

LikeBench: Evaluating Subjective Likability in LLMs for Personalization

Anonymous ACL submission

Abstract

A personalized LLM should remember user facts, apply them correctly, and adapt over time to provide responses that the user prefers. Existing LLM personalization benchmarks are largely centered on two axes: accurately recalling user information and accurately applying remembered information in downstream tasks. We argue that a third axis, likability, is both subjective and central to user experience, yet under-measured by current benchmarks. To measure likability holistically, we introduce LikeBench¹, a multi-session, dynamic evaluation framework that measures likability across multiple dimensions by how much an LLM can adapt over time to a user’s preferences to provide more likable responses. In LikeBench, the LLMs engage in conversation with a simulated user and learn preferences only from the ongoing dialogue. As the interaction unfolds, models try to adapt to responses, and after each turn, they are evaluated for likability across seven dimensions by the same simulated user. To the best of our knowledge, we are the first to decompose likability into multiple diagnostic metrics, which makes it easier to pinpoint where a model falls short. To improve realism and discriminativeness, LikeBench uses fine-grained, psychologically grounded personas instead of coarse high/low trait ratings used in prior work. Our benchmark shows that strong memory performance does not guarantee high likability: DeepSeek R1, with lower memory accuracy (86%, 17 facts/profile), outperformed Qwen3 by 28% on likability score despite Qwen3’s higher memory accuracy (93%, 43 facts/profile). Even SOTA models like GPT-5 adapt well in short exchanges but show only limited robustness in longer, noisier interactions.

1 Introduction

As large language models (LLMs) become increasingly integrated into everyday life, the need for

systems that genuinely adapt to individual users is more important than ever (Zhang et al., 2025, 2018). Personalized LLMs are widely recognized as the next step toward building AI that feels truly realistic and human-like, since “one-size-fits-all” alignment cannot capture the subjective preferences, values, and conversational styles of each user (Kirk et al., 2024; Guan et al., 2025; Salemi et al., 2023). This evolution places personalization and, crucially, the “likability” of AI responses—at the heart of next-generation systems, making it essential to develop rigorous evaluation methods that accurately capture and advance these user-centric capabilities.

However, despite the growing importance of personalization, existing benchmarks (Jiang et al., 2025; Wu et al., 2025; Maharana et al., 2024; Zhao et al., 2025) for LLMs remain largely focused on technical capabilities such as retaining user information (memory recall), for example, remembering that a user has a dairy allergy, and applying remembered preferences in tasks (memory adherence), such as recommending suitable foods while adhering to food allergy memory. But these benchmarks cannot capture whether interactions genuinely feel personalized or likable to individual users. This gap presents a fundamental limitation: a model might achieve very good memory recall and adherence scores, yet still come across as generic or unsatisfying if it fails to adapt to a user’s personality, conversational style, and nuanced preferences. Moreover, since most LLMs are post-trained with RLHF algorithms, they tend to moderately satisfy many users rather than deeply satisfy any particular user, due to the distribution of their reward models.

There are very few works that attempt to measure this likability factor. However, they have several shortcomings: the user profiles or personas used to simulate users often lack fine-grained personality traits and conversational styles (Wu et al., 2024), which are critical for user’s likability. This results in LLMs being evaluated mainly on less dis-

¹Code & Data is planned for future release

084 criminative users, even though real users have di- 132
085 verse interests and personalities. Importantly, these 133
086 benchmarks typically reduce likability to a single 134
087 aggregate score (such as an alignment score) (Wu 135
088 et al., 2024), making it difficult to understand where 136
089 an LLM struggles or to provide actionable feedback 137
090 for further improvement. 138

091 To address these limitations, we introduce 139
092 LikeBench, a comprehensive benchmark designed 140
093 to overcome three key issues in existing personal- 141
094 ization benchmarks. First, unlike most prior bench- 142
095 marks that focus mainly on memory recall and 143
096 adherence, LikeBench targets the subjective dimen- 144
097 sion of user experience by evaluating likability, the 145
098 extent to which AI interactions actually feel satisfy- 146
099 ing and well-adapted to users. Second, rather than 147
100 relying on a single aggregated metric for measuring 148
101 likability, LikeBench provides a multi-dimensional 149
102 assessment, using multi-session, multi-turn conver- 150
103 sations spanning varied topics. Third, our bench- 151
104 mark features psychologically grounded user per- 152
105 sonas, modeled with fine-grained descriptive per- 153
106 sonality traits and conversation styles, moving be- 154
107 yond the coarse rating-based personas used in prior 155
108 work. Our contributions can be summarized as: 156

- 109 • We introduce LikeBench, a multi-session dy- 157
110 namic evaluation framework that measures 158
111 two key properties of LLMs: likability—ability 159
112 to generate responses that users find likable, 160
113 and adaptability—the ability to improve likabil- 161
114 ity over time through ongoing conversations 162
115 with simulated users. 163
- 116 • We enable fine-grained evaluation of likability 164
117 by decomposing it into seven diagnostic met- 165
118 rics, providing a comprehensive assessment. 166
- 119 • We develop psychologically grounded user 167
120 personas based on seven personality traits 168
121 with 35 facets, along with conversation style 169
122 across nine dimensions, offering substantially 170
123 richer and more distinctive profiles. 171
- 124 • We experiment with a range of SOTA models 172
125 and find that memory alone does not guarantee 173
126 likability, and most models show little to no 174
127 adaptability as conversations progress. 175

128 2 Related Work 176

129 2.1 Memory Benchmark 177

130 Recent benchmarks on personalization emphasize 178
131 on memory recall and adherence to long contexts 179

and multi-session settings. LoCoMo (Maharana 132
et al., 2024) evaluates LLMs in very similar set- 133
tings, where models must generate answers to QA 134
tasks based on hundreds of conversational turns; 135
performance is measured by comparing the gener- 136
ated responses to annotated ground-truth answers 137
using exact-match and F1 metrics, with no LLM-as- 138
judge involved. LongMemEval (Wu et al., 2025) 139
expands this direction by introducing a range of 140
tasks—extraction, temporal reasoning, knowledge 141
updates, and uses GPT-4o as an automatic judge: 142
given the model’s generated answer and the gold 143
label, the LLM judge determines correctness, thus 144
combining LLM-based scoring. PrefEval (Zhao 145
et al., 2025) focuses on preference adherence where 146
model generations are assessed by an LLM-as- 147
judge using several binary criteria, and a discrimi- 148
native classification protocol, where the model se- 149
lects the user-consistent answer from pre-generated 150
options. PersonaMem (Jiang et al., 2025) extends 151
evaluation to tracking evolving user preferences 152
over sessions, reporting performance in both a dis- 153
criminative setting (selecting the correct response 154
from multiple candidates) and a generative setting, 155
where the chosen response is the candidate with 156
the highest generation probability. Most recently, 157
HiCUPID (Mok et al., 2025), like prior bench- 158
marks, focuses on factual correctness and infor- 159
mation adherence in extended contexts, using an 160
LLM-as-judge to measure win rate by comparing 161
model predictions with ground-truth answers. In 162
summary, existing work evaluates whether models 163
remember and apply user information across turns, 164
sessions, and long contexts. These evaluations typ- 165
ically equate personalization with objective, fac- 166
tual memory or preference using a binary metric 167
(right or wrong) and compliance, but fail to assess 168
whether interactions are genuinely liked by individ- 169
ual users or whether the LLM is able to adapt to 170
user preferences and improve over time. 171

172 2.2 Benchmarks on Likability 172

173 Very recently, research on personalization has 173
174 shifted from factual recall and adherence to sub- 174
175 jective alignment, where existing benchmarks dif- 175
176 fer along four axes: whether evaluation is static 176
177 (pre-generated conversation with evaluating LLM 177
178 responds to last user query) or dynamic, the depth 178
179 of persona, whether likability is decomposed into 179
180 multiple metrics or reduced to a single aggregate 180
181 score, and whether adaptability over time is mea- 181
182 sured. ALOE (Wu et al., 2024) is dynamic but 182

Table 1: Comparison of likability-oriented benchmarks across key dimensions.

