 622
623
624
625

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*. 626
627
628
629
630

631	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,	Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu	687
632	Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias	Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxi-	688
633	Plappert, Jerry Tworek, Jacob Hilton, Reiichiro	ang Huang, Weilin Zhao, et al. 2024. Minicpm:	689
634	Nakano, et al. 2021. Training verifiers to solve math	Unveiling the potential of small language models	690
635	word problems. <i>arXiv preprint arXiv:2110.14168</i> .	with scalable training strategies. <i>arXiv preprint</i>	691
		<i>arXiv:2404.06395</i> .	692
636	Jay DeYoung, Sarthak Jain, Nazneen Fatema Rajani,	Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei	693
637	Eric Lehman, Caiming Xiong, Richard Socher, and	Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu,	694
638	Byron C Wallace. 2020. Eraser: A benchmark to	Chuanheng Lv, Yikai Zhang, Yao Fu, et al. 2024.	695
639	evaluate rationalized nlp models. In <i>Proceedings</i>	C-eval: A multi-level multi-discipline chinese evalua-	696
640	<i>of the 58th Annual Meeting of the Association for</i>	tion suite for foundation models. <i>Advances in Neural</i>	697
641	<i>Computational Linguistics</i> , pages 4443–4458.	<i>Information Processing Systems</i> , 36.	698
642	Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin,	Sarthak Jain and Byron C. Wallace. 2019. Attention is	699
643	Shengding Hu, Zhiyuan Liu, Maosong Sun, and	not Explanation . In <i>Proceedings of the 2019 Con-</i>	700
644	Bowen Zhou. 2023. Enhancing chat language models	<i>ference of the North American Chapter of the Asso-</i>	701
645	by scaling high-quality instructional conversations.	<i>ciation for Computational Linguistics: Human Lan-</i>	702
646	In <i>Proceedings of the 2023 Conference on Empiri-</i>	<i>guage Technologies, Volume 1 (Long and Short Pa-</i>	703
647	<i>cal Methods in Natural Language Processing</i> , pages	<i>pers)</i> , pages 3543–3556, Minneapolis, Minnesota.	704
648	3029–3051.	Association for Computational Linguistics.	705
649	GitHub. 2023a. About github copilot	Yunjie Ji, Yong Deng, Yan Gong, Yiping Peng, Qiang	706
650	chat. https://docs.github.com/	Niu, Lei Zhang, Baochang Ma, and Xiangang Li.	707
651	en/copilot/github-copilot-chat/	2023. Exploring the impact of instruction data	708
652	about-github-copilot-chat .	scaling on large language models: An empirical	709
653	GitHub. 2023b. Copilot. https://github.com/	study on real-world use cases. <i>arXiv preprint</i>	710
654	features/copilot .	<i>arXiv:2303.14742</i> .	711
655	Rose Guingrich and Michael SA Graziano. 2023. Chat-	Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-	712
656	bots as social companions: How people perceive con-	sch, Chris Bamford, Devendra Singh Chaplot, Diego	713
657	sciousness, human likeness, and social health benefits	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	714
658	in machines. <i>arXiv preprint arXiv:2311.10599</i> .	laume Lample, Lucile Saulnier, et al. 2023. Mistral	715
659	Chaoqun He, Renjie Luo, Shengding Hu, Yuanqian	7b. <i>arXiv preprint arXiv:2310.06825</i> .	716
660	Zhao, Jie Zhou, Hanghao Wu, Jiajie Zhang, Xu Han,	Hyunwoo Kim, Jack Hessel, Liwei Jiang, Peter West,	717
661	Zhiyuan Liu, and Maosong Sun. 2024. Ultraeval:	Ximing Lu, Youngjae Yu, Pei Zhou, Ronan Le Bras,	718
662	A lightweight platform for flexible and compre-	Malihe Alikhani, Gunhee Kim, Maarten Sap, and	719
663	hensive evaluation for llms. <i>arXiv preprint</i>	Yejin Choi. 2022. Soda: Million-scale dialogue dis-	720
664	<i>arXiv:2404.07584</i> .	tillation with social commonsense contextualization .	721
665	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,	In <i>Proceedings of the 2022 Empirical Methods in</i>	722
666	Mantas Mazeika, Dawn Song, and Jacob Steinhardt.	<i>Natural Language Processing</i> .	723
667	2020. Measuring massive multitask language under-	Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai	724
668	standing. In <i>International Conference on Learning</i>	Zhao, Yeyun Gong, Nan Duan, and Timothy Bald-	725
669	<i>Representations</i> .	win. 2023a. Cmmlu: Measuring massive multitask	726
670	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul	language understanding in chinese. <i>arXiv preprint</i>	727
671	Arora, Steven Basart, Eric Tang, Dawn Song, and	<i>arXiv:2306.09212</i> .	728
672	Jacob Steinhardt. 2021. Measuring mathematical	Raymond Li, Yangtian Zi, Niklas Muennighoff, Denis	729
673	problem solving with the math dataset. In <i>Thirty-</i>	Kocetkov, Chenghao Mou, Marc Marone, Christo-	730
674	<i>fifth Conference on Neural Information Processing</i>	pher Akiki, LI Jia, Jenny Chim, Qian Liu, et al. 2023b.	731
675	<i>Systems Datasets and Benchmarks Track (Round 2)</i> .	Starcoder: may the source be with you! <i>Transactions</i>	732
676	Jennifer Hill, W Randolph Ford, and Ingrid G Farreras.	<i>on Machine Learning Research</i> .	733
677	2015. Real conversations with artificial intelligence:	Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori,	734
678	A comparison between human–human online conver-	Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and	735
679	sations and human–chatbot conversations. <i>Comput-</i>	Tatsunori B. Hashimoto. 2023c. AlpacaEval: An	736
680	<i>ers in human behavior</i> , 49:245–250.	automatic evaluator of instruction-following models.	737
681	Or Honovich, Thomas Scialom, Omer Levy, and Timo	https://github.com/tatsu-lab/alpaca_eval .	738
682	Schick. 2023. Unnatural instructions: Tuning lan-	Yanran Li, Hui Su, Xiaoyu Shen, Wenjie Li, Ziqiang	739
683	guage models with (almost) no human labor. In <i>Pro-</i>	Cao, and Shuzi Niu. 2017. DailyDialog: A manually	740
684	<i>ceedings of the 61st Annual Meeting of the Associ-</i>	labelled multi-turn dialogue dataset . In <i>Proceedings</i>	741
685	<i>ation for Computational Linguistics</i> , pages 14409–	<i>of the Eighth International Joint Conference on Nat-</i>	742
686	14428.	<i>ural Language Processing</i> , pages 986–995, Taipei,	743

