

# Can LLMs Generate High-Quality Task-Specific Conversations?

Anonymous ACL submission

## Abstract

This paper introduces a parameterization framework for controlling conversation quality in large language models. We explore nine key parameters across six dimensions that enable precise specification of dialogue properties. Through experiments with state-of-the-art LLMs, we demonstrate that parameter-based control produces statistically significant differences in generated conversation properties. Our approach addresses challenges in conversation generation, including topic coherence, knowledge progression, character consistency, and control granularity. The framework provides a standardized method for conversation quality control with applications in education, therapy, customer service, and entertainment. Future work will focus on implementing additional parameters through architectural modifications and developing benchmark datasets for evaluation.

## 1 Introduction

Generative AI systems can now produce coherent text, images, and code at scale (OpenAI et al., 2024b; Rombach et al., 2022; Chen et al., 2021). However, for multi-turn dialogue, model behavior remains difficult to specify and verify. Even when outputs appear plausible, properties such as topic coherence, knowledge progression, and persona stability can vary substantially across turns. This gap motivates a controllability interface that is both expressive enough to describe desired dialogue properties and measurable enough to evaluate whether those properties were realized.

In this paper, we study whether current LLMs admit a parameterized control interface for multi-turn conversation generation. Our objective is generation-facing: can we expose a set of dialogue attributes as explicit parameters, and does varying these parameters induce recoverable, systematic changes in the generated transcripts? Con-

cretely, we treat end-to-end conversation generation as a controlled simulation setting, where the transcript itself must manifest the configured properties, rather than relying on post-hoc prompting heuristics. (Zhang et al., 2020; Roller et al., 2021)

In domains such as business advising, civic services, and educational support, researchers often require long, context-rich dialogues that reflect specific user backgrounds, knowledge levels, and interaction styles. Yet existing simulator designs typically express these properties as natural-language instructions, leaving the intended attributes underspecified and difficult to measure. As a result, it is often unclear whether observed gains come from true controllability, model quality, or prompt artifacts. (Li et al., 2023; Huang et al., 2020; Bender et al., 2021)

We address this by introducing a parameter schema paired with a recoverability-based evaluation protocol. The schema makes conversation properties explicit and systematically variable; the protocol tests whether those properties are actually realized by asking independent judges to recover the parameters from transcripts. This framing directly operationalizes parameter adherence as the central validity criterion for a controllable generation interface, and separates controllability from incidental improvements in generic text quality. (Zheng et al., 2023)

The key contributions of this paper are:

1. A taxonomy of 35 conversation parameters with 9 dominating factors organized into six dimensions that capture essential aspects of multi-turn conversation control.
2. Analysis of parameter impact and redundancy under our recoverability-based evaluation protocol, identifying a core subset of high-impact parameters and documenting where additional parameters provide diminishing or overlapping gains.

- 082 3. A recoverability-based evaluation protocol  
083 that operationalizes parameter adherence as a  
084 transcript-level criterion for controllable gener-  
085 ation.
- 086 4. A benchmark study showing that modern  
087 LLMs can implement a subset of these pa-  
088 rameters through prompt conditioning, with  
089 measurable and statistically significant behav-  
090 ioral differences when parameters are varied.
- 091 5. An external validity design that tests whether  
092 recovered parameters transfer across dialogue  
093 regimes beyond simulator-only corpora.

094 In this paper, we introduce a parameterization  
095 framework for LLM-based conversation generation.  
096 Unlike unstructured prompting approaches, this  
097 parameterization enables precise specification of  
098 conversation properties that can be systematically  
099 varied, measured, and optimized. This approach  
100 builds upon prior work in controlled text genera-  
101 tion (Keskar et al., 2019; Dathathri et al., 2019;  
102 Khalifa et al., 2021) but extends these techniques  
103 specifically for multi-turn dialogue contexts with  
104 novel parameter dimensions.

## 105 2 Related Work

106 A large body of work studies controllable genera-  
107 tion by conditioning on attributes or control vari-  
108 ables, including early approaches to conditional  
109 language modeling and plug-in control mecha-  
110 nisms (Keskar et al., 2019; Dathathri et al., 2019;  
111 Khalifa et al., 2021). More recent work on steerabil-  
112 ity and preference-conditioned generation shows  
113 that LLM behavior can be moved along inter-  
114 pretable axes and optimized for diverse user prefer-  
115 ences, but is often evaluated at the utterance level or  
116 with short contexts (Liang et al., 2024; Wang et al.,  
117 2024). In contrast, our target is *transcript-level*  
118 *control*: parameters such as smoothness, focus,  
119 knowledge gap, and decision style must remain  
120 stable (or evolve in a prescribed way) across many  
121 turns, and success is defined by whether the config-  
122 ured attributes are recoverable from the resulting  
123 multi-turn conversation.

124 LLM-driven simulators are increasingly used to  
125 generate task-specific dialogues when real conver-  
126 sations are scarce, especially in low-resource set-  
127 tings. Recent work demonstrates practical frame-  
128 works for synthetic dialogue generation, but also  
129 highlights a central risk for prompt-based control:

apparent improvements may reflect template arti-  
facts rather than portable conversational proper-  
ties (Suresh et al., 2025; Castillo-Bolado et al.,  
2024). This concern motivates our external validity  
design, where we evaluate whether parameters re-  
covered from real corpora can be used to condition  
continuation generation in directionally consistent  
ways across different dialogue regimes.

Conversation quality is known to be multidimen-  
sional, and prior work documents persistent gaps  
between automatic metrics and human judgments,  
especially as context length grows (Deriu et al.,  
2020; Mehri and Eskenazi, 2020; See et al., 2019;  
Zheng et al., 2023). A growing set of benchmarks  
now explicitly targets multi-turn capabilities, inter-  
active agent behavior, and long-horizon instruction  
following, providing evidence that strong single-  
turn performance does not guarantee robust multi-  
turn behavior (Kwan et al., 2024; He et al., 2024;  
Sirdeshmukh et al., 2025; Jia et al., 2025; Castillo-  
Bolado et al., 2024). Our evaluation objective dif-  
fers from generic capability scoring: we measure  
whether a *specified control setting* is realized in  
the transcript, operationalized via parameter recov-  
erability and directionally consistent metric shifts  
under controlled perturbations.

Maintaining stable persona traits and role-  
specific behavior over extended interactions re-  
mains a known failure mode in dialogue gener-  
ation, including issues with entity tracking, coref-  
erence, and contradiction over multiple turns (Li  
et al., 2016; Zhang et al., 2018; Xu et al., 2022;  
Dziri et al., 2022). Recent benchmarks focusing on  
dynamic role-play and persona-conditioned inter-  
action further emphasize that long-horizon consis-  
tency is difficult even when local responses appear  
plausible (Yuan et al., 2025; Tan et al., 2025). Our  
parameter schema explicitly includes user identity  
and decision-making style controls, and our ad-  
herence metrics quantify where models succeed  
and where they fail under long contexts and cross-  
domain settings.

Finally, recent work proposes robustness-  
oriented multi-turn benchmarks and holistic evalu-  
ation in spoken dialogue regimes, emphasizing dis-  
tribution shift, error accumulation, and interaction-  
level brittleness (Xiang et al., 2025; Yan et al.,  
2025). These findings align with our motivation  
for external validity testing: if a control interface  
is meaningful, it should remain interpretable and  
recoverable beyond a single synthetic generation  
setup.

### 3 Methods

For this exploratory study, we selected nine parameters that has the highest impact from the 35 factors of conversation quality, which are spread across the six dimensions.

*Turn:* The number of turns of the conversation.

*Industry Context:* The initial field of this conversation.

*Knowledge Gap Level:* The prior knowledge the entrepreneur has of the conversation’s field. This is a method used in (Baskar et al., 2025) to measure the model’s knowledge alignment with the entrepreneur. We define the gap as a 1-5 integer value, where 1 refers to an expert with a deep understanding of the domain, and 5 refers to a complete novice with minimal business knowledge about their ideas.

*Smoothness Factor:* A grade A-F indicating conversation flow, with A referring to a perfectly flowing conversation with logical transitions, and F referring to a highly disjointed conversation with random topic jumping.

*Focus Level:* A grade 1-5 indicating how focused the entrepreneur is on this conversation. 1 refers to free-flowing, wide-ranging conversation covering many aspects, and 5 refers to laser-focused on specific details of implementation.

*Identity:* The initial setting of the entrepreneur’s background, which is used by (Aher et al., 2023) to simulate gender and racial diversity.

*Technical Language Level:* A 0-1 float number indicating the level of technical language the entrepreneur is using in the conversation. Similar methods were used in (Scarlatos et al., 2025) to trace knowledge levels in system-user conversation.

*Formality Level:* A 0-1 float number indicating the formal phrase usage in the conversation.

*Decision-Making Style:* The style of response the entrepreneur treats the system’s response. It can be one of analytical, Intuitive, consultative, or impulsive.

The exact definition of other parameters used in the prompt, the precise definition of the value of each parameter, and examples can be found in the Appendix.

**Prompt Engineering** The data set is created by constructing parameterized prompts that combine three key components: a base conversation generation prompt specifying the business advisory scenario, detailed parameter definitions for each

dimension, and the specific parameter values for each conversation instance. For each experimental condition, we systematically vary the parameter values while maintaining consistent entrepreneur background profiles and industry contexts. The final prompt is fed to the target LLM to generate complete multi-turn conversations. The full prompt structure with an example implementation is in the Appendix.

**Model Selection** We evaluate four state-of-the-art LLMs: *Gemini-2.5-pro* (Comanici et al., 2025), *Claude-3.7-sonnet* (Bai et al., 2022), *o3, o4-mini* (OpenAI, 2025), with other smaller or open-source LLMs: *Deepseek-r1* (DeepSeek-AI et al., 2025), *gpt-4o-mini* (OpenAI et al., 2024a), *Llama3.1:70b* (Grattafiori et al., 2024). Unless otherwise noted: temperature = 0.7, top- $p$  = 0.95, max-tokens = 8192, no custom stop tokens. We generate by default 20 samples per random seed.