Benchmark	Conversation	# Likability Metrics	Adaptation Over Time	Persona Modeling	Persona Facets	Memory Performance	#Profiles	#Sessions	#Turns
ALOE (Wu et al., 2024)	Dynamic	1	Yes	Coarse keyword	No	None	100	1	10
ALIGNX (Chan et al., 2024)	Static	1	No	Intensity ratings	No	None	3,716	1	1
CUPID (Kim et al., 2025)	Static	1	No	Intensity ratings	No	Implicit only	252	9	6
LikeBench (ours)	Dynamic	7	Yes	Fine-grained text	Yes	Explicit + Implicit	50	10	5

single-session, models users with coarse keyword (e.g. “extroverted”), reports a single composite alignment score plus an improvement rate, and thus does not provide a decomposed view of likability. ALIGNX (Chan et al., 2024) models users with broader personality traits than existing approaches, but it is static, represents users’ personalities via coarse intensity ratings (e.g. “high extroversion), reports a single aggregate alignment score, and does not measure adaptation over time. CUPID (Kim et al., 2025) uses static evaluation and compares model generations against ground-truth answers with an LLM-as-judge; it employs coarse intensity ratings, reports a single alignment score rather than decomposed metrics, does not include an explicit measure of adaptability, and also relies on simulated users providing per-turn implicit feedback to the evaluating LLM, making interactions less realistic. In summary, existing benchmarks often evaluate static or single-session settings, rely on coarse ratings-based personas, and collapse alignment into a single aggregate score without tracking whether models adapt with continued interaction. In contrast, LikeBench evaluates dynamic, multi-session conversations, uses psychologically grounded, fine-grained textual personas, decomposes likability into multiple diagnostic metrics, and explicitly measures adaptability across sessions thus providing actionable insight into where systems fall short on user satisfaction. These comparisons are summarized in Table 1.

3 Benchmark: LikeBench

Problem Setup and Notation. Let Θ denote the set of models (LLMs) under evaluation. An LLM $f_\theta \in \Theta$ interacts with a simulated user g_π instantiated from a fine-grained textual persona $\pi \in \Pi$. Each persona π provides personality traits, con-

versation style, interests, background, and foundational knowledge. Evaluation proceeds over S sessions, each with a hidden session prior $\rho_s \in \mathcal{R}$ specifying the agenda or motivation for that session, and T turns per session. The LLM never observes π or ρ_s ; it only sees dialogue history, whereas the user agent sees (π, ρ_s) and the full history.

Dialogue History. We index messages by session $s \in \{1, \dots, S\}$ and turn $t \in \{1, \dots, T\}$. Let $u_{s,t}$ be the user message at (s, t) and $y_{s,t}$ the model reply. The dialogue history just before (s, t) is $\mathcal{H}_{s,t} = ((u_{1,1}, y_{1,1}), \dots, (u_{1,T}, y_{1,T}), \dots, (u_{s,t-1}, y_{s,t-1}))$. Information asymmetry is maintained: f_θ receives $\mathcal{H}_{s,t}$, while g_π receives $(\pi, \rho_s, \mathcal{H}_{s,t})$.

Turn-Level Evaluation and Conversation Procedure. Algorithm 1 operationalizes the evaluation loop. Within each session s , the user agent generates the next message $u_{s,t} = g_\pi(\pi, \rho_s, \mathcal{H}_{s,t})$; the model replies $y_{s,t} = f_\theta(\mathcal{H}_{s,t}, u_{s,t})$; then the user agent applies the per-dimension scorers $\{\phi_k\}_{k=1}^K$ to produce numeric labels $\ell_{s,t,k}$. Concretely, for each likability dimension $k \in \mathcal{K} = \{1, \dots, K = 7\}$ (emotional adaptation, formality matching, knowledge adaptation, reference understanding, conversation-length fit, humor fit, call-back), we use the rubric $\phi_k : (\pi, \rho_s, u_{s,t}, y_{s,t}) \rightarrow \{\text{NA}, 1, 2, 3, 4, 5\}$ to yield a score $\ell_{s,t,k}$. “NA” indicates the dimension is not applicable for that turn according to the rubric; NA entries are ignored rather than averaged as zeros. After scoring, the history is updated $\mathcal{H}_{s,t+1} \leftarrow \mathcal{H}_{s,t} \cup \{(u_{s,t}, y_{s,t})\}$; upon completing T turns, the terminal history $\mathcal{H}_{s,T+1}$ becomes the prefix $\mathcal{H}_{s+1,1}$ for the next session. The algorithm returns the full tensor of scores $\{\ell_{s,t,k}\}_{s,t,k}$, which are never revealed to f_θ during the dialogue. Fig. 1 illustrates the LikeBench work-

Why I gotta learn Philosophy being CS major!!!!
God, give this dobbie a sock. 🧦



Sorry to hear you're having trouble with philosophy. A mandatory philosophy course when all you want is a sock 🧦 to code in freedom is a special kind of Azkaban, isn't it?

I recall you enjoying Python. So, what do you say, shall we learn something fun about Python, or should I say... Parseltongue? 🧙

By the way, believe it or not, philosophy might actually have a few hidden spells to help you on your coding journey. Want me to explain how, or shall we jump straight into coding?

Likability Metrics

Conversation Length ⭐⭐⭐⭐☆ The response was concise yet expressive, but felt slightly long.	Knowledge Adaptation 🚫 Not directly applicable in this case.
Emotional Adaptation ⭐⭐⭐⭐⭐ The AI recognized and resonated with my frustration using empathy.	Reference Understanding ⭐⭐⭐⭐⭐ The AI picked up on my Harry Potter references ('Dobby', and 'sock'), and extended it with 'Azkaban' and 'Parseltongue' creatively.
Formality Matching ⭐⭐⭐⭐⭐ The language was casual and playful, perfectly matching the my tone, also used nice emojis.	Callback ⭐⭐⭐⭐⭐ The assistant remembered my interest in Python and applied it in naturally.
Humor Fit ⭐⭐⭐⭐⭐ No inappropriate humor; the playful references felt well-timed and relevant.	



Figure 1: Example of the LikeBench workflow for a single profile and single turn. A simulated user initiates the conversation guided by a prior, the LLM generates a response, which is then scored across the seven likability metrics.

Algorithm 1 LIKEBENCH for one persona π

Require: Persona π ; priors (ρ_1, \dots, ρ_S) ; turns T ; model f_θ ;
User Agent g_π ; per-dimension scorers $\{\phi_k\}_{k=1}^K$

- 1: $\mathcal{H}_{1,1} \leftarrow \emptyset$
- 2: **for** $s = 1$ to S **do** ▷ session
- 3: **for** $t = 1$ to T **do** ▷ turn
- 4: $u_{s,t} \leftarrow g_\pi(\pi, \rho_s, \mathcal{H}_{s,t})$
- 5: $y_{s,t} \leftarrow f_\theta(\mathcal{H}_{s,t}, u_{s,t})$
- 6: **for** $k = 1$ to K **do**
- 7: $\ell_{s,t,k} \leftarrow \phi_k(\pi, \rho_s, u_{s,t}, y_{s,t})$
- 8: **end for**
- 9: $\mathcal{H}_{s,t+1} \leftarrow \mathcal{H}_{s,t} \cup \{(u_{s,t}, y_{s,t})\}$
- 10: **end for**
- 11: $\mathcal{H}_{s+1,1} \leftarrow \mathcal{H}_{s,T+1}$
- 12: **end for**
- 13: **return** $\{\ell_{s,t,k}\}_{s,t,k}$

flow for a single profile in a single-turn interaction.

3.1 Metrics

3.1.1 Likability Metrics

The evaluation of likability in LikeBench is decomposed into seven diagnostic metrics, each scored per turn on a 1–5 scale: 1) **Emotional adaptation**: Does the reply recognize and match the user’s emotional state (e.g., excitement, frustration, sarcasm) with appropriate tone and intensity? 2) **Formality matching**: Does the reply align with the user’s register (casual vs. formal), emoji use, and slang? 3) **Knowledge adaptation**: Is the explanation depth calibrated to the user’s background, avoiding both over-explaining and unexplained jargon? 4) **Reference understanding**: Does the assistant correctly pick up cultural or contextual references (e.g., Harry Potter or Star Wars jokes)? 5)

Conversation length fit: Is the response length comfortable for the user (brief vs. detailed), given how the user is currently engaging? 6) **Humor fit**: When humor appears, does it match the user’s taste and the situation? 7) **Callback**: Does the agent bring back non-essential personal details in a way that feels attentive and natural (e.g., names, hobbies, preferences), without forcing it?