744	Taiwan. Asian Federation of Natural Language Processing.	Sofia Serrano and Noah A. Smith. 2019. <i>Is attention interpretable?</i> In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2931–2951, Florence, Italy. Association for Computational Linguistics.	795
745			796
746	Microsoft. 2024. Write code without the keyboard. https://githubnext.com/projects/copilot-voice/ .		797
747			798
748			799
749	Yi Mou and Kun Xu. 2017. The media inequality: Comparing the initial human-human and human-ai social interactions. <i>Computers in Human Behavior</i> , 72:432–440.	Murray Shanahan, Kyle McDonell, and Laria Reynolds. 2023. Role play with large language models. <i>Nature</i> , 623(7987):493–498.	800
750			801
751			802
752			
753	Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. 2021. Webgpt: Browser-assisted question-answering with human feedback. <i>arXiv preprint arXiv:2112.09332</i> .	Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. Character-llm: A trainable agent for role-playing. In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> .	803
754			804
755			805
756			806
757		Gabriel Skantze. 2021. Turn-taking in conversational systems and human-robot interaction: a review. <i>Computer Speech & Language</i> , 67:101178.	807
758			808
			809
759	OpenAI. 2023a. Chatgpt can now see, hear, and speak. https://openai.com/blog/chatgpt-can-now-see-hear-and-speak .	Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, et al. 2023. Challenging big-bench tasks and whether chain-of-thought can solve them. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 13003–13051.	810
760			811
761			812
762	OpenAI. 2023b. Introducing chatgpt. https://openai.com/blog/chatgpt#OpenAI .		813
763			814
764	OpenAI. 2024. Hello gpt-4o. https://openai.com/index/hello-gpt-4o/ .	Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2023. Alpaca: A strong, replicable instruction-following model. https://crfm.stanford.edu/2023/03/13/alpaca.html .	815
765			816
766	Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. <i>arXiv preprint arXiv:2304.03277</i> .	Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. 2023. Gemini: a family of highly capable multimodal models. <i>arXiv preprint arXiv:2312.11805</i> .	817
767			818
768			819
769	Iryna Pentina, Tyler Hancock, and Tianling Xie. 2023. Exploring relationship development with social chatbots: A mixed-method study of replika. <i>Computers in Human Behavior</i> , 140:107600.		820
770			821
771			822
772			823
773	Chen Qian, Xin Cong, Cheng Yang, Weize Chen, Yusheng Su, Juyuan Xu, Zhiyuan Liu, and Maosong Sun. 2023. Communicative agents for software development. <i>arXiv preprint arXiv:2307.07924</i> .	Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. 2023. Fine-tuning language models for factuality. In <i>The Twelfth International Conference on Learning Representations</i> .	824
774			825
775			826
776			827
777	Qwen. 2024. Introducing qwen1.5. https://qwenlm.github.io/blog/qwen1.5/ .	Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>arXiv preprint arXiv:2307.09288</i> .	828
778			829
779	Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, I. Evtimov, Joanna Bitton, Manish P Bhatt, Cristian Cantón Ferrer, Aaron Grattafori, Wenhan Xiong, Alexandre D’efossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. <i>Code llama: Open foundation models for code</i> . <i>ArXiv</i> , abs/2308.12950.	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	830
780			831
781			832
782			833
783			834
784			835
785			836
786			837
787			838
788			839
789	Rachel Rudinger, Aaron Steven White, and Benjamin Van Durme. 2018. Neural models of factuality. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 731–744.	Cunxiang Wang, Xiaoze Liu, Yuanhao Yue, Xiangru Tang, Tianhang Zhang, Cheng Jiayang, Yunzhi Yao, Wenyang Gao, Xuming Hu, Zehan Qi, et al. 2023a. Survey on factuality in large language models: Knowledge, retrieval and domain-specificity. <i>arXiv preprint arXiv:2310.07521</i> .	840
790			841
791			842
792			843
793			844
794			845
			846
			847
			848
			849