**Baseline** We use prompt-based simulation using Claude Model *claude-3.7-sonnet* (Bai et al., 2022) as our baseline since it has the best performance among all other vanilla LLMs. (see the Appendix for baseline model comparison). Baseline results are produced using only the target turn, a random initial character setting with a brief background and previous experience with no special prompts or parameters, and rely solely on the LLM’s ability to generate outputs.

**Evaluation Methods** We assess our simulator along five axes using the same generated corpora: (i) *Topic Diversity* (topic counts and entropy), (ii) *Parameter Adherence* (recovery of set parameters from transcripts), (iii) *Topic Drift* (embedding-based distance to initial topic), (iv) *Character Properties Stability* (persona consistency and contradiction/memory errors), and (v) *Entity Revisit Rate* (and usefulness). Operational definitions and prompts for each axis appear in Appendix F. Human protocols and agreement statistics are summarized in the main text and fully detailed in Appendix A. To mitigate circularity in LLM-as-judge evaluation, we calibrate the judge against human labels on a held-out subset and report human-judge correspondence, and replicate core recoverability results with a consensus on multiple judges from a different provider when the generator under evaluation matches the default judge family.

For each task, we create a set of simulators and control the parameters to generate task-specific con-

283 conversations from the same topic.

284 *Topic Diversity* We use a random seed to create  
285 800 entrepreneurs’ background data. The gener-  
286 ated parameters are then injected into the prompt  
287 and fed into each LLM. We estimate topic diver-  
288 sity by clustering conversation topic summaries  
289 using sentence embeddings with cosine distance  
290 and a min cluster size set per dataset to maintain  
291 comparable granularity. For transparency, we ad-  
292 ditionally report sensitivity across a small range  
293 of clustering thresholds in Appendix F. We also  
294 compared the diversity of the topics by entropy:  
295  $H(X) = -\sum_{i=1}^n p(x_i) \log p(x_i)$ , where  $p(x_i)$  is  
296 the probability of topic  $x_i$ .

297 *Parameter Adherence* We generate 200 en-  
298 trepreneurs’ background data and randomized con-  
299 versation parameters. These are fed into each LLM  
300 across four different conversation lengths: 5, 10,  
301 15, and 20 turns, resulting in a total of 800 con-  
302 versations. The evaluation employs a hybrid human-  
303 LLM assessment framework in which both human  
304 annotators and *Claude-sonnet-3.7* serve as judges.

305 The evaluation protocol provides judges with  
306 only the conversation transcript, requiring them  
307 to infer the original parameters based on prede-  
308 fined parameter definitions. For numerical pa-  
309 rameters (all on the 1–5 Likert scale), adherence  
310 is measured by mean squared error:  $MSE =$   
311  $\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2$ , where  $y_i$  is the set value and  
312  $\hat{y}_i$  is the inferred value.

313 *Topic Drift* We generate 200 20-turn en-  
314 trepreneur conversations, each with smoothness  
315 factor set to A (highest topic adherence) and F  
316 (lowest topic adherence), along with 200 baseline  
317 conversations without smoothness factor control,  
318 resulting in 600 total conversations for topic drift  
319 analysis. The smoothness factor parameter con-  
320 trols the degree to which conversations maintain  
321 thematic coherence versus allowing natural topic  
322 exploration and deviation from the original busi-  
323 ness concept.

324 We compute topic drift as cosine distance be-  
325 tween the embedding of the entrepreneur utter-  
326 ance at turn  $t$  and the embedding of the opening  
327 topic cue:  $Drift(t) = 1 - \cos(\mathbf{e}_t, \mathbf{e}_0)$ , where  $\mathbf{e}_t =$   
328  $Embed(utterance_t)$  and  $\mathbf{e}_0 = Embed(topic\ cue)$ .  
329 Unless otherwise noted, embeddings use all-  
330 MiniLM-L6-v2 with mean pooling and  $\ell_2$  normal-  
331 ization.

332 *Character Properties Stability* We generate 500  
333 20-turn conversations with both the entrepreneur’s  
334 formality and technical levels randomized between

335 0 and 1, then 500 more with the formality param-  
336 eter omitted and another 500 with the technical  
337 parameter omitted.

338 Character stability is evaluated by:

339 *Formality Level:* Formality is determined by a  
340 composite of vocabulary sophistication, sentence  
341 structure, and pronoun usage.

342 *Technical Language Level:* The technical level  
343 is determined by the density of the domain ter-  
344 minology, the complexity of the concepts, and  
345 the usage of jargon. The final stability score  
346 is calculated by the  $1 - 0.5(\text{Formality Error} +$   
347  $\text{Technical Level Error})$

348 *Entity Revisit Rate* We generate 100 en-  
349 trepreneurs’ background information with Knowl-  
350 edge Gap Level parameters ranging from 1-5,  
351 where this parameter measures the knowledge dis-  
352 parity between the user’s existing background and  
353 their proposed business concept. Each entrepreneur  
354 profile is used to generate conversations in four dif-  
355 ferent lengths (5, 10, 15, and 20 turns).

356 The evaluation is done by first extracting NER  
357 and core concepts using BERT. We then track when  
358 previously mentioned entities reappear in subse-  
359 quent turns in the conversation. The concept of a  
360 recall rate is calculated as  $\frac{1}{T-1} \sum_{t=2}^T |\text{Entities}_t \cap$   
361  $\bigcup_{i=1}^{t-1} \text{Entities}_i|$ , where  $\text{Entities}_t$  represents the set  
362 of entities mentioned at turn  $t$ , and  $T$  is the total  
363 duration of the conversation.

364 **External Datasets for Validity** A concern in  
365 prompt-based controllable generation is that ap-  
366 parent control may reflect prompt artifacts in a  
367 synthetic setting rather than a portable interface  
368 that corresponds to properties observed in real  
369 dialogues. To test the validity of our parameter  
370 schema beyond simulator-only corpora, we evalu-  
371 ate our framework on three publicly accessible,  
372 multi-turn dialogue datasets that represent distinct  
373 dialogue regimes. We use the Fin-Ally (Das  
374 et al., 2025), MultiWOZ 2.2 (Zang et al., 2020),  
375 and STARv2 (Zhao et al., 2023) as our datasets.  
376 We use these datasets to answer two validity ques-  
377 tions. First, the parameter recoverability on real  
378 data. Given only a human-written transcript and  
379 parameter setting, can independent judges produce  
380 the similar dialogue that maintains the same con-  
381 versation and user parameters? This tests whether  
382 our parameters correspond to observable proper-  
383 ties of real dialogues. Second, the cross-domain  
384 portability. When we condition generation on a real  
385 dialogue prefix from each dataset and vary a single

Model	Embedding diversity
claude	0.2912
deepseek-r1	0.4161
o3	0.3360
o4-mini	0.2830
gpt-4.1	0.3436
gpt-4o-mini	0.2085
gemini	0.3747
llama3.1:70b	0.0576
baseline	0.1075

Table 1: Embedding diversity (sentence embedding).

Model	Topic diversity	Topic entropy
claude	111	4.469
deepseek-r1	143	5.275
o3	136	4.464
o4-mini	154	5.311
gpt-4.1	140	4.578
gpt-4o-mini	84	3.859
gemini	141	5.266
llama3	5	0.888
baseline	35	2.985

Table 2: Topic diversity and topic entropy.

parameter while holding all other settings fixed, do we observe systematic and directionally consistent changes in transcript-level measurements across domains? These tests are designed to separate true controllability from dataset-specific prompting effects, and to ground our simulator evaluation in external corpora with different conversational structures.

**Human Evaluation Protocols** We reuse the exact conversation from our automatic evaluations. Human studies cover parameter adherence, smoothness factor, knowledge-gap sensitivity, decision-making style classification, concept-revisit rate. Each item is rated by 2 blinded raters, with stratified sampling by model, length, and parameter level. We report  $MSE/\tau$  for numeric parameters and macro-F1/ $\kappa$  for categorical labels. For external datasets, annotators label a subset of the dataset with summaries and key parameters to test transfer and recoverability. Full task definitions, rater rubrics, sampling details, and statistical procedures appear in Appendix A.

## 4 Experiments

**Parameter adherence varies across models with improving accuracy over extended conversations.** A central concern with prompt-based controllable generation is that apparent control may be an artifact of the synthetic prompt distribution rather than a portable interface that corresponds

Parameter	Fin-Ally	STARv2	MultiWOZ
Formality (ICC)	0.8119	0.7632	0.7407
Technicality (ICC)	0.7876	0.7495	0.7211
Focus (ICC)	–	0.7084	0.6834
Smoothness (ICC)	–	0.6764	0.6572
Decision Style ( $\kappa$ )	0.8320	–	–
Knowledge Gap (ICC)	0.7795	–	–

Table 3: Human agreement on external corpora (two blinded raters). Fin-Ally is short in the released CSV, so we omit long-horizon labels.

Dataset	Setting	Foc. MSE↓	Smth. MSE↓	Drift↓
STARv2	Unctrl.	0.9433	0.8954	0.4141
STARv2	Param.	<b>0.7587</b>	<b>0.7015</b>	<b>0.3276</b>
MultiWOZ	Unctrl.	0.9117	0.8696	0.3909
MultiWOZ	Param.	<b>0.7340</b>	<b>0.6818</b>	<b>0.3026</b>

Table 4: External validity via prefix-to-continuation generation. “Param.” conditions generation on parameters recovered from the real prefix. Drift uses the same embedding-based topic drift metric as Sec. 3.

to properties present in real dialogue. We therefore complement our simulator-only adherence curves with two external evaluations that mirror our main protocol while isolating different failure modes. Whether the parameters are recoverable from human-written transcripts under our definitions, and whether conditioning on those recovered parameters produces directional and measurable changes when continuing real multi-turn dialogues.