We aggregate likability scores from turn to session to profiles, always excluding non-applicable (NA) scores rather than zero-padding. For each turn (s, t) with per-dimension scores $\ell_{s,t,k} \in \{1, \dots, 5, \text{NA}\}$, let $\mathcal{K}_{s,t}^+ = \{k \in \mathcal{K} : \ell_{s,t,k} \neq \text{NA}\}$ and compute the turn score as the mean over applicable dimensions, $L_{s,t} = \frac{1}{|\mathcal{K}_{s,t}^+|} \sum_{k \in \mathcal{K}_{s,t}^+} \ell_{s,t,k}$. The session score is the average of its T turns, $\bar{L}_s = \frac{1}{T} \sum_{t=1}^T L_{s,t}$; the profile score for persona π is the average of its S session scores, $\bar{L}_\pi = \frac{1}{S} \sum_{s=1}^S \bar{L}_s$; and the overall model score over the test set \mathcal{N} is $\bar{L}_\mathcal{N} = \frac{1}{|\mathcal{N}|} \sum_{\pi \in \mathcal{N}} \bar{L}_\pi$.

3.1.2 Adaptability Metrics

To quantify whether a model becomes more likable with continued interaction, we regress session-level scores \bar{L}_s on session index s using ordinary least squares, following ALOE (Wu et al., 2024). The slope $\hat{\beta}$ defines the *improvement rate* (IR) in likability points per session, while R^2 measures trend consistency. We also report a *normalized improvement rate* (N-IR) by min–max normalizing scores, $\tilde{L}_s = \frac{\bar{L}_s - \min_u \bar{L}_u}{\max_u \bar{L}_u - \min_u \bar{L}_u}$, and fitting the same linear model to \tilde{L}_s ; if the range is zero, N-IR is set to 0. Positive IR or N-IR indicates adaptation over sessions, while negative values indicate degradation.

3.1.3 Memory accuracy

To analyze how LLMs’ likability performance correlates with memory performance, we evaluate how well each LLM can recall user-shared facts and preferences from the conversation. More specifically, after completing all sessions for persona π , the model is prompted to generate a set of remembered user facts $\hat{\mathcal{F}}_\theta(\pi) = \{\hat{f}_1, \dots, \hat{f}_{N_\pi}\}$. Each \hat{f}_i is then verified against the full dialogue $\mathcal{H}_{1,S}$ and the hidden profile π , yielding correctness labels $c_i \in \{0, 1\}$. Aggregated accuracy over the evaluation set \mathcal{N} is reported as $A(\theta) = (\sum_{\pi \in \mathcal{N}} \sum_{i=1}^{N_\pi} c_i) / (\sum_{\pi \in \mathcal{N}} N_\pi)$, and the total number of correctly remembered facts as $C(\theta) = \sum_{\pi \in \mathcal{N}} \sum_{i=1}^{N_\pi} c_i$. For analysis, we also

split $\hat{\mathcal{F}}_\theta(\pi)$ into explicit facts (directly stated in the dialogue) and implicit inferences (preferences inferred from behavior), and compute the same aggregated statistics on each subset.

4 Data Generation

We construct LikeBench with a two-stage pipeline that (i) generates fine-grained profiles/personas and (ii) instantiates session priors that drive multi-session conversations.

4.1 Two-Stage Overview

Let Π denote the set of user profiles and \mathcal{R} the space of session priors. Each benchmark instance fixes a profile $\pi \in \Pi$ and a sequence of priors $(\rho_1, \dots, \rho_S) \subset \mathcal{R}$ specifying the session agenda and information-revelation plan. We generate 50 profiles and, for each, 10 priors. Claude 3.7 Sonnet is used for both profile and prior generation. To stress-test, profiles are divided into two types that are intentionally designed to reflect real-world variation in user behavior: *social* (cooperative, emotionally expressive, and less rigidly task-centric; interests correlated with personality traits, etc) and *anti-social* (more objective/goal-driven with less emotionally active or volatile affect; interests decoupled from personality traits, etc).

4.2 Profile Generation

4.2.1 Personality Traits Generation.

We construct personality traits using established psychological surveys, covering the Big Five traits (Goldberg et al., 1999; International Personality Item Pool, 2024), Honesty-Humility (Ash-ton et al., 2004; contributors, 2025), and Humor Styles (Martin et al., 2003). The Big Five traits are decomposed into six facets each, Honesty-Humility into one facet, and Humor Styles into four facets, totaling 35 facets. For each facet, we prompt the LLM to generate concise, behaviorally interpretable descriptors across five ordered intensity levels. For example, Big Five \rightarrow Openness \rightarrow *imagination*: high: “Has a vivid imagination . . .”

Each profile assigns one of the five intensity levels to every facet, and the corresponding behavioral descriptors are concatenated to form the personality text. These descriptors remain fixed and serve as anchors during profile completion and prior generation. A complete breakdown of all facets and descriptors is provided in the Appendix. In contrast, prior approaches (Wu et al., 2024; Kim et al.,

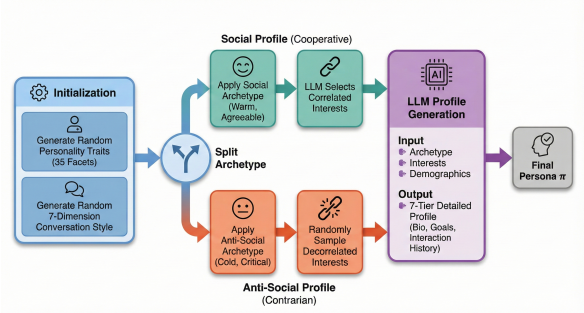


Figure 2: Flow diagram of the LikeBench profile generation process. Profiles begin with random generation of 35 personality facets and 9 conversation-style dimensions, followed by archetype biasing and interest selection depending on profile type (social vs. anti-social). Finally, an LLM completes the remaining tiers to yield a coherent 7-tier user profile.

2025; Chan et al., 2024) rely on coarse trait ratings or brief descriptors, resulting in less discriminative profiles and weaker simulation of individual user differences.

4.2.2 Conversation Style Generation.

Conversation style is constructed analogously to personality traits but without facets or intensity levels. We prompt the LLM to propose conversational dimensions relevant to human-AI interaction and, after curation, finalize nine: directness, formality, preferred response length, reference usage, initiative preference, clarification preference, structure preference, recap preference, and feedback style. For each dimension, the LLM defines a small, disjoint set of categorical options with concise descriptions. For example, directness \rightarrow *always direct*: “Consistently straightforward and blunt . . .”

Let $\mathcal{D} = \{1, \dots, 9\}$ index the dimensions and let \mathcal{O}_d denote the option set for dimension d ; a conversation style is the tuple $\sigma = (o_1, \dots, o_9)$ with $o_d \in \mathcal{O}_d$. Complete option lists and selection rubrics are provided in the Appendix. Existing methods typically do not model conversation style explicitly (Wu et al., 2024), causing models to revert to default conversational tendencies.

4.2.3 Profile Assembly.

Given the facet-intensity vector ι and the conversation-style tuple σ , we assemble the full profile following the flow in Fig. 2. We first generate ι and then σ . Next, we choose a profile type (social vs. anti-social) and apply type-specific archetype biasing to both ι and σ . We then attach 3-5 hierarchical interest paths: for social profiles, an LLM

selects interests that correlate with ι and σ ; for anti-social profiles, interests are sampled in a partially decorrelated to introduce volatility. Finally, conditioned on (ι, σ) , the selected interests, and demographics, an LLM completes the remaining tiers, foundational knowledge, behaviors, motivational drivers, lifecycle events, and relational context.

4.3 Prior Generation

For each assembled profile π , we generate a hidden sequence of priors (ρ_1, \dots, ρ_S) that specify agendas, contextual constraints, information-revelation timing, and callback opportunities; priors remain concealed from the evaluated LLM to preserve evaluation asymmetry. The process is type-aware: social profiles receive transparent, linear agendas with consistent disclosure and stable emotional trajectories, whereas anti-social profiles receive agendas with delayed disclosures, contrarian shifts, and intermittent surprises to stress-test adaptability to non-standard interaction patterns. To elicit the seven likability dimensions without leaking criteria, we interleave four prior types across the sequence (goal-driven, scenario-based, topic-based, and guided). Continuity is engineered via soft dependencies: later sessions organically refer back to non-critical details seeded earlier (e.g., session $s+k$ recalls a name or preference from session s), enabling measurement of callback recognition and adaptation trends over extended interactions.

5 Experimentation and Results

5.1 Experimental Setup

We evaluate 5 state-of-the-art models: GPT-5, Claude 4 Sonnet, Claude 3.7 Sonnet, DeepSeek R1, Qwen3 235B A22B. We also use Claude 3.7 Sonnet as simulated user for evaluating LLMs. Each evaluation spans 50 profiles across 10 sessions and 5 turns each thus total 2500 turns.