850 Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Man-
851 dlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and
852 Anima Anandkumar. 2023b. Voyager: An open-
853 ended embodied agent with large language models.
854 In *Intrinsically-Motivated and Open-Ended Learning*
855 *Workshop@ NeurIPS2023*.

856 Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa
857 Liu, Noah A Smith, Daniel Khashabi, and Hannaneh
858 Hajishirzi. 2023c. Self-instruct: Aligning language
859 models with self-generated instructions. In *Proceed-*
860 *ings of the 61st Annual Meeting of the Association*
861 *for Computational Linguistics*, pages 13484–13508.

862 Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley.
863 2023. *Baize: An open-source chat model with*
864 *parameter-efficient tuning on self-chat data*. In *Pro-*
865 *ceedings of the 2023 Conference on Empirical Meth-*
866 *ods in Natural Language Processing*, pages 6268–
867 6278, Singapore. Association for Computational Lin-
868 guistics.

869 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali
870 Farhadi, and Yejin Choi. 2019. Hellaswag: Can a
871 machine really finish your sentence? In *Proceedings*
872 *of the 57th Annual Meeting of the Association for*
873 *Computational Linguistics*, pages 4791–4800.

874 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
875 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
876 Zhuohan Li, Dacheng Li, Eric Xing, et al. 2024.
877 Judging llm-as-a-judge with mt-bench and chatbot
878 arena. *Advances in Neural Information Processing*
879 *Systems*, 36.

880 Qi Zhou, Bin Li, Lei Han, and Min Jou. 2023. Talking
881 to a bot or a wall? how chatbots vs. human agents
882 affect anticipated communication quality. *Computers*
883 *in Human Behavior*, 143:107674.

884 Dean H Zimmerman and Candace West. 1996. 9. sex
885 roles, interruptions and silences in conversation. In
886 *Towards a Critical Sociolinguistics*, page 211. John
887 Benjamins.