Table 3 reports inter-rater agreement on a labeled subset of three public corpora using the same parameter rubrics as Sec. 3. Agreement is consistently high for style-linked dimensions with strong surface realizations, including Formality and Technicality. Across both business-oriented dialogue (Fin-Ally) and task-oriented dialogue (STARv2/MultiWOZ). Discourse-level properties that require longer-range judgments, such as Focus and Smoothness, achieve moderate but stable agreement on STARv2 and MultiWOZ. This can reflect the expected annotator ambiguity when a transcript mixes local coherence with occasional topic exploration. Fin-Ally is short in the released dataset, so we do not report long-horizon labels there; instead, we focus on parameters that can be inferred reliably from local behavior, where Decision-Making Style reaches strong agreement and Knowledge Gap remains recoverable.

Table 4 then evaluates whether this recoverability can be translated into portable control in a realistic deployment setting. We use a prefix-to-continuation protocol that given a human-written

447 dialogue prefix from STARv2 or MultiWOZ, we  
 448 generate a continuation either without control (Unc-  
 449 ctrl.) or conditioned on the parameters recovered  
 450 from that prefix (Param.). If parameterization is  
 451 valid beyond simulator corpora, it should reduce  
 452 adherence error and stabilize topic tracking even  
 453 when the context is out-of-distribution relative to  
 454 our synthetic prompts. This is exactly what we ob-  
 455 serve. The result shows the parameterized approach  
 456 effectively reduce the error rate in both datasets.  
 457 The consistency of these gains across two distinct  
 458 dialogue regimes supports the interpretation that  
 459 our interface captures transferable dialogue proper-  
 460 ties rather than prompt-specific artifacts.

461 **Simulators have bias on topic selection, and may**  
 462 **not generate a diverse pool of topics.** The sim-  
 463 ulators can be classified into two broad camps ac-  
 464 cording to their approach to exploring the subject  
 465 matter, as shown in Table 7. Advanced models such  
 466 as *Gemini-2.5-pro* and *DeepSeek-R1* exhibit supe-  
 467 rior topic diversification capabilities. These models  
 468 demonstrate a more uniform attention distribution  
 469 across thematic domains, closely approximating  
 470 human-like conversational breadth. In contrast, less  
 471 capable models like *GPT-4o-mini* produce more  
 472 constrained topic distributions, while lightweight  
 473 models such as *Llama3.1:70b* show severe limita-  
 474 tions with only 5 distinct topics.

475 The baseline approach without parameterization  
 476 yields poor diversity metrics. Mid-tier systems  
 477 occupy an intermediate position, with respectable  
 478 topic coverage but exhibiting concentration pat-  
 479 terns around familiar conceptual clusters. This  
 480 shows model architectures can explore diverse the-  
 481 matic spaces while maintaining coherent conversa-  
 482 tional flow.

483 We also examine sentence diversity by calculat-  
 484 ing semantic diversity through the cosine similar-  
 485 ity of embeddings generated by all-MiniLM-L6-  
 486 v2. (Table 1). The embedding diversity rankings  
 487 partially diverge from topic-level diversity mea-  
 488 sures, suggesting that models may employ different  
 489 strategies for achieving variation, and they may use  
 490 similar words or add additional definitions (e.g.,  
 491 AI-driven business vs. non-AI-driven) to express  
 492 different topics.

493 Beyond quantitative diversity measures, we ob-  
 494 serve systematic biases in topic selection patterns.  
 495 For example, when generating food-related busi-  
 496 ness scenarios, models frequently default to vegan  
 497 or health-conscious options regardless of user spec-

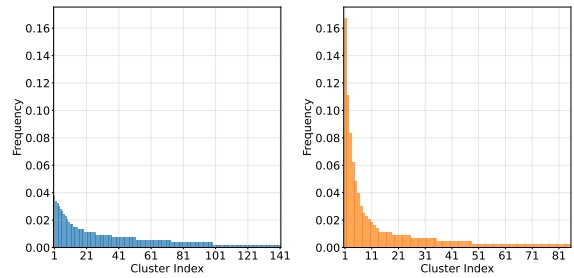


Figure 1: Topic frequency distributions of the *gpt-4o-mini* (orange) and *gemini-2.5-pro* (blue). Clusters are sorted in descending frequency. More advanced model (*gemini-2.5-pro*) produced a more diverse topic compared to the less advanced model (*gpt-4o-mini*)

498 ifications. This tendency toward "safe" or socially  
 499 desirable recommendations indicates inherent train-  
 500 ing biases that may limit the authenticity of gener-  
 501 ated conversations.

502 **Adding Smoothness factor improves topic cor-**  
 503 **relation.** Adding a smoothness factor to simulate  
 504 the conversation flow not only creates a diversifi-  
 505 ed conversation but also improves the model’s  
 506 adherence to the main topic. (Figure 2). Both  
 507 the small and the more advanced models can im-  
 508 prove adherence to the main topic after setting a  
 509 high smoothness factor, and advanced models can  
 510 successfully create a more significant difference  
 511 between high and low smoothness factors. Without  
 512 the smoothness factor, the model can only provide  
 513 a conversation that has low correlation to the given  
 514 topic. Human-labeled results show high agree-  
 515 ment between human annotators and LLM judges  
 516 ( $\kappa = 0.91$ ).

517 **Parameter adherence varies across models with**  
 518 **improving accuracy over extended conversa-**  
 519 **tions.** Analysis of parameter adherence across  
 520 conversation turns reveals substantial differences  
 521 in model capabilities, with most parameters show-  
 522 ing improved accuracy as conversations progress.  
 523 As shown in Figures 3(a-c), advanced models such  
 524 as Claude and Gemini demonstrate superior param-  
 525 eter implementation, with MSE errors for the focus  
 526 level, the knowledge gap level, and the experience  
 527 level decreasing from initial values to more accu-  
 528 rate parameter representation over 20 conversations  
 529 compared to the ground truth. This improvement  
 530 pattern suggests that models require several turns  
 531 to fully establish and maintain specified param-  
 532 eter values. The evaluation of the decision-making  
 533 style (Figure 3 (b)) shows binary classification ac-

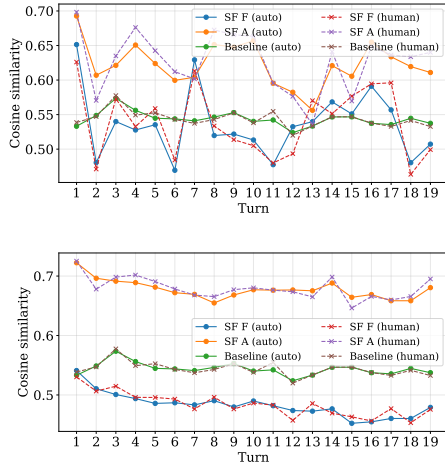


Figure 2: The cosine similarity between the entrepreneur’s utterance to the main topic in different smoothness factors for two models: (a) *gpt-4o-mini*, (b) *claude-3.7-sonnet*. *claude-3.7-sonnet* is showing a high separation between the highest and lowest smoothness factor, showing better understanding and adherence to parameters.

534 accuracy, where advanced models achieve 0.8-1.0 ac-  
 535 curacy rates while lighter models like *gpt-4o-mini*  
 536 struggle to maintain consistent classification per-  
 537 formance. The smoothness factor analysis (Figure  
 538 3(d)) demonstrates that parameter control effec-  
 539 tiveness varies by model architecture, with Claude  
 540 maintaining clear parameter differentiation while  
 541 smaller models show less distinct parameter im-  
 542 plementation regardless of specified values.

543 **Knowledge gap parameters influence concept re-**  
 544 **visit patterns in advanced models.** The relation-  
 545 ship between Knowledge Gap Level and concept  
 546 revisit behavior reveals substantial differences in  
 547 advanced models’ adaptation capabilities, as shown  
 548 in Figure 8. *Gemini-2.5-pro* exhibits a clear inverse  
 549 correlation between knowledge gap and revisit rate,  
 550 with highly knowledgeable users (Level 1) show-  
 551 ing revisit rates of approximately 0.5-0.6, while  
 552 novice users (Level 5) demonstrate lower revisit  
 553 rates around 0.1-0.2 across all conversation lengths.  
 554 This pattern aligns with pedagogical theory, where  
 555 experts benefit from reinforcement of complex con-  
 556 cepts, while beginners require more linear infor-  
 557 mation introduction. Conversely, Claude shows a  
 558 lower differentiation between knowledge gap lev-  
 559 els, but a higher differentiation over turns. This  
 560 shows that some models cannot correctly simulate  
 561 a conversation with low revisit rates.

562 With a high knowledge gap level, all models

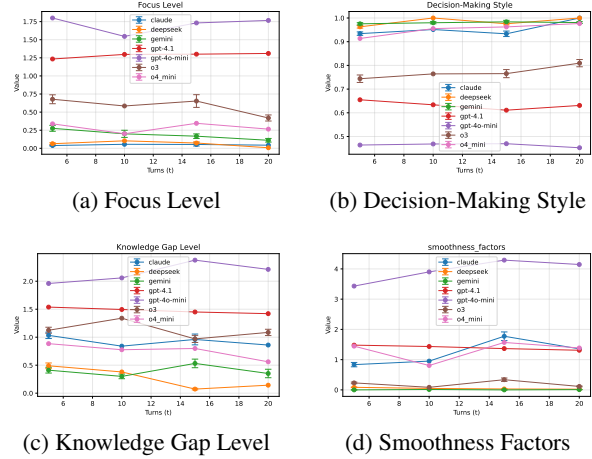


Figure 3: Model metric curves vs. conversation turns. The comparison between human evaluation can be found in App. E

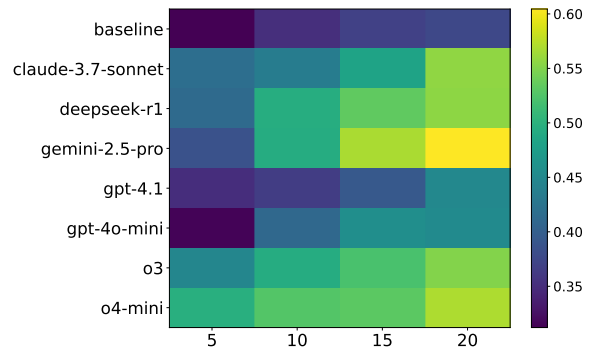


Figure 4: Concept-revisit rate by turns for each model with knowledge gap level of user set to 1 (most knowledgeable). All models exhibit a higher revisit rate with turn progression.