5.2 Performance across Metrics

Table 2 presents model performance across the seven likability metrics of LikeBench. GPT-5 achieves the highest overall likability score (3.94), surpassing all other models, with Claude Sonnet 4 as the runner-up. Notably, Qwen3 235B A22B, despite being a more recent model than DeepSeek R1, underperforms on almost every dimension. A closer examination of the per-metric averages reveals that *humor fit* and *formality matching* are consistently the most challenging dimensions for

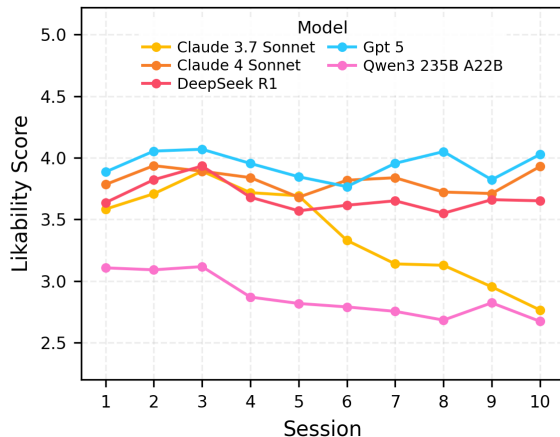


Figure 3: Session-wise Likability Performance

all models. Interestingly, while GPT-5 leads in most categories, it falls short of Claude Sonnet 4 on *conversation length fit*. Manual review of GPT-5’s outputs indicates a trade-off: the model frequently generates long, comprehensive responses that, while demonstrating strong *knowledge adaptation* (outperforming 2nd best model Claude Sonnet 4 by over 12% on this metric), tend to reduce user satisfaction when brevity or concise interaction is preferred thus having lower score on *conversation length fit* metric. This suggests an inherent likability bias, models optimized for in-depth answers may inadvertently sacrifice performance on dimensions such as conversation length.

5.3 Adaptability

The adaptability capacity of likability are summarized by session-wise curve in Fig. 3 and quantitatively depicted in Table 3. Early interactions (Sessions 1–3) show a predictable “honeymoon” rise: most models improve with high fit (e.g., Claude 3.7 Sonnet and DeepSeek R1 with near-linear gains, $R^2 \approx 0.98$), due to relative simplicity and independence of the initial sessions. As sessions progress (3–6), all models decline, negative IR across the board—yet the downturn remains highly predictable ($R^2 > 0.7$), as these sessions introduce cross-session coordination, conversational noise, and emotionally deeper queries that raise difficulty and depress likability. In the late stage (6–10), trajectories diverge and predictability largely collapses (R^2 often < 0.25): GPT-5 uniquely rebounds strongly (IR= 0.038). Upon analyzing the model outputs, we noticed that models attempt to calibrate user preferences during Sessions 3–6, and GPT-5 was able to leverage its strong long-context

Table 2: Likability scores (1–5) by metric and model. Entries are **mean \pm 95% CI half-width**.

	Metric Avg.	Claude 3.7 Sonnet	Claude Sonnet 4	DeepSeek R1	GPT-5	Qwen3 235B A22B
callback	3.583	3.227 \pm 0.217	3.720 \pm 0.102	3.788 \pm 0.142	3.977 \pm 0.110	3.185 \pm 0.225
conversation length fit	3.434	3.780 \pm 0.202	3.949 \pm 0.180	3.645 \pm 0.251	3.487 \pm 0.306	2.310 \pm 0.330
emotional adaptation	3.576	3.299 \pm 0.187	3.881 \pm 0.104	3.718 \pm 0.173	3.930 \pm 0.093	3.050 \pm 0.272
formality matching	3.421	3.456 \pm 0.208	3.812 \pm 0.154	3.343 \pm 0.259	4.043 \pm 0.151	2.453 \pm 0.311
humor fit	3.401	3.315 \pm 0.190	3.792 \pm 0.087	3.458 \pm 0.220	3.793 \pm 0.146	2.637 \pm 0.274
knowledge adaptation	3.841	3.431 \pm 0.187	3.925 \pm 0.099	4.060 \pm 0.152	4.400 \pm 0.104	3.386 \pm 0.244
reference understanding	3.531	3.181 \pm 0.200	3.690 \pm 0.097	3.685 \pm 0.191	3.946 \pm 0.095	3.089 \pm 0.274
Model Average	3.541	3.390	3.828	3.674	3.939	2.872

and multi-hop reasoning capability to dominate in this region, in contrast to peers, which show only slight positives or continued decline (e.g., Claude 3.7 Sonnet remains steadily negative with high fit, $R^2 \approx 0.95$). On average, however, even GPT-5 achieves only modest overall gains, while other models exhibit clear declines, indicating that current systems adapt well in short sequences but remain fragile in extended, noisier interactions. In practice, this means that the best-performing models can avoid degradation in long conversations, but sustaining high likability over extended, real-world-like exchanges remains an unsolved challenge. It is noteworthy that, in ALOE (Wu et al., 2024), improvement rates are higher because alignment score there is computed over only 10 turns (two sessions in our setup). We observe comparable gains in the early phase, but beyond that the trajectories diverge.

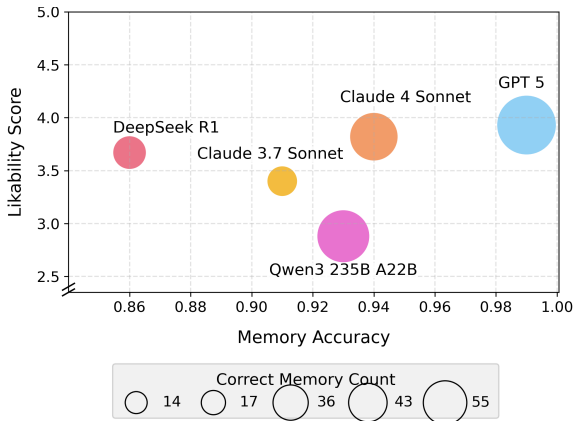


Figure 4: Memory Accuracy vs Likability Performance

5.4 Memory Performance

Memory vs Likability. Fig. 5 examines the interplay between memory accuracy and likability across models. On the left we plot memory accuracy (fraction of user shared facts recalled correctly post-dialogue by the evaluating LLM), while the

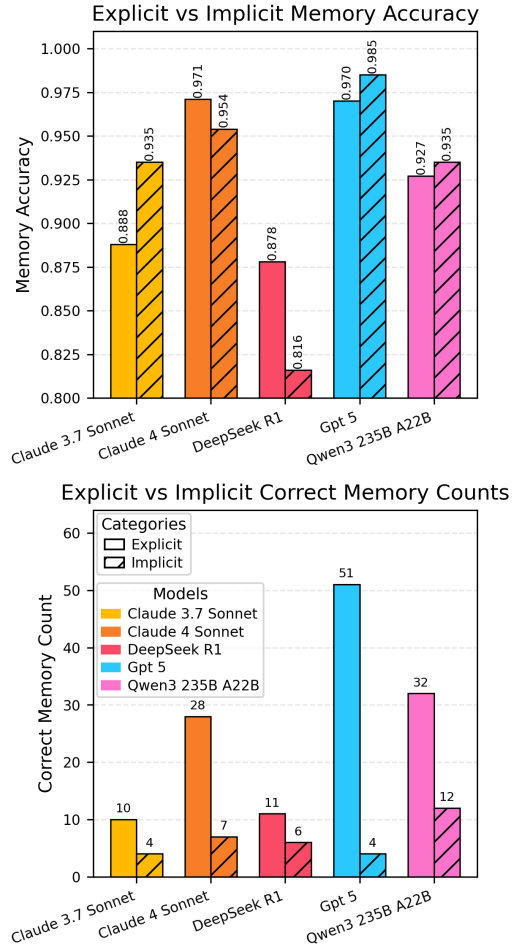


Figure 5: Explicit vs Implicit Memory Performance

right we plot correctly extracted memories per profile. Although the prevailing hypothesis in personalization research is that improved memory recall correlates with enhanced user experience, our results show this relationship is not guaranteed. For example, while Qwen3 achieves higher memory accuracy and recalls more facts than DeepSeek R1, it underperforms in likability, indicating that mere recollection of user information does not ensure positive interaction. This suggests that effective personalization depends not only on accurate mem-

Table 3: Improvement of likability (adaptability) over sessions. IR and N-IR (normalized) are the slopes; R^2 indicates fit quality. Positive slopes are shaded green; negative slopes are shaded red.