A Transition Sentences 888

To generate a sentence with coherent context, we
utilize ChatGPT to rewrite the template below,
which replaces {sentence_a} and {sentence_b}
with one transition sentence and new content re-
spectively. 889
890
891
892
893

Fuse the two sentences smoothly and
replace [topic] with the topic of sentence
two.

Sentence one "{sentence_a}"

Sentence two "{sentence_b}"

Give me your answer only, no other
words. Give me your answer only, no other
words.

A.1 Generation Termination Transition Sentences 894

1. <Empty string> 895
896
2. I need to cut you off right now; this is urgent. 897
3. Excuse me, I need to interject for a moment. 898
4. Sorry to interrupt, but I have something im- 899
portant to add. 900
5. Excuse me, may I interrupt for a moment? 901
6. I'm sorry to break in, but there's something 902
important I need to address. 903
7. I apologize for interrupting, but I'd like to 904
interject for a moment. 905
8. I'm sorry to interrupt, but I have a quick point 906
to make. 907
9. I appreciate your input, but I need a moment 908
of silence now. 909
10. I'm sorry to interrupt, but I really need some 910
quiet time to focus. 911
11. Enough talking! I need you to be quiet now. 912

A.2 Regeneration Transition Sentences 913

1. I may not have expressed myself clearly. What 914
I meant was [topic] 915
2. I think there might be a bit of confusion. Let 916
me clarify [topic] 917
918

919	3. I appreciate your input, but I was hoping for more details on [topic]	5. I was just reading about [topic]. What are your thoughts on that?	961
920			962
921	4. I think there might be a misunderstanding. What I'm really looking for is [topic]	6. By the way, speaking of something else.	963
922			
923	5. I may not have explained myself clearly. Let me rephrase the question. What are your thoughts on [topic]?	7. That reminds me, have you heard about [topic]?	964
924			965
925			
926	6. Actually, the correct information is [topic]. Could you share your perspective on that?	8. Can we shift gears for a moment and talk about [topic]?	966
927			967
928	7. I'm a bit confused because what you mentioned contradicts the information I have. Can we go over this again?	9. I've been curious about [topic]. Have you ever considered it?	968
929			969
930		10. I was thinking about [topic]. What are your thoughts on that?	970
931	8. I'm sorry, but that information seems to be incorrect. Let me clarify the question, and please provide the accurate details regarding [topic].		971
932		11. Now, shifting gears to a different subject, have you ever explored [topic]?	972
933			973
934			
935	9. I'm sorry, but that's not accurate. The correct information is [topic]. It's essential to have the correct details for our discussion.	12. Moving on to a different topic, have you ever considered [topic]?	974
936			975
937		13. Changing the subject, have you ever thought about [topic]?	976
938	10. I appreciate your effort in responding, but I think there might be a misunderstanding. What I intended to convey was [topic]. Let's revisit the topic to ensure we're on the same page.		977
939		14. Switching gears, let's talk about [topic]	978
940			
941		15. On a different note, have you ever thought about [topic]?	979
942			980
943	11. I see there might be some confusion. Let me clarify my point further to ensure we're on the same page. What I meant was [topic]. Can we discuss this to make sure we have a mutual understanding?	16. Speaking of which, have you ever considered exploring [topic]?	981
944			982
945		17. Changing the subject, let's now delve into [topic]	983
946			984
947			
948	12. There seems to be a misunderstanding. I meant [topic]. Let's align our understanding.	18. Shifting gears a bit, let's talk about [topic]	985
949			
950	13. No.		
951	14. Oh, No.		
952	15. No, you are wrong.		
953	A.3 Dialogue Reset Transition Sentences	B Questionnaire Details	986
954	1. <Empty string>	B.1 Subject Instruction	987
955	2. That's interesting, and speaking of [topic], have you ever...?	Before the experiment, we inform each participant of the subject instruction. The whole instruction is listed below:	988
956			989
957	3. I was just thinking about [topic], what are your thoughts on that?		990
958		1. This experiment requires subjects to have conversations with chat models. The content does not involve any dangerous remarks or have an impact on the subjects' physical and mental health.	991
959	4. That's fascinating! On a different note, have you ever thought about [topic]?		992
960		2. This test includes two parts: chatting and interacting with the models and filling out the questionnaire.	993
			994
			995
			996
			997
			998