563 show a higher revisit rate compared to the base-  
 564 line. (Figure 6). Advanced models, including  
 565 *o3*, *gpt-4.1*, and *Claude-3.7-sonnet*, maintain high  
 566 character consistency scores that improve over ex-  
 567 tended conversations, while mid-tier models show  
 568 respectable but more variable performance. The  
 569 baseline approach demonstrates significantly lower  
 570 consistency. This suggests that sophisticated pa-  
 571 rameter implementation requires substantial model  
 572 capacity to fully understand and adhere to the pa-  
 573 rameters, but all models can obtain a significant  
 574 level of performance increase.

575 **Character parameters are stable across all mod-**  
 576 **els.** The character parameter study shows that  
 577 all models can reach high parameter stability over  
 578 turns, although more advanced models have bet-  
 579 ter performance (Figure 6). All models exhibit  
 580 improved stability trajectories over conversation

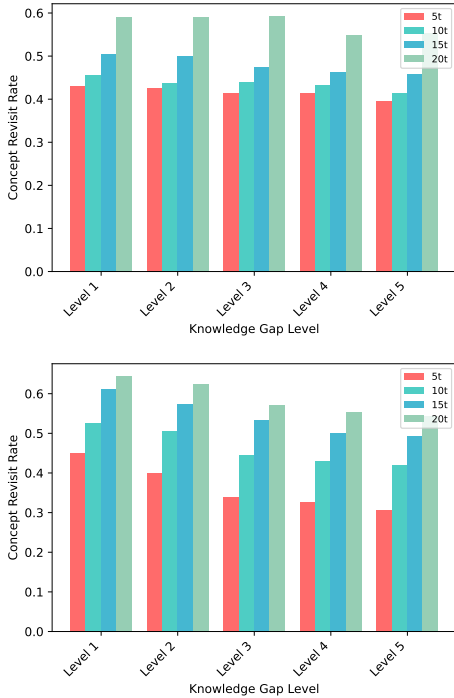


Figure 5: Concept-revisit rate by knowledge-gap level for two models: (a) Claude, (b) Gemini. Knowledge Gap Level 1 is the smallest knowledge gap, and Knowledge Gap Level 5 is the highest. *Gemini-2.5-pro* shows a more significant difference when modifying Knowledge Gap Level.

length, with consistency scores rising from initial values. This could be because the model does not have enough context initially, but the performance stabilizes after 15 turns.

We also performed an ablation analysis presented in Table 5, where we test the error of the model when only the formality parameters of the model or technical parameters are given. The result shows that the combined parameter implementation yields benefits exceeding the sum of individual components in both models. This suggests that adding more specified parameters to the model may further increase the model’s capability of simulating complex conversations.

## 5 Conclusion

We create a comprehensive parameterization framework for controlling LLM-based conversation generation, demonstrating both the potential and limitations of current approaches to fine-grained dialogue control. Our experiments with the simulator show that advanced models can effectively differentiate between parameter values and maintain improving consistency over long conversations.

Model	Turns	Formality	Technical	Full
claude-3.7	5	0.280	0.252	0.206
claude-3.7	10	0.305	0.265	0.205
claude-3.7	15	0.298	0.258	0.192
claude-3.7	20	0.292	0.252	0.184
o3	5	0.255	0.212	0.173
o3	10	0.222	0.175	0.143
o3	15	0.215	0.162	0.131
o3	20	0.212	0.155	0.130

Table 5: Average performance errors for Formality Only, Technical Only, and Full Parameters across varying conversation turns. Human evaluation shows a correlation of Pearson’s  $r = 0.938$ . See App. E for full human results and agreement analyses

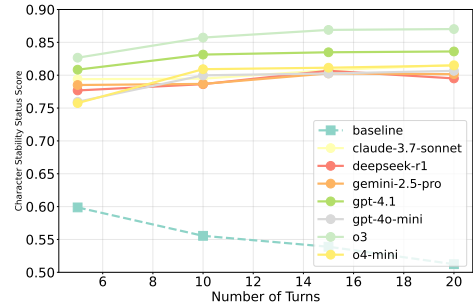


Figure 6: Character Parameter stability over turns, the baseline has a decreasing stability over turns, while all other models with character properties show an increase in stability score. The agreement between the human and LLM judge at baseline is  $\kappa = 0.9316$ , and  $\kappa = 0.9241$  for *claude-3.7-sonnet*

## 6 Limitations

Several issues are unaddressed in this exploratory study. We only provide the necessary parameters for conversation generation, not an exhaustive set of parameters that covers all aspects. More parameters could be added to the prompt since we have already proven that interconnected parameters can improve conversation quality.

A fine-tuned LLM with human-labeled conversation parameters as a dataset may increase the simulator’s sensitivity to intermediate values. We are using the default temperature settings. More analysis could be made on different parameter settings and fine-tuned open-source LLMs.

Parameterized settings cannot increase the model’s factual accuracy. Adding a factual accuracy parameter can prompt the LLM to provide incorrect information, but they are also not sensitive enough to intermediate parameters and does not decrease the hallucination rate compared to the vanilla model. A RAG-based approach is still needed to decrease the simulator’s hallucination.

## 7 Ethical Considerations

We adhere to the ACL Code of Ethics and will submit the Responsible NLP checklist.

Our contribution is a *simulation* framework for controllable, multi-turn conversations. Risks include persuasive misuse, misrepresentation of synthetic dialogues as human, and style emulation of specific groups. We label released samples as *synthetic*, recommend visible and machine-readable provenance, prohibit deceptive use, and discourage targeting real individuals.

All API models are used under their providers' standard research terms; Llama-3.1-70B is used under its open-source license. No user PII is sent; prompts contain synthetic content only. We vary parameters that affect coherence/focus to study controllability, not to generate harmful content. Degradation is limited to stylistic/structural properties in benign domains, with disallowed-content checks. We stress that parameterization does *not* ensure truthfulness; downstream use should add retrieval/verification.

The *Identity* parameter can reinforce stereotypes if entangled with behavior. We decouple identity from behavioral controls, randomize combinations, and encourage reporting group-wise summaries and disparity checks. The framework must not be used to infer protected attributes.

Corpora are model-generated; no PII is collected or released. Any artifacts (prompts, code, metrics) will be filtered for PII/copyright and released under a research license forbidding deceptive use.

Annotators worked only with synthetic text, provided informed consent, could withdraw, and were fairly compensated; no demographics were collected.

We provide prompts/parameter grids, seeds, and evaluation code for replication and report computing to limit environmental impact.

Simulators aid pre-deployment testing but do not replace human studies; final evaluations require real users under appropriate ethical review and consent.

## 8 LLM usage

We use LLM mainly for paraphrasing sentences and checking grammar mistakes.

## References

- Gati V Aher, Rosa I Arriaga, and Adam Tauman Kalai. 2023. Using large language models to simulate multiple humans and replicate human subject studies. pages 337–371.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, and 1 others. 2022. [Constitutional ai: Harmlessness from ai feedback](#).
- Sarvesh Baskar, Tanmay Tulsidas Verelakar, Srinivasan Parthasarathy, and Manas Gaur. 2025. [From guessing to asking: An approach to resolving the persona knowledge gap in llms during multi-turn conversations](#).
- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? pages 610–623.
- David Castillo-Bolado, Joseph Davidson, Finlay Gray, and Marek Rosa. 2024. [Beyond prompts: Dynamic conversational benchmarking of large language models](#). *Preprint*, arXiv:2409.20222.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, and 1 others. 2021. [Evaluating large language models trained on code](#).
- Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blistein, Ori Ram, and 1 others. 2025. [Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities](#).
- Sarmistha Das, Priya Mathur, Ishani Sharma, Sriparna Saha, Kitsuchart Pasupa, and Alka Maurya. 2025. [Fin-ally: Pioneering the development of an advanced, commonsense-embedded conversational ai for money matters](#). *Preprint*, arXiv:2509.24342.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. 2019. Plug and play language models: A simple approach to controlled text generation.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, and 1 others. 2025. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#).
- Jan Deriu, Alvaro Rodrigo, Arantxa Otegi, Guillermo Echeгойen, Sophie Rosset, Eneko Agirre, and Mark Cieliebak. 2020. Survey on evaluation methods for dialogue systems. *Artificial Intelligence Review*, 54(1):755–810.