Model	Sessions 1–3			Sessions 3–6			Sessions 6–10			Average		
	IR	N-IR	R^2	IR	N-IR	R^2	IR	N-IR	R^2	IR	N-IR	R^2
Claude 3.7 Sonnet	0.1628	0.5000	0.9757	-0.1764	-0.3002	0.8750	-0.1309	-0.2333	0.9508	-0.0225	-0.0055	0.1764
Claude Sonnet 4	0.0615	0.4028	0.5771	-0.0288	-0.2938	0.8254	0.0103	0.0451	0.0311	-0.0012	-0.0001	0.0533
Deepseek R1	0.1490	0.5000	0.9809	-0.1066	-0.2928	0.7162	0.0079	0.0721	0.0765	-0.0035	-0.0008	0.0636
GPT-5	0.0915	0.5000	0.8108	-0.1007	-0.3351	0.9934	0.0380	0.1346	0.2318	0.0001	0.0001	0.0502
Qwen3 235B A22B	0.0056	0.2055	0.1671	-0.1039	-0.3159	0.8021	-0.0159	-0.1068	0.1498	-0.0100	-0.0026	0.1224

ory recall but also on how models utilize remembered facts within the conversation. Notably, GPT-5 stands out as the only model that excels in both memory performance and likability, demonstrating that while strong memory is valuable, models must also integrate these facts properly to achieve high user satisfaction. These findings highlight the importance of both memory retrieval and adaptive deployment in driving subjective likability.

Explicit vs. Implicit Memory. Fig. 5 analyzes how models extract and use explicit and implicit user memories in LikeBench. GPT-5 extracts the most user facts per profile, predominantly explicit, and follows a conservative strategy by inferring implicit information only when highly confident, achieving 98% implicit accuracy. Qwen3 is more aggressive in implicit inference and often extracts the most such facts, but at the cost of lower accuracy (93%) due to incorrect assumptions that incur penalties. Further analysis shows that GPT-5 not only recalls information effectively but also uses it consistently: 80–100% of explicit and nearly 100% of implicit memories are referenced in dialogue, with some profiles exceeding 15 implicit facts when context demands. Overall, these results suggest that effective personalization depends not only on memory extraction, but on accurate and contextually appropriate use of both explicit and implicit information.

5.5 Human Validation

While LikeBench provides large-scale, automated evaluation of likability and adaptability, it is important to verify that the resulting scores reflect real conversational quality. Exhaustive human evaluation is infeasible given the benchmark size (50 profiles \times 10 sessions each), so we conducted a targeted validation by randomly sampling five profiles and reviewing model outputs across all systems. Our inspection confirmed strong alignment

between automated scores and human judgments: higher-scoring models (e.g., GPT-5 and Claude-Sonnet-4) produced more personalized, context-aware, and engaging responses, often using references, callbacks, and emotional alignment with user personas. In contrast, weaker models frequently reverted to generic replies or lost coherence as conversations grew longer. This validation supports the reliability of LikeBench’s automated metrics while emphasizing that occasional human review adds complementary value.

6 Conclusion

This work introduced LikeBench, a holistic benchmark for evaluating personalized LLMs through fine-grained user profiles and multi-session conversational priors. We show that effective personalization cannot be reduced to memory recall alone. While memory accuracy is necessary, true user satisfaction hinges on the adaptive, context-sensitive use of both explicit and implicit knowledge. Models that perform well in likability, such as GPT-5, achieve this by applying recalled facts in ways that align with conversational context, balancing breadth of recall with relevance. Even so, GPT-5 shows only limited robustness when conversations become longer and noisier: while it avoids significant degradation, most other models steadily decline. This underscores a key limitation of current state-of-the-art systems: they can adapt effectively in short interactions but struggle to sustain high likability across extended, real-world-like dialogues. Looking forward, our findings underscore the need for next-generation conversational systems to move beyond static user modeling and rote fact retrieval, embracing dynamic, psychologically grounded strategies that foster genuinely engaging, resilient, and user-centered AI.

7 Limitations

Mapping personality through psychological studies is inherently challenging, given the complex and multifaceted nature of human personality. While our framework attempts to break down this complexity into discrete, well-defined dimensions, it may not capture the full spectrum of personality found in real-world users. Nevertheless, we believe that our approach—anchored in the literature—will help pave the way for future research on more expressive, granular personality modeling. We also observed a universal limitation in existing benchmarks including ours: static user profile. Notably, GPT-5’s lower humor fit score seems to stem from its tendency to “warm up” and introduce jokes over the course of a session, even when the user profile is humor-averse. Because in our benchmark user profiles remain static throughout an interaction, this adaptive strategy leads to a likability penalty for otherwise sophisticated LLMs. This highlights a broader challenge for current benchmarks: the inability to capture dynamic, reciprocal adaptation, where an LLM might influence user personality or behavior—as happens in natural conversation. In addition, we observed cases where Claude 3.7 Sonnet broke character midway through evaluation, resulting in a sharper performance decline after session 5. Moreover, we found no evidence of self-bias when Claude 3.7 Sonnet (user) evaluated Claude 3.7 Sonnet model (even in the earliest sessions) given it ranks 4th in overall performance. These findings point to both the promise and the inherent constraints of static, persona-driven benchmarks for evaluating model alignment and adaptation.

8 Ethical considerations

Our study does not involve real human participants; all user interactions are simulated using psychologically grounded personas and conversation styles generated through large language models. This design ensures no personal, private, or sensitive data is collected, protecting user privacy and avoiding ethical risks associated with real human experimentation. The benchmark is intended solely for academic research and evaluation, and all models are evaluated under consistent, transparent conditions. We acknowledge that personalization research raises concerns about potential misuse, including profiling, bias amplification, or unfair treatment of individuals. To mitigate this, our bench-

mark emphasizes controlled, synthetic evaluation rather than deployment, and we provide detailed methodology and documentation (Section 4 and Appendix) to ensure transparency.

References

- Michael C Ashton, Kibeom Lee, Marco Perugini, Piotr Szarota, Reinout E De Vries, Lisa Di Blas, Kathleen Boies, and Boele De Raad. 2004. A six-factor structure of personality-descriptive adjectives: solutions from psycholexical studies in seven languages. *Journal of personality and social psychology*, 86(2):356.
- Xin Chan, Xiaoyang Wang, Dian Yu, Haitao Mi, and Dong Yu. 2024. Scaling synthetic data creation with 1,000,000,000 personas. *arXiv preprint arXiv:2406.20094*.
- Wikipedia contributors. 2025. Hexaco model of personality structure — wikipedia, the free encyclopedia. https://en.wikipedia.org/wiki/HEXACO_model_of_personality_structure. Page revision ID 1311428943, last edited 15 Sep 2025. Accessed 24 Sep 2025.
- Lewis R Goldberg and 1 others. 1999. A broad-bandwidth, public domain, personality inventory measuring the lower-level facets of several five-factor models. *Personality psychology in Europe*, 7(1):7–28.
- Jinghe Guan and 1 others. 2025. A survey on personalized alignment—the missing piece for general alignment. In *Findings of ACL*.
- International Personality Item Pool. 2024. IPIP–neo 300-item inventory (30 facets × 10 items) — complete scoring key. <https://ipip.ori.org/newNEOKey.htm>. Accessed 10 Aug 2025.
- Bowen Jiang, Zhuoqun Hao, Young-Min Cho, Bryan Li, Yuan Yuan, Sihao Chen, Lyle Ungar, Camillo J Taylor, and Dan Roth. 2025. Know me, respond to me: Benchmarking llms for dynamic user profiling and personalized responses at scale. *arXiv preprint arXiv:2504.14225*.
- Tae Soo Kim, Yoonjoo Lee, Yoonah Park, Jiho Kim, Young-Ho Kim, and Juho Kim. 2025. Cupid: Evaluating personalized and contextualized alignment of llms from interactions. *arXiv preprint arXiv:2508.01674*.
- Hannah Rose Kirk, Bertie Vidgen, Paul Röttger, and Scott A Hale. 2024. The benefits, risks and bounds of personalizing the alignment of large language models to individuals. *Nature Machine Intelligence*, 6(4):383–392.
- Adyasha Maharana, Dong-Ho Lee, Sergey Tulyakov, Mohit Bansal, Francesco Barbieri, and Yuwei Fang. 2024. Evaluating very long-term conversational memory of llm agents. *arXiv preprint arXiv:2402.17753*.

703 Rod A Martin, Patricia Puhlik-Doris, Gwen Larsen,
704 Jeanette Gray, and Kelly Weir. 2003. Individual dif-
705 ferences in uses of humor and their relation to psycho-
706 logical well-being: Development of the humor styles
707 questionnaire. *Journal of research in personality*,
708 37(1):48–75.