999	3. The models are voice input and output modes that support multiple rounds of dialogue. At the end of each dialogue, you can press the new conversation button to start a new round of conversation.	collect their knowledge and usage of LLMs and voice assistants, which is tightly related to our research topic. As for the evaluation of the two chat models, we utilize their experience. The participants permit all those characteristics and experience information collection for research purposes only.	1043
1000			1044
1001			1045
1002			1046
1003			1047
1004	4. The models are English models and only support English dialogue.		1048
1005			1049
1006	5. There are two types of models, A and B. You must have at least 10 conversations with each model.		1050
1007			1051
1008			1052
1009	6. We have included some questions to start the conversation, just for reference.		1053
1010			1054
1011	7. This test mainly evaluates the performance of the two models in terms of response speed, human-likeness, faithfulness, factuality, and overall experience.		
1012			
1013			
1014			
1015	8. After the chat, fill out the questionnaire.		

B.2 Questionnaire

1016			
1017	1. Score the model’s response speed to evaluate whether the model can respond to your request.		
1018			
1019			
1020	2. Score the faithfulness of the model’s answers to evaluate whether the model understands your question, follows your instructions, and whether the answer is relevant to your chat topic.		
1021			
1022			
1023			
1024			
1025	3. Score the factuality of the model’s answers and evaluate whether the content of the answers is correct.		
1026			
1027			
1028	4. Score the human-likeness of the model’s answers and evaluate whether the conversation process between you and the model is close to the feeling of daily communication between people and whether the conversation process is smooth.		
1029			
1030			
1031			
1032			
1033			
1034	5. Score the overall experience of the model.		

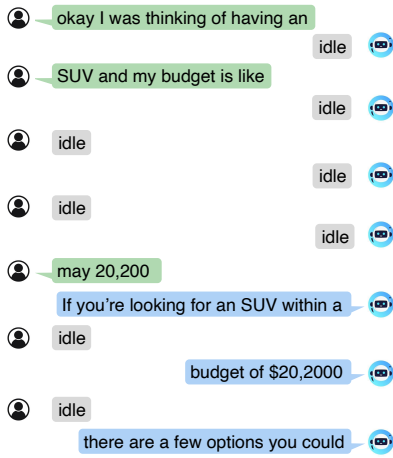
C Explanation of Ethical Concerns

All participants are recruited from a partner company. Those experiments are conducted during their working hours and we do not pay them additionally.

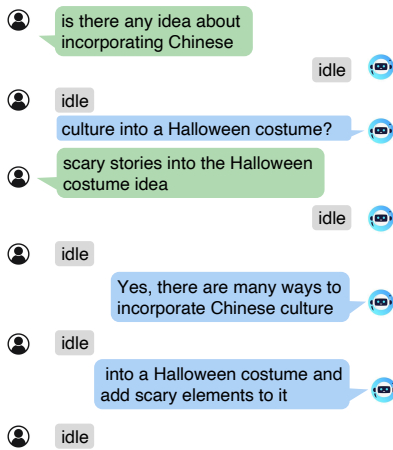
In the human-evaluation experiment, we collect basic demographic characteristics information: gender, age, and educational qualification. We also

D Case Demonstration

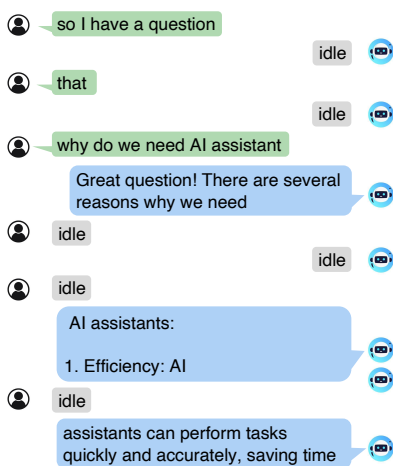
Here are some cases of conversation segments between the MiniCPM-duplex and human users. In Figure 7, the duplex model generates a response until it obtains enough information from the user.



(a) Case A



(b) Case B



(c) Case C

Figure 7: User study cases.