725	Nouha Dziri, Ehsan Kamaloo, Sivan Milton, Osmar Zai-	OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher,	779
726	iane, Mo Yu, Edoardo Ponti, and Siva Reddy. 2022.	Adam Perelman, Aditya Ramesh, Aidan Clark,	780
727	On the origin of hallucinations in conversational mod-	AJ Ostrow, Akila Welihinda, Alan Hayes, Alec	781
728	els: Is it the datasets or the models? pages 5271–	Radford, Aleksander Mađry, Alex Baker-Whitcomb,	782
729	5285.	Alex Beutel, and 1 others. 2024a. <a href="#">Gpt-4o system</a>	783
		<a href="#">card</a> .	784
730	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	OpenAI. 2025. Introducing OpenAI o3 and o4-mini.	785
731	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	Accessed: 2025-07-29.	786
732	Dahle, Aiesha Letman, Akhil Mathur, and 1 others.		
733	2024. <a href="#">The llama 3 herd of models</a> .		
734	Yun He, Di Jin, Chaoqi Wang, Chloe Bi, Karishma	OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal,	787
735	Mandyam, Hejia Zhang, Chen Zhu, Ning Li, Tengyu	Lama Ahmad, and 1 others. 2024b. <a href="#">Gpt-4 technical</a>	788
736	Xu, Hongjiang Lv, Shruti Bhosale, Chenguang Zhu,	<a href="#">report</a> . <i>Preprint</i> , arXiv:2303.08774.	789
737	Karthik Abinav Sankararaman, Eryk Helenowski,		
738	Melanie Kambadur, Aditya Tayade, Hao Ma, Han	Stephen Roller, Emily Dinan, Naman Goyal, Da Ju,	790
739	Fang, and Sinong Wang. 2024. <a href="#">Multi-if: Benchmarking</a>	Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott,	791
740	<a href="#">llms on multi-turn and multilingual instructions</a>	Eric Michael Smith, Y-Lan Boureau, and Jason West-	792
741	<a href="#">following</a> . <i>Preprint</i> , arXiv:2410.15553.	ston. 2021. Recipes for building an open-domain	793
		chatbot. pages 300–325.	794
742	Minlie Huang, Xiaoyan Zhu, and Jianfeng Gao. 2020.	Robin Rombach, Andreas Blattmann, Dominik Lorenz,	795
743	Challenges in building intelligent open-domain dia-	Patrick Esser, and Bjorn Ommer. 2022. High-	796
744	log systems. 38(3):1–32.	resolution image synthesis with latent diffusion mod-	797
		els. pages 10684–10695.	798
745	Qi Jia, Ye Shen, Xiujie Song, Kaiwei Zhang, Shibo	Alexander Scarlatos, Ryan S Baker, and Andrew Lan.	799
746	Wang, Dun Pei, Xiangyang Zhu, and Guangtao Zhai.	2025. Exploring knowledge tracing in tutor-student	800
747	2025. <a href="#">One battle after another: Probing llms’ limits</a>	dialogues using llms. pages 249–259.	801
748	<a href="#">on multi-turn instruction following with a benchmark</a>		
749	<a href="#">evolving framework</a> . <i>Preprint</i> , arXiv:2511.03508.		
750	Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney,	Abigail See, Stephen Roller, Douwe Kiela, and Jason	802
751	Caiming Xiong, and Richard Socher. 2019. <a href="#">Ctrl: A</a>	Weston. 2019. <a href="#">What makes a good conversation?</a>	803
752	<a href="#">conditional transformer language model for control-</a>	<a href="#">how controllable attributes affect human judgments</a> .	804
753	<a href="#">lable generation</a> .		
754	Muhammad Khalifa, Hady Elsahar, and Marc Dymet-	Ved Sirdeshmukh, Kaustubh Deshpande, Johannes	805
755	man. 2021. A distributional approach to controlled	Mols, Lifeng Jin, Ed-Yeremai Cardona, Dean Lee,	806
756	text generation.	Jeremy Kritz, Willow Primack, Summer Yue, and	807
		Chen Xing. 2025. <a href="#">Multichallenge: A realistic multi-</a>	808
757	Wai-Chung Kwan, Xingshan Zeng, Yuxin Jiang, Yufei	<a href="#">turn conversation evaluation benchmark challenging</a>	809
758	Wang, Liangyou Li, Lifeng Shang, Xin Jiang, Qun	<a href="#">to frontier llms</a> . <i>Preprint</i> , arXiv:2501.17399.	810
759	Liu, and Kam-Fai Wong. 2024. <a href="#">Mt-eval: A multi-</a>		
760	<a href="#">turn capabilities evaluation benchmark for large lan-</a>	Sathya Krishnan Suresh, Wu Mengjun, Tushar Pranav,	811
761	<a href="#">guage models</a> . <i>Preprint</i> , arXiv:2401.16745.	and Eng Siong Chng. 2025. <a href="#">Diasynth: Synthetic</a>	812
		<a href="#">dialogue generation framework for low resource dia-</a>	813
762	Jiwei Li, Michel Galley, Chris Brockett, Georgios Sp-	<a href="#">logue applications</a> . <i>Preprint</i> , arXiv:2409.19020.	814
763	ithourakis, Jianfeng Gao, and Bill Dolan. 2016. A		
764	persona-based neural conversation model. pages 994–	Juntao Tan, Liangwei Yang, Zuxin Liu, Zhiwei Liu,	815
765	1003.	Rithesh Murthy, Tulika Manoj Awalgaoonkar, Jian-	816
		guo Zhang, Weiran Yao, Ming Zhu, Shirley Kokane,	817
766	Yunxiang Li, Zihan Li, Kai Zhang, Ruilong Dan, and	Silvio Savarese, Huan Wang, Caiming Xiong, and	818
767	You Zhang. 2023. Chatdoctor: A medical chat model	Shelby Heinecke. 2025. <a href="#">Personabench: Evaluat-</a>	819
768	fine-tuned on a large language model meta-ai (llama)	<a href="#">ing ai models on understanding personal informa-</a>	820
769	using medical domain knowledge. <i>arXiv preprint</i>	<a href="#">through accessing (synthetic) private user data</a> .	821
770	<i>arXiv:2303.14070</i> .	<i>Preprint</i> , arXiv:2502.20616.	822
771	Xun Liang, Hanyu Wang, Yezhaohui Wang, Shichao	Haoxiang Wang, Yong Lin, Wei Xiong, Rui Yang,	823
772	Song, Jiawei Yang, Simin Niu, Jie Hu, Dan Liu,	Shizhe Diao, Shuang Qiu, Han Zhao, and Tong	824
773	Shunyu Yao, Feiyu Xiong, and Zhiyu Li. 2024. <a href="#">Con-</a>	Zhang. 2024. <a href="#">Arithmetic control of LLMs for di-</a>	825
774	<a href="#">trollable text generation for large language models:</a>	<a href="#">verse user preferences: Directional preference align-</a>	826
775	<a href="#">A survey</a> . <i>Preprint</i> , arXiv:2408.12599.	<a href="#">ment with multi-objective rewards</a> . In <i>Proceedings</i>	827
		<a href="#">of the 62nd Annual Meeting of the Association for</a>	828
776	Shikib Mehri and Maxine Eskenazi. 2020. <a href="#">Usr: An</a>	<a href="#">Computational Linguistics (Volume 1: Long Papers)</a> ,	829
777	<a href="#">unsupervised and reference free evaluation metric for</a>	pages 8642–8655, Bangkok, Thailand. Association	830
778	<a href="#">dialog generation</a> .	for Computational Linguistics.	831

832	Hao Xiang, Tianyi Tang, Yang Su, Bowen Yu, An Yang, Fei Huang, Yichang Zhang, Yaojie Lu, Hongyu Lin, Xianpei Han, Jingren Zhou, Junyang Lin, and Le Sun. 2025. <a href="#">Rmtbench: Benchmarking llms through multi-turn user-centric role-playing</a> . <i>Preprint</i> , arXiv:2507.20352.	886
833		887
834		888
835		889
836		890
837		
838	Jing Xu, Arthur Szlam, and Jason Weston. 2022. Beyond goldfish memory: Long-term open-domain conversation. pages 5180–5197.	891
839		892
840		893
841	Ruiqi Yan, Xiquan Li, Wenxi Chen, Zhikang Niu, Chen Yang, Ziyang Ma, Kai Yu, and Xie Chen. 2025. <a href="#">Uro-bench: Towards comprehensive evaluation for end-to-end spoken dialogue models</a> . <i>Preprint</i> , arXiv:2502.17810.	894
842		895
843		896
844		897
845		
846	Dingbo Yuan, Yipeng Chen, Guodong Liu, Chenchen Li, Chengfu Tang, Dongxu Zhang, Zhenkui Wang, Xudong Wang, and Song Liu. 2025. Dmt-rolebench: A dynamic multi-turn dialogue based benchmark for role-playing evaluation of large language model and agent. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 39, pages 25760–25768.	898
847		899
848		900
849		901
850		902
851		903
852		904
853	Xiaoxue Zang, Abhinav Rastogi, Srinivas Sunkara, Raghav Gupta, Jianguo Zhang, and Jindong Chen. 2020. <a href="#">Multiwoz 2.2 : A dialogue dataset with additional annotation corrections and state tracking base-lines</a> . <i>Preprint</i> , arXiv:2007.12720.	905
854		
855		
856		
857		
858	Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. Personalizing dialogue agents: I have a dog, do you have pets too? pages 2204–2213.	907
859		908
860		
861		
862	Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020. <a href="#">Dialogpt: Large-scale generative pre-training for conversational response generation</a> .	909
863		910
864		911
865		912
866		913
867	Jeffrey Zhao, Yuan Cao, Raghav Gupta, Harrison Lee, Abhinav Rastogi, Mingqiu Wang, Hagen Soltau, Izhak Shafran, and Yonghui Wu. 2023. <a href="#">Anytod: A programmable task-oriented dialog system</a> . <i>Preprint</i> , arXiv:2212.09939.	914
868		915
869		916
870		917
871		918
872	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. <i>Advances in neural information processing systems</i> , 36:46595–46623.	919
873		920
874		921
875		922
876		923
877		
878	<b>A Human Evaluation Protocols</b>	
879	<b>A.1 Participants</b>	
880	We recruited five experienced university student volunteers with 3 different expertise areas located in the US (each >6 months of prior annotation work) for the human study. Annotators completed a brief screening task. No personally identifying information was collected. The study was reviewed	
881		
882		
883		
884		
885		
	as non-human-subjects research under institutional guidance because annotators evaluated synthetic text only and provided no private information; nevertheless, we followed standard consent and data-minimization practices.	886
		887
		888
		889
		890
	<b>Limitation.</b> With $N=10$ raters this pilot primarily establishes feasibility and inter-rater reliability bounds; it does not provide population-level estimates. We therefore report uncertainty and treat the human study as confirmatory for construct validity and judge calibration, not as a definitive benchmark.	891
		892
		893
		894
		895
		896
		897
	<b>A.2 Blinding</b>	898
	Annotators were strictly blind to (i) the generating model identity/version, (ii) all parameter settings, and (iii) the study hypotheses. Each task UI displayed only the initial topic cue (one-line description) and the conversation transcript in a structured JSON-like view. Conversation length was visible by necessity but randomized across items.	899
		900
		901
		902
		903
		904
		905
	<b>A.3 Tasks and Rater Rubrics</b>	906
	All tasks used written rubrics with anchors and examples.	907
		908
	<b>• Parameter Adherence.</b>	909
	– <i>Focus Level (1–5)</i> : We use the exact same definition as in the prompt, with examples for each score level.	910
		911
	– <i>Smoothness (A–F→1–6)</i> : Anchors include turn-to-turn linkage checks.	912
		913
	– <i>Knowledge Gap Level (1–5)</i> : inferred user expertise; 1 (expert) to 5 (novice). Anchors emphasize initiative, prerequisite vocabulary, and abstraction level.	914
		915
	– <i>Formality / Technicality (0.0–1.0 sliders)</i> : defined via lexicon sophistication, syntax, register; and domain-term density/jargon use. Anchored examples provided.	916
		917
		918
	– <i>Decision-Making Style (categorical)</i> : {analytical, intuitive, consultative, risk-averse, impulsive}. Definitions with exemplars; single-best-label instruction.	919
		920
		921
		922
		923
		924
		925
		926
	<b>• Character Consistency Audit.</b> Per-turn tagging of violations: <i>contradiction</i> (factual/persona), <i>memory error</i> , <i>style drift</i> . Raters mark presence/absence and leave a short note at first occurrence.	927
		928
		929
		930
		931
		932