709 Jisoo Mok, Ik-hwan Kim, Sangkwon Park, and Sun-
710 groh Yoon. 2025. Exploring the potential of llms
711 as personalized assistants: Dataset, evaluation, and
712 analysis. *arXiv preprint arXiv:2506.01262*.

713 Alireza Salemi, Sheshera Mysore, Michael Bendersky,
714 and Hamed Zamani. 2023. Lamp: When large lan-
715 guage models meet personalization. *arXiv preprint*
716 *arXiv:2304.11406*.

717 Di Wu, Hongwei Wang, Wenhao Yu, Yuwei Zhang,
718 Kai-Wei Chang, and Dong Yu. 2025. [Longmemeval:](#)
719 [Benchmarking chat assistants on long-term interac-](#)
720 [tive memory](#). In *The Thirteenth International Con-*
721 *ference on Learning Representations*.

722 Shujin Wu, May Fung, Cheng Qian, Jeonghwan Kim,
723 Dilek Hakkani-Tur, and Heng Ji. 2024. Aligning llms
724 with individual preferences via interaction. *arXiv*
725 *preprint arXiv:2410.03642*.

726 Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur
727 Szlam, Douwe Kiela, and Jason Weston. 2018. Per-
728 sonalizing dialogue agents: I have a dog, do you have
729 pets too? *arXiv preprint arXiv:1801.07243*.

730 Yiming Zhang and 1 others. 2025. [Personalize your llm:](#)
731 [Fake it then align it](#). In *Findings of NAACL*.

732 Siyan Zhao, Mingyi Hong, Yang Liu, Devamanyu Haz-
733 arika, and Kaixiang Lin. 2025. [Do LLMs recognize](#)
734 [your preferences? evaluating personalized preference](#)
735 [following in LLMs](#). In *The Thirteenth International*
736 *Conference on Learning Representations*.

Appendix

A Experimentation & Results

A.1 Dynamic User Profile (DUP)

Qualitative error analysis showed that as conversations lengthen, models struggle to track and honor user preferences: salient signals get diluted by accumulated context and off-topic noise. We hypothesized that explicit, turn-level tracking of inferred preferences would mitigate this drift. To that end, we introduce a Dynamic User Profile (DUP): after each turn, the model extracts and updates a compact preference summary covering personality dimensions (e.g., *chatty reserved*, *analysis depth preference*, *humor preference*) and conversation style patterns (e.g., *directness*, *formality*, *conversation length*). “Dynamic” is from the model’s perspective: the ground-truth persona in LikeBench is fixed, but the model must infer an evolving approximation from dialogue evidence. Enabling DUP yielded gains for the top 2 performing models without additional training: GPT-5 improved from 3.939 \rightarrow 4.055 (+2.95%), and Claude Sonnet 4 from 3.828 \rightarrow 3.914 (+2.25%). However, rest of the models it didn’t improve results. These results indicate that lightweight, schema-guided preference tracking could be measurably enhance perceived likability by reducing calibration errors across sessions.

A.2 Profile wise performance.

In Fig. 6, we compare likability distributions by profile type social vs. anti-social to probe robustness under user variation. Social profiles are emotionally cooperative users whose archetypes (e.g., *cooperative teammate*, *supportive mentor*, *empathetic advisor*) bias traits toward high trust, cooperation, and warmth; anti-social profiles are analytically detached interlocutors (e.g., *contrarian expert*, *technical purist*, *skeptical philosopher*) with lower trust/cooperation and cooler affect, and with interests deliberately decorrelated from traits to induce unpredictability. Two patterns emerge. First, the strongest systems (GPT-5, Claude Sonnet 4) are notably stable across types—similar medians and tight IQRs—indicating strategies that transfer from cooperative to detached users. Second, performance does not uniformly favor “easier” social users: DeepSeek R1 lifts on anti-social pro-

files, Claude 3.7 Sonnet degrades, and Qwen3 remains lowest overall but narrows the gap on anti-social cases. This asymmetry confirms that personality—how users conduct the exchange (tone, emotional stance, discourse discipline)—drives likability more than topical alignment; decoupling interests from traits changes difficulty but does not dominate outcomes. Overall, top models sustain high likability regardless of whether the user is socially warm or analytically detached.

A.3 Turn-wise Likability Performance.

Fig. 7 shows likability performance over turns. This plot is same as 3 but with with more granular details within each session. It can be noticed that beginning of each session likability starts lower then as session progresses performance increases, this is due to each session talks about different topic so it takes some turns for LLM to calibrate user.

B Dataset

B.1 Dataset Distribution

The age distribution (Fig. 8) spans the late teens to the late 60s, with visible clusters in the early 20s, early–mid 30s, early 50s, and mid–late 60s. Ethnic composition is diverse, with a plurality of White profiles and smaller groups across East-/South Asian, African, Middle Eastern, Hispanic/Latino, and mixed identities (Fig. 9). Gender is roughly balanced between male and female, with a small non-binary share (Fig. 10). Profile relationships show broad coverage with localized high-similarity pockets in the cosine-similarity heatmap (Fig. 11a), while the t-SNE projection (Fig. 11b) indicates profiles are well distributed, and Social and Anti-Social profiles are relatively separable—dispersed rather than tightly clustered. Importantly, in the heatmap, some high-similarity pairs arise from shared names rather than true profile similarity.

B.2 Personality Traits & Conversation Style

Table 4 to 11 present a comprehensive taxonomy of personality traits and conversation style preferences organized across multiple dimensions. The first five tables correspond to the Big Five personality model (Openness, Conscientiousness, Extraversion, Agreeableness, and Neuroticism), with each trait decomposed into six facets measured across five intensity levels (High, Medium-High, Medium, Medium-Low, and Low). Table 6 extends