933	• <b>KGL Sensitivity.</b> Blind inference of KGL	(iii) time-on-task filters (remove items below 10th	980
934	with a scaffolding checklist (evidence of pre-	percentile viewing time unless justification note	981
935	requisite knowledge, abstraction level, need	present), and (iv) duplication prevention.	982
936	for definitions, rate of requests for clarifica-		
937	tion).		
938	• <b>Entity Revisit Usefulness.</b> For a sampled	<b>A.7 Metrics and Statistical Analysis</b>	983
939	subset of turns, raters label entity revisits as	• <b>Agreement.</b> For numeric scales we report	984
940	{reinforce, advance, clarify, superficial}, with	Intraclass Correlation Coefficient (ICC(2,k))	985
941	examples.	and Krippendorff’s $\alpha$ ; for categoricals (De-	986
942		cision Style) we report Cohen’s $\kappa$ . All with	987
943	• <b>Controlled Degradation.</b> For items gener-	nonparametric 95% CIs via 1,000 bootstrap	988
944	ated under low-quality settings (e.g., Smooth-	replicates clustered by item.	989
945	ness=F), raters answer a binary check: “Does		
946	the transcript manifest the specified degrada-	• <b>Mixed-Effects Models.</b> We fit linear/logistic	990
	tion?” with short justification.	mixed-effects models with fixed effects for	991
		parameter level, model, and turn length, and	992
		random intercepts for prompt and rater. Model	993
		selection used AIC with maximal random-	994
		effects structure feasible for convergence.	995
947	For the external datasets, the annotators are given	• <b>Multiple Testing.</b> We control the false discov-	996
948	the original dialogue, and in addition to annotating	ery rate across families of related hypotheses	997
949	the original dataset based on the 9 main parameters,	via Benjamini–Hochberg ( $q=0.10$ ) and report	998
950	they are also prompted to provide a summarization	both raw and adjusted $p$ -values when shown.	999
951	of the dialogue as the original prompt to the model		
952	for future generation tasks.		
953	<b>A.4 Sampling and Power</b>	<b>A.8 LLM-Judge Calibration (“Judge Card”)</b>	1000
954	We stratified by <i>model</i> (all models under test), <i>turn</i>	We calibrate an LLM judge on a rater-labeled sub-	1001
955	<i>length</i> (5/10/15/20), and <i>parameter level</i> (extremes	set (70/30 train/holdout split, stratified by task).	1002
956	+ midpoints where applicable). Within each stratum,	For each metric we fit a bias-correction model (or-	1003
957	we uniformly sampled conversation instances	inary least squares for numeric, multinomial logistic	1004
958	from the automatic pool. Given $N=2$ raters, the	for categorical) mapping LLM scores to human	1005
959	study is underpowered for small effects; our <i>a priori</i>	scores. We report: (i) human–LLM correlation	1006
960	goal was precision of reliability (95% CI width	(Pearson/Spearman) on the holdout set, (ii) calibra-	1007
961	$\leq 0.20$ for $\kappa$ in high-agreement regimes) and coarse	tion plots, (iii) post-calibration RMSE / macro-F1,	1008
962	effect-size estimation. We therefore report boot-	and (iv) recommended use constraints (e.g., avoid	1009
963	strap CIs and avoid null-hypothesis claims when	mid-level Smoothness adjudication without human	1010
964	CIs are wide.	review). To prevent family overlap bias, we avoid	1011
965	<b>A.5 Rater Recruitment, Training, and</b>	using a judge from the same provider as any com-	1012
966	<b>Procedure</b>	pared generator in sensitivity analyses and replicate	1013
967	Annotators completed: (i) a 30-minute tutorial with	with an alternative judge to assess robustness.	1014
968	screenshots, (ii) a guided practice set of 12 items	<b>A.9 Data and Artifact Availability</b>	1015
969	with gold explanations, (iii) a qualification quiz	We release (under a research license): anonymized	1016
970	(pass threshold: $\geq 80\%$ agreement with gold ratio-	prompts for annotation, rater rubrics/instructions,	1017
971	nals). Items were presented in randomized order	gold examples, and the code used for sampling,	1018
972	with per-rater unique shuffles. Model and param-	UI rendering, agreement computation, bootstrap,	1019
973	eter metadata were removed from the payload and	and mixed-effects modeling. Generated transcripts	1020
974	filenames. Each session included regular, unlabeled	are released where license permits; otherwise we	1021
975	gold checks.	provide hashes and per-turn feature summaries.	1022
976	<b>A.6 Quality Control and Exclusions</b>	<b>B Selecting the High-Impact Parameters</b>	1023
977	We applied: (i) gold-check filters (remove rater-	We first define a neutral configuration $C_{\text{neutral}}$ over	1024
978	session segments with $<70\%$ gold agreement), (ii)	all 35 parameters. For each parameter $p$ , we specify	1025
979	attention checks (detect straight-lining on sliders),		

two extreme but valid values ( $p_{\text{low}}, p_{\text{high}}$ ) consistent with our parameter definitions, while keeping all remaining parameters fixed at their neutral values.

We then sample  $M$  entrepreneur background profiles (we use  $M=200$  in our planned experiment, matching the scale used for parameter adherence in Sec. 3). For each background  $b$  and each parameter  $p$ , we generate a matched pair of conversations:

$c_{p,b}^{\text{low}}$  : generated with  $c_{\text{neutral}}$  but with  $p \leftarrow p_{\text{low}}$ ,  
 $c_{p,b}^{\text{high}}$  : generated with  $c_{\text{neutral}}$  but with  $p \leftarrow p_{\text{high}}$ .

All runs use the same model, the same random-seed policy, and the same fixed conversation length to remove confounds. This yields, for each parameter  $p$ , a dataset

$$D(p) = \{(b, c_{p,b}^{\text{low}}, c_{p,b}^{\text{high}}) \mid b \in \{1, \dots, M\}\},$$

in which any systematic differences between  $c_{p,b}^{\text{low}}$  and  $c_{p,b}^{\text{high}}$  can be attributed to the single parameter  $p$ .

**Behavioral effect size.** Given the single-parameter datasets  $D(p)$ , we quantify how much each parameter changes conversation behavior using the same metrics as in Sec. 3. For a given parameter  $p$  and metric  $A$ , we compute the signed difference for each matched pair:

$$\Delta_{p,b}^A = A(c_{p,b}^{\text{high}}) - A(c_{p,b}^{\text{low}}).$$

We summarize its magnitude by the mean effect size:

$$E^A(p) = E_b[\Delta_{p,b}^A].$$

To compare parameters across metrics with different scales, we  $z$ -normalize effect sizes per metric:

$$\tilde{E}^A(p) = \frac{E^A(p) - \mu_A}{\sigma_A},$$

where  $\mu_A$  and  $\sigma_A$  are the mean and standard deviation of  $\{E^A(p)\}_{p=1}^{35}$  for metric  $A$ .

**Controllability via parameter recovery.** We also measure how controllable each parameter is, using the same parameter-recovery idea as in our main study. For each  $p$ , we pool all conversations in  $D(p)$  and ask *Claude-3.7-sonnet* to infer whether a given transcript was generated with  $p_{\text{low}}$  or  $p_{\text{high}}$  (without seeing the parameter values). This yields a classification accuracy  $\text{Acc}(p)$  for distinguishing

Parameter	Mean $ \tilde{E}(p) $	$C(p)$	$S(p)$
<b>Knowledge Gap Level</b>	1.393	0.9245	<b>1.288</b>
<b>Smoothness Factor</b>	1.329	0.9264	<b>1.231</b>
<b>Focus Level</b>	1.090	0.9082	<b>0.9971</b>
<b>Technical Language Level</b>	0.9557	0.9265	<b>0.8855</b>
<b>Formality Level</b>	0.8747	0.9457	<b>0.8272</b>
<b>Decision-Making Style</b>	0.8024	0.9097	<b>0.7319</b>
<b>Identity (User)</b>	0.8192	0.8459	<b>0.6929</b>
<b>Industry Context</b>	0.7171	0.8159	<b>0.5851</b>
<b>Complexity Progression</b>	0.5462	0.7530	0.4113
Practical-Theoretical Balance	0.5040	0.7271	0.3664
Framework	0.4786	0.7063	0.3380
Feedback Reception	0.4418	0.6752	0.2983
Emotional Journey	0.4266	0.6438	0.2746
Factual Accuracy	0.4036	0.5940	0.2397

Table 6: Single-parameter sensitivity ranking. We remove *Turns* since it is always included as a base setting in all experiments.

low vs. high settings. We convert this into a  $[0, 1]$  controllability score:

$$C(p) = \max(0, 2 \cdot (\text{Acc}(p) - 0.5)),$$

so that random guessing ( $\text{Acc}=0.5$ ) maps to  $C(p)=0$  and perfect identification maps to  $C(p)=1$ .

**Overall sensitivity index.** For each parameter, we define an overall sensitivity index that combines how much it changes the conversation with how reliably that change is recoverable from the transcript:

$$S(p) = C(p) \cdot \frac{1}{|\mathcal{A}(p)|} \sum_{A \in \mathcal{A}(p)} \tilde{E}^A(p),$$

where  $\mathcal{A}(p)$  is the subset of metrics theoretically relevant to parameter  $p$ .

**Intermediate values.** For intermediate parameter values, we agree with the reviewers' reading and explicitly report this limitation in Sec. ??: current models reliably separate extreme settings, but show weak discrimination among mid-range settings. Our goal is to empirically document this limitation, not to claim that the framework currently yields perfect fine-grained control.

## C Prompts

In this section, we present the prompt used for conversation generation.

## C.1 Raw Prompt

### Raw Prompt

Create a K-turn conversation between an AI adviser and an entrepreneur trying to work on <A business field>. In the conversation, the AI adviser is an informed business coach in a Small Business Development Corporation, and the entrepreneur is a < entrepreneur's demographic background > with a focus on <entrepreneur's idea>.

## C.2 Parameterized Prompt

Below is the complete prompt to the LLM for parameterized conversation generation:

### Parameterized Prompt

Create a K-turn conversation between an AI adviser and an entrepreneur trying to work on <A business field>. In the conversation, the AI adviser is an informed business coach in a Small Business Development Corporation, and the entrepreneur is a < entrepreneur's demographic background > with a focus on <entrepreneur's idea>.

Core structural parameters that define the conversation's basic framework:

- **Purpose:** The primary intent of the conversation
  - advisory: Problem-solving and guidance-focused dialogue
  - educational: Knowledge transfer and learning-oriented
  - exploratory: Discovery and brainstorming-centered
  - evaluative: Assessment and critique-focused
- **Turns:** Total number of conversation turns (exchanges between participants)
- **Turn Balance:** Distribution of conversation contributions between participants (expressed as ratio, e.g., "55:45" means user speaks 55% of turns, advisor 45%)
- **Arc:** Overall narrative structure of the conversation
  - problem-solution: Identifies issues and develops solutions
  - exploration-conclusion: Broad investigation leading to specific outcomes
  - question-answer: Sequential inquiry and response pattern
  - build-refine: Iterative development and improvement process

- **Initiator:** Which participant starts the conversation
  - user: Entrepreneur begins with question or problem
  - assistant: Advisor opens with inquiry or observation

- user: Entrepreneur begins with question or problem
- assistant: Advisor opens with inquiry or observation

- **Topic Scope:** Array of subject areas that may be covered during the conversation (e.g., ["food business", "marketing", "operations"])

Parameters defining the characteristics and relationship between conversation participants:

- **Knowledge Gap Level (KGL)**

- 1: Expert with deep understanding of business domain
- 2: Advanced practitioner with solid foundational knowledge and some specialized expertise
- 3: Moderate familiarity with business concepts
- 4: Basic understanding with significant knowledge gaps requiring guidance
- 5: Complete novice with minimal business knowledge about their ideas

- **Assistant Parameters:**

- **Identity:** Role and background description (e.g., "experienced business advisor with small business expertise")
- **Consistency Level:** How consistently the assistant maintains their role and expertise (0.0 = highly variable, 1.0 = perfectly consistent)

- **User Parameters:**

- **Identity:** Role and background description (e.g., "early-stage food business entrepreneur")
- **Focus Level (FL)**
  - \* 1: Free-flowing, wide-ranging conversation covering many aspects
  - \* 2: Mostly broad discussion with occasional deep dives into specific areas
  - \* 3: Balanced focus with some exploration of tangential topics
  - \* 4: Primarily focused on core issues with minimal tangential exploration
  - \* 5: Laser-focused on specific details of implementation
- **Prior Knowledge Level:** User's existing expertise in the domain (1 = complete novice, 2 = limited

knowledge, 3 = moderate level understanding, 4 = extensive previous experience, 5 = expert level)

– **Decision-Making Style (DMS)**

- \* **Analytical:** Focuses on data, metrics, and logical analysis
- \* **Intuitive:** Relies on gut feeling and personal judgment
- \* **Consultative:** Seeks multiple perspectives before deciding
- \* **Risk-averse:** Prioritizes minimizing potential downsides
- \* **Impulsive:** Makes quick decisions without extensive deliberation

– **Feedback Reception (FR)**

- \* **Receptive:** Eagerly accepts and builds upon advice
- \* **Balanced:** Considers advice thoughtfully with moderate acceptance
- \* **Skeptical:** Questions most suggestions, needs convincing
- \* **Resistant:** Pushes back against most advice, difficult to persuade

Parameters controlling how knowledge is delivered and educational objectives are achieved:

- **Framework:** Educational methodology employed
  - socratic: Question-driven discovery learning
  - didactic: Direct instruction and explanation
  - collaborative: Joint problem-solving approach
  - experiential: Learning through practical examples and scenarios
- **Practical-Theoretical Balance:** Ratio of practical application to theoretical concepts (0.0 = purely theoretical, 1.0 = purely practical)
- **Complexity Progression:** Array showing how conceptual difficulty increases throughout the conversation (e.g., [0.3, 0.5, 0.7, 0.8] indicates gradual complexity increase)
- **Industry Context:** Specific sector or domain focus (e.g., "food-business", "technology", "healthcare")

Parameters governing interpersonal interactions and emotional progression:

- **Formality:** Level of professional versus casual communication (0.0 = highly casual, 1.0 = highly formal)

- **Emotional Journey:** Array of emotional states and their intensities throughout the conversation

- Each entry contains an emotion and intensity level (0.0 = minimal, 1.0 = maximum)
- Example: ["uncertainty": 0.8, "curiosity": 0.7, "confusion": 0.5, "understanding": 0.6, "confidence": 0.7]

- **Relationship Development:** How much the participant relationship evolves during the conversation (0.0 = static relationship, 1.0 = significant relationship building)

- **Disagreement Handling:** Approach to managing conflicting viewpoints

- diplomatic: Respectful acknowledgment and gentle correction
- direct: Clear, straightforward disagreement
- avoidant: Minimizing or redirecting conflict
- collaborative: Working together to resolve differences

Parameters controlling language use and communication style:

- **Technical Language Level:** Degree of specialized terminology and jargon (0.0 = plain language only, 1.0 = highly technical)
- **Question Types:** Distribution of different inquiry styles
  - **Closed:** Yes/no or specific factual questions
  - **Open:** Broad, exploratory questions requiring detailed responses
  - **Rhetorical:** Questions posed for emphasis rather than response
  - **Clarifying:** Questions seeking to understand or confirm information
  - Values should sum to 1.0 (e.g., "closed": 0.2, "open": 0.5, "rhetorical": 0.1, "clarifying": 0.2)
- **Response Style:** Communication characteristics
  - **Conciseness:** Brevity versus elaboration (0.0 = very verbose, 1.0 = extremely concise)
  - **Directness:** Straightforward versus indirect communication (0.0 = highly indirect, 1.0 = completely direct)
  - **Formality:** Professional versus casual language (0.0 = very casual, 1.0 = highly formal)

Parameters ensuring quality and comprehensiveness of conversation content:

- **Factual Accuracy:** Degree of correctness in information provided (0.0 = potentially inaccurate, 1.0 = verified accuracy)
- **Example Specificity:** Level of detail in illustrations and case studies (0.0 = general examples, 1.0 = highly specific, detailed examples)
- **Stakeholder Perspectives:** Array of viewpoints to be considered during the conversation (e.g., ["customer", "supplier", "regulator", "competitor"])

When generating conversations using these parameters:

1. Begin by establishing participant identities and knowledge levels
2. Follow the specified conversation arc while maintaining turn balance
3. Progress complexity according to the defined progression array
4. Incorporate emotional journey elements at appropriate conversation points
5. Ensure content addresses multiple stakeholder perspectives
6. Maintain consistency with linguistic pattern specifications
7. Adapt technical language level to participant knowledge asymmetry

Before conversation generation, validate that:

- All numerical parameters fall within specified ranges (0.0-1.0)
- Question type distributions sum to 1.0
- Turn balance ratios are mathematically consistent
- Complexity progression shows logical advancement
- Stakeholder perspectives are relevant to industry context

Generated conversations should follow this structure:

```
{
  "metadata": {
    "participantRoles": {...},
    "conversationArc": "...",
    "totalTurns": n
  },
  "conversation": [
```

```
{
  "turn": 1,
  "speaker": "user|assistant",
  "content": "...",
  "emotionalState": "...",
  "complexityLevel": 0.x
},
...
],
"analysis": {
  "parameterAdherence": {...},
  "learningObjectivesMet": [...],
  "stakeholderPerspectivesCovered": [...]
}
}
```

Here is an example input about a user's background:

```
{
  "conversationParameters": {
    "fundamentals": {
      "purpose": "advisory",
      "turns": 12,
      "turnBalance": "55:45",
      "arc": "problem-solution",
      "initiator": "user",
      "topicScope":
        ["food business",
         "marketing", "operations"]
    },
    "participants": {
      "knowledgeGapLevel": 3,
      "assistant": {
        "identity":
          "experienced business advisor",
        "consistencyLevel": 0.85
      },
      "user": {
        "identity":
          "early-stage food"
          "business entrepreneur",
        "focusLevel": 3,
        "priorKnowledgeLevel": 0.4,
        "decisionMakingStyle": "analytical",
        "feedbackReception": "receptive"
      }
    },
    "learningApproach": {
      "framework": "socratic",
      "practicalTheoreticalBalance": 0.7,
      "complexityProgression":
        [0.3, 0.5, 0.7, 0.8],
      "industryContext": "food-business"
    },
    "conversationDynamics": {
      "formality": 0.7,
      "emotionalJourney": [
        {"uncertainty": 0.8},
        {"curiosity": 0.7},
        {"understanding": 0.6},
        {"confidence": 0.7}
      ],
      "relationshipDevelopment": 0.5,
      "disagreementHandling": "diplomatic"
    },
    "linguisticPatterns": {
      "technicalLanguageLevel": 0.6,
      "questionTypes": {
        "closed": 0.2,
```