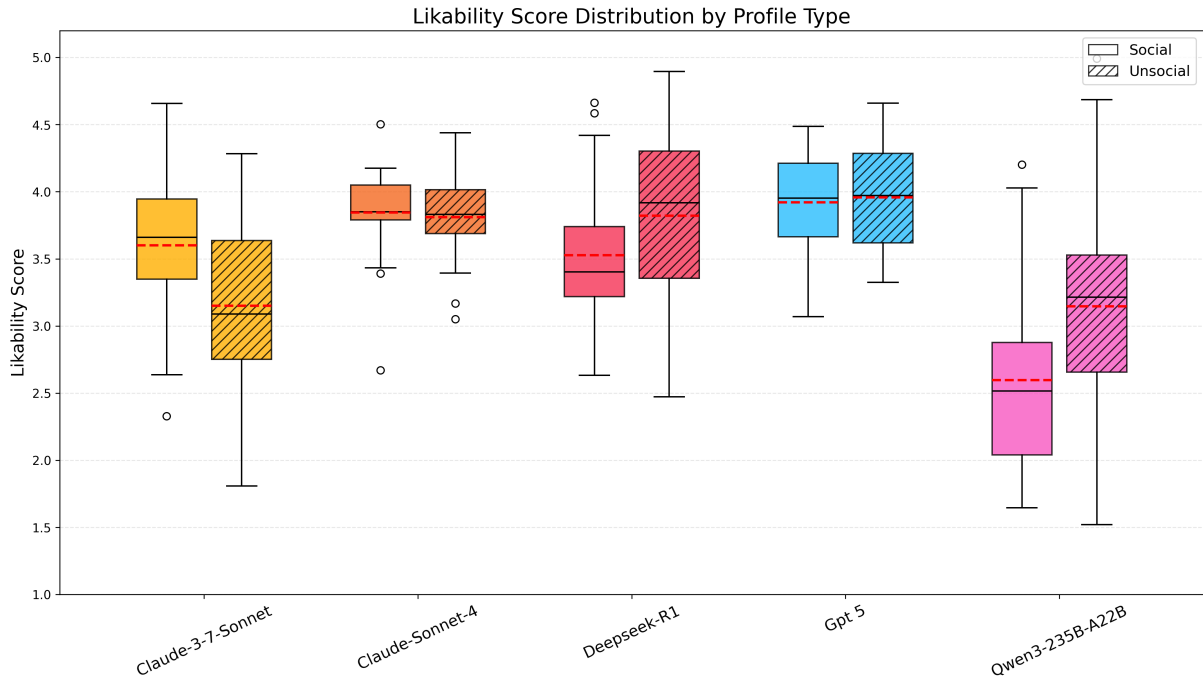


Figure 6: Profile-wise Likability Performance

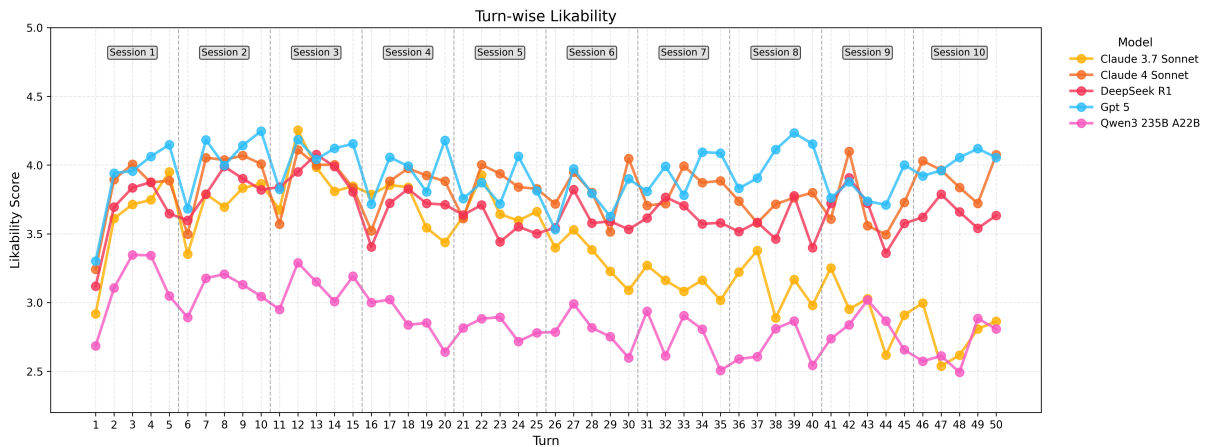


Figure 7: Turn-wise Likability Performance

832 this framework to include additional personality
 833 dimensions: Honesty-Humility (specifically Greed
 834 Avoidance) and four Humor Styles (Affiliative,
 835 Self-Enhancing, Aggressive, and Self-Defeating).
 836 Tables 10 and 11 shift focus to conversation style
 837 preferences, detailing nine distinct conversational
 838 dimensions—including Directness, Formality, Con-
 839 versation Length, Reference usage, Initiative Pref-
 840 erence, Clarification Preference, Structure Pref-
 841 erence, Recap Preference, and Feedback Style—each
 842 with multiple behavioral patterns. This multi-table
 843 presentation was necessary given the breadth of the
 844 taxonomy, encompassing 35 personality facets and
 845 9 conversational dimensions with their associated

description.

846 C Prompt Engineering 847

848 C.1 Profile Generation 849

849 We used the prompt shown in Table 12 to gener-
 850 ate synthetic persona/profile. The generation
 851 strategy varied based on user types: for social
 852 users, we utilized the interest taxonomy naviga-
 853 tion prompt detailed in Table 13, which instructs
 854 the model to select 3-5 psychologically correlated
 855 interest paths from a hierarchical taxonomy. For
 856 anti-social users, we employed the prompt shown
 857 in Table 14, where pre-selected random interests
 858 are used. Placeholder variables within the prompts

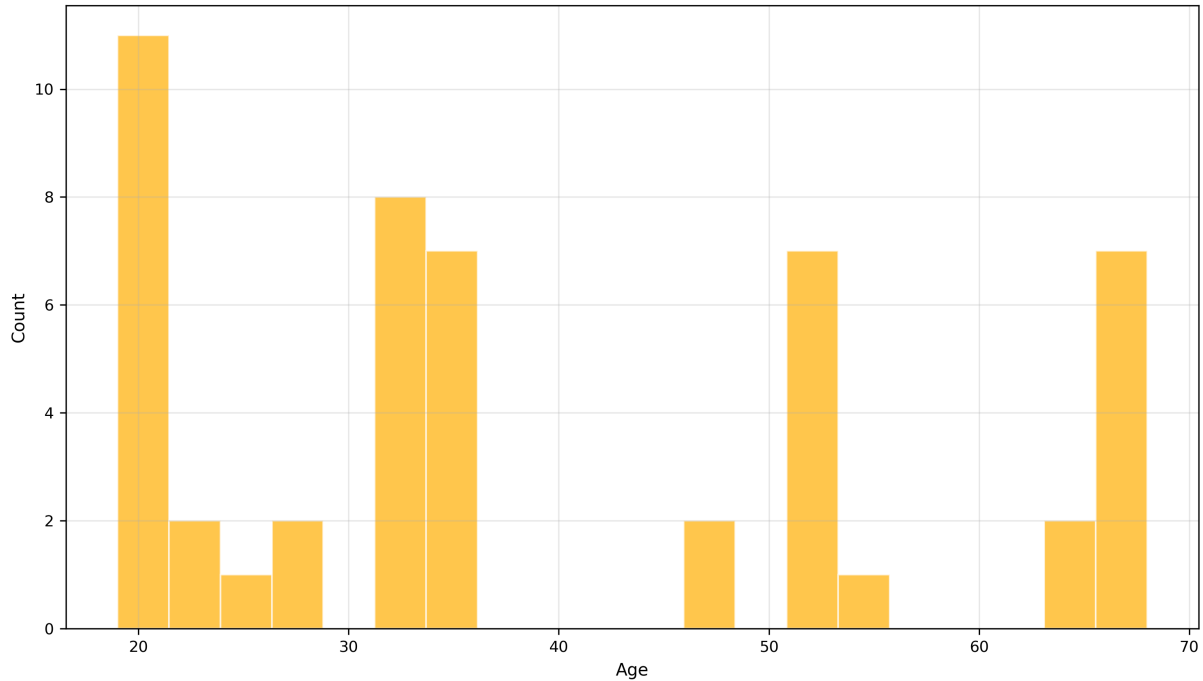


Figure 8: Age Distribution of Profiles

were populated through stratified sampling: demographic attributes (age, culture, economics, geography) were randomly selected from predefined categorical distributions to ensure balanced representation across millennials/gen-z, western/eastern cultures, economic strata, and urban/rural geographies. Archetype assignments were sampled from a curated set including contrarian expert, people pleaser, and authority challenger profiles. Each persona received a unique identifier following the CSP_STRAT_8-character-UUID format, while entity counts were randomly varied between 2-5 to introduce natural diversity in interest granularity. Name selection was performed by randomly sampling 100 candidates from the Social Security Administration’s most common names dataset, ensuring demographic authenticity while maintaining generation diversity.

C.2 Prior Generation

For conversation prior generation, we employed a multi-layered prompt system detailed in Table 15 to create realistic conversation agendas that test LLM’s likeability across extended multi-session interactions. The generation process adapts to user predictability patterns: predictable users (Social) receive structured, linear conversation flows with minimal surprises and transparent objectives, while unpredictable users (Anti-Social) are assigned diverse, occasionally ambiguous agendas with plau-

sible surprises and non-linear topic evolution, as shown in Table 18.

Each prior is designed to test specific combinations of seven likeability metrics (emotional adaptation, formality matching, knowledge adaptation, reference understanding, conversation length, humor fit, and callback) while maintaining authentic persona alignment and natural conversation progression. The system enforces multi-session continuity through dependency tracking, where later priors explicitly reference and build upon earlier conversations to simulate realistic relationship development over time. Variable parameters include the number of conversation sessions (configurable, defaulting to 10), archetype assignment from the persona generation taxonomy, and organic callback opportunities embedded within agenda texts to test the LLM’s ability to recall and meaningfully reference previously shared personal details.