Model	Topic diversity	Topic entropy
claude	25	2.366
deepseek-r1	18	2.195
o3	27	2.493
o4-mini	33	2.880
gpt-4.1	31	2.762
gpt-4o-mini	12	1.012
gemini-2.5-pro	28	2.511
llama3.1:70b	5	0.810
claude-3.7-sonnet	35	2.985

Table 7: Topic diversity and topic entropy of baseline models.

```

    "open": 0.5,
    "rhetorical": 0.1,
    "clarifying": 0.2
  },
  "responseStyle": {
    "conciseness": 0.5,
    "directness": 0.6,
    "formality": 0.7
  }
},
"contentAttributes": {
  "factualAccuracy": 0.9,
  "exampleSpecificity": 0.6,
  "stakeholderPerspectives":
  ["customer", "supplier",
  "regulator", "competitor"]
}
}
}

```

### C.3 Conversation Parameters Structure

The conversation generator operates using a hierarchical parameter system organized into six main categories: Fundamentals, Participants, Learning Approach, Conversation Dynamics, Linguistic Patterns, and Content Attributes.

## D More Results

### D.1 Baseline Performance Comparison

We compare the performance of different models in terms of topic diversity and topic entropy when given the baseline prompt. (Table 7). The result shows *claude-3.7-sonnet* has the best topic diversity, and smaller models like *llama3.1:70b* have the same poor performance compared to the parameterized version.

### D.2 More Parameter Adherence Results

**Experience Level** We categorize the experience level using the prior knowledge level in the original prompt and calculate the MSE between the actual and predicted value. All models show a decrease in MSE with higher turns. (Figure 8)

### D.3 Human-Evaluated Results for Table 5

We asked human raters to score the per-condition average errors corresponding to Table 5 (*Formality Only, Technical Only, Full*) across turn lengths. Errors were judged using the same rubric (Sec. A) and aggregated to match the table layout. Below we report the human-evaluated averages and alignment with the automatic metrics.

Model	Turns	Formality	Technical	Full
claude-3.7	5	0.283	0.255	0.208
claude-3.7	10	0.302	0.268	0.207
claude-3.7	15	0.300	0.260	0.193
claude-3.7	20	0.294	0.254	0.186
o3	5	0.257	0.214	0.174
o3	10	0.223	0.176	0.144
o3	15	0.216	0.163	0.132
o3	20	0.213	0.156	0.131

Table 8: Human-evaluated average performance errors corresponding to Table 5.

**Correspondence with automatic metrics.** We measure correspondence between Table 5 and Table 8 across all cells.

- Pearson  $r = 0.938$  (95% CI: 0.922–0.952)
- Spearman  $\rho = 0.934$  (95% CI: 0.916–0.948)
- ICC(2,k) = 0.931 (95% CI: 0.907–0.948)

These values indicate high alignment ( $> 0.93$ ) between human ratings and the automatic metrics used in Table 5.

**Feedback Reception** The measurement of feedback reception is categorized into four types described in the prompt, and the result is calculated based on the rate of correct classification. The response indicates that some advanced models achieve a very high level of accuracy by combining a mixture of LLM and human decision-making, demonstrating that these models can accurately simulate the user’s sentiment based on a description. Other advanced models and small models show less optimal results in this role-playing setting. (Figure 8)

**While simulators can generate good responses, they may fail to create bad ones** While models demonstrate clear differentiation between extreme parameter values in focus levels (Level 1 vs Level 5), they exhibit poor sensitivity to intermediate parameter settings. In Figure 7, all three models show relatively flat performance curves across the

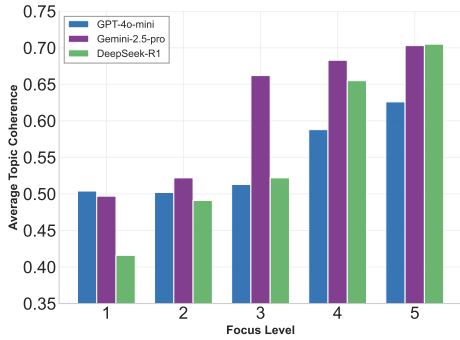


Figure 7: Average topic coherence between turns, models have parameter sensitivity issues on intermediate values, but all models can differentiate the lowest and highest value.

middle range (Levels 2-4), with topic coherence scores clustering around 0.45-0.55 regardless of the specified focus level. This suggests that models can successfully implement "very focused" versus "very unfocused" conversation styles but struggle to generate nuanced variations in between.

Similar behavior is observed in Figures 2 and 6. In Figure 2, both models show only marginally lower cosine similarity scores compared to the baseline, failing to achieve the expected degradation specified by smoothness factor  $F$  (*Highly disjointed with random topic jumping*). In Figure 6, *claude-3.7-sonnet* demonstrates minimal differentiation between knowledge gap levels 1 and 5, while *gemini-2.5-pro* exhibits comparable limitations, conflating performance across levels 3-5 despite maintaining clear separation between the extreme values (levels 1 and 5).

The insensitivity of the parameter may be due to a lack of fine-tuning. With only the definitions for each level provided to the LLM, models can only rely on pre-trained representations to map abstract parameter descriptions to concrete output behaviors. Given sufficient examples of intermediate quality levels between "highly focused" and "completely unfocused" conversations, the model could possibly provide a more distinguishable result. Further, post-training alignment procedures through RLHF further reinforce the model's tendency to produce helpful, coherent responses, creating systematic resistance to generating lower-quality content regardless of parameter specifications, which lowers the model's ability to generate poor-quality conversations.

## E Supplemental Human Evaluation

### E.1 A/B Topic Coherence vs. Baseline (Win-Rate)

Model	Win-rate vs. Baseline
Gemini-2.5-pro	0.8911
Claude-3.7-sonnet	0.8458
Deepseek-r1	0.8465
o3	0.8347
o4-mini	0.8410
gpt-4o-mini	0.7026

Table 9: Human A/B coherence win-rates.

### E.2 Parameter Adherence (Human vs. Automatic)

Parameter	Pearson $r$	Human $\leftrightarrow$ Auto agreement
Focus	0.914	ICC=0.901
Smoothness	0.944	ICC=0.932
KGL	0.931	ICC=0.919
Formality	0.953	ICC=0.941
Technicality	0.946	ICC=0.934
Decision Style	-	$\kappa = 0.918$ , macro-F1=0.882

Table 10: Human-automatic correspondence by parameter.

### E.3 KGL Sensitivity (Human Recovery Accuracy)

Model	Exact match	$\pm 1$ tolerance
Gemini-2.5-pro	0.6314	0.9108
Claude-3.7-sonnet	0.6025	0.8969
Deepseek-r1	0.6122	0.8703
o3	0.5867	0.8668
o4-mini	0.6160	0.8826
gpt-4o-mini	0.4701	0.7995
Baseline	0.3846	0.7321

Table 11: Human inference of KGL from transcripts.

## F Evaluation Tasks

Here, we first introduce our evaluation tasks and explain the methods in Section 3.

**Topic Diversity** The conversation needs a topic to start. After setting the topic area before the simulation, LLM will pick a subtopic based on the configured parameters to best suit the entrepreneur's background. In this task, we compare the distributions of topics mentioned by the simulator.

**Parameter Adherence** To evaluate whether the conversation generated follows the given parameters, we evaluate the difference between the settled

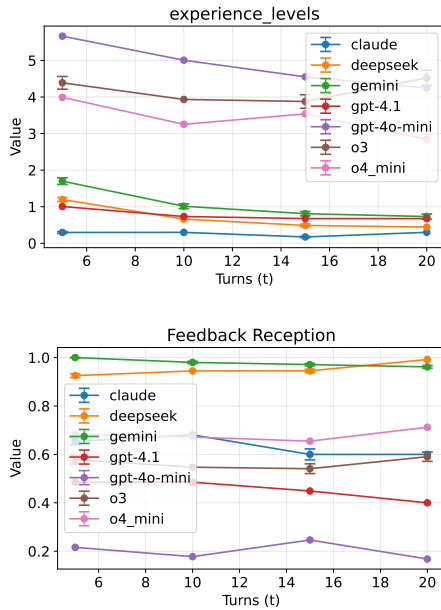


Figure 8: Additional figures on parameter adherence

parameters vs. the inferred parameters given only the generated conversation.

**Topic Drift** Natural dialogue often involves gradual topic transitions that can lead to substantial drift from the original subject matter, making thematic coherence throughout extended conversations a challenge. We measure the semantic distance between conversation segments to quantify how far the dialogue deviates from its initial topic focus. We calculate sentence embedding to compute cosine similarity scores between the opening conversational topic and subsequent dialogues, tracking the drift over turns.

**Character Properties Stability** Consistent character portrayal across conversation turns is essential for believable dialogues, yet current LLMs often exhibit personality inconsistencies that undermine conversation quality. This evaluation measures character stability by analyzing linguistic markers, decision-making patterns, and domain expertise demonstrations throughout generated conversations. We measure deviations between the character’s behavior in conversation versus their given background or parameters.

**Entity Revisit Rate** Effective conversations demonstrate sophisticated information management by strategically reintroducing previously mentioned entities, concepts, and topics, creating coherent narrative threads rather than generating un-

related information. We quantify how frequently and effectively the conversation simulator references earlier elements by tracking named entities and key concepts from earlier turns, then analyzing whether their subsequent appearances serve meaningful conversational purposes.

1245  
1246  
1247  
1248  
1249  
1250