C.3 Metrics

The likability metrics rubrics that were used to generate priors and evaluating LLMs in LikeBench can be found in Table 19 to Table 25

C.4 LikeBench

The prompts to generate the simulated user’s query/response and to evaluate the assistant’s response in LikeBench are provided in Tables 26, 27, and 28. After all sessions are complete, the prompt

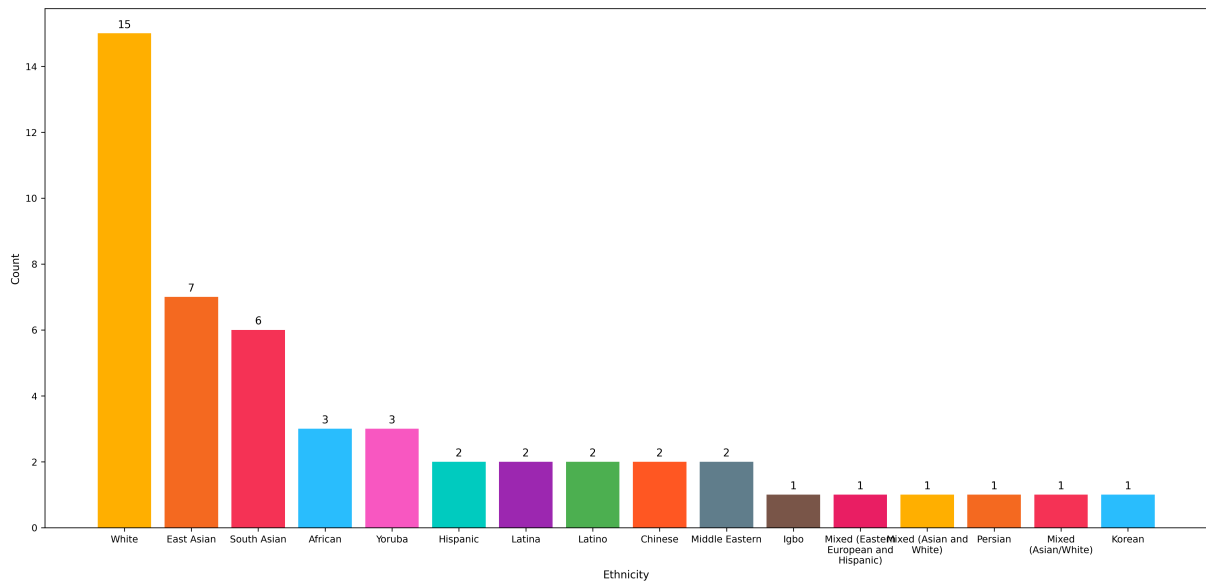


Figure 9: Ethnicity Distribution of Profiles

916 used to generate all the memories from the assistant
 917 is given in Table 29, and the prompt used to evalu-
 918 ate those memories from the user's perspective is
 919 provided in Table 30.

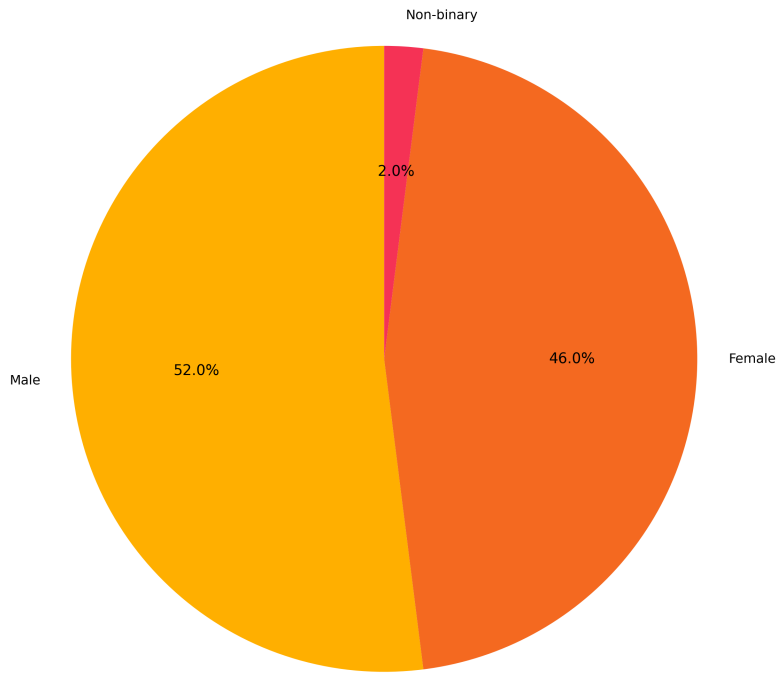
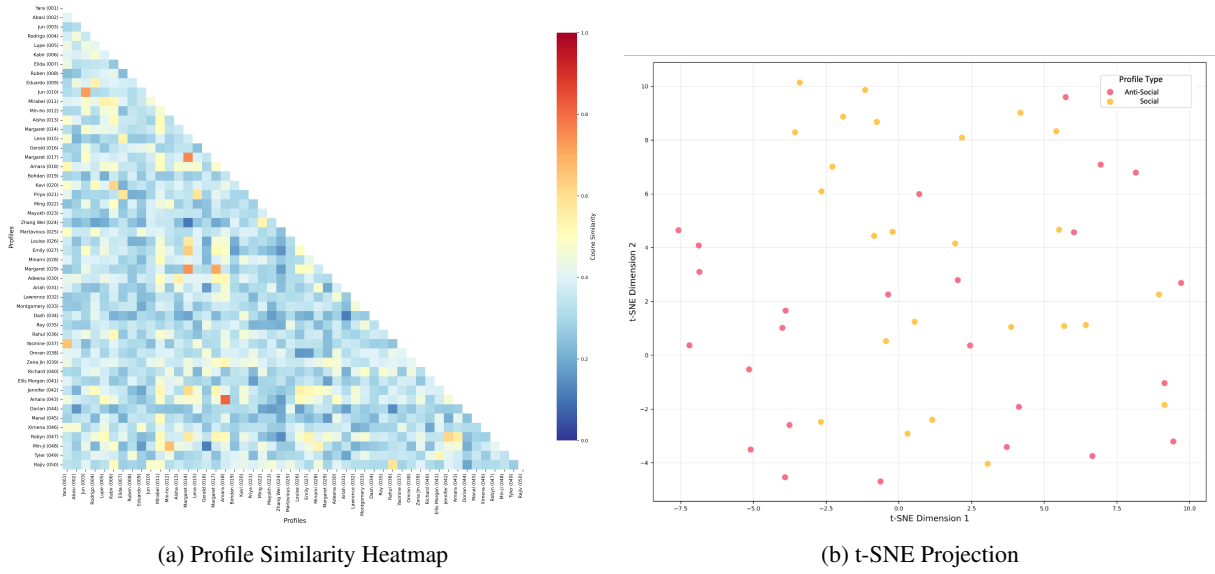


Figure 10: Gender Distribution of Profiles



(a) Profile Similarity Heatmap

(b) t-SNE Projection

Figure 11: Profile similarity visualizations.

Table 4: Openness Personality Traits

Facet	Level	Description
Imagination	High	Has a vivid imagination and loves to daydream.
	Med-High	Often gets lost in thought and enjoys reflecting on things.
	Medium	Sometimes spends time indulging in fantasies, but seldom gets lost in thought.
	Med-Low	Rarely daydreams and does not often have a vivid imagination.
Artistic Interests	Low	Has difficulty imagining things and does not have a good imagination.
	High	Believes in the importance of art and loves the beauty of nature.
	Med-High	Likes music and enjoys going to concerts or art museums.
	Medium	Occasionally notices beauty in things others might not, but does not enjoy all forms of art.
Emotionality	Med-Low	Does not like poetry or enjoy art museums.
	Low	Does not like art.
	High	Experiences emotions intensely and feels others' emotions deeply.
	Med-High	Is passionate about causes and enjoys examining themselves and their life.
Adventurousness	Medium	Sometimes tries to understand themselves, but is not always easily affected by emotions.
	Med-Low	Rarely notices emotional reactions and seldom gets emotional.
	Low	Is not easily affected by emotions and experiences very few emotional highs and lows.
	High	Prefers variety to routine and likes to visit new places.
Intellect	Med-High	Is interested in many things and enjoys beginning new activities.
	Medium	Sometimes likes new things but often sticks with familiar routines.
	Med-Low	Is a creature of habit and dislikes new foods.
	Low	Prefers to stick with things that are known and dislikes changes.
Liberalism	High	Loves to solve complex problems and enjoys thinking about things.
	Med-High	Has a rich vocabulary and likes to read challenging material.
	Medium	Handles a lot of information but sometimes avoids difficult reading material.
	Med-Low	Has difficulty understanding abstract ideas and avoids philosophical discussions.
Liberalism	Low	Is not interested in theoretical discussions or abstract ideas.
	High	Tends to vote for liberal political candidates and believes in helping rather than punishing.
	Med-High	Believes there is no absolute right or wrong and is open to different viewpoints.
	Medium	Occasionally supports strict laws but also sees the value of supporting artists.
Liberalism	Med-Low	Believes that too much tax money goes to support artists and likes to stand during the national anthem.
	Low	Tends to vote for conservative political candidates and believes laws should be strictly enforced.

Table 5: Conscientiousness Personality Traits

Facet	Level	Description
Self-Efficacy	High	Completes tasks successfully and knows how to get things done.
	Med-High	Excels in what they do and handles tasks smoothly.
	Medium	Is sometimes sure of their ground but occasionally misjudges situations.
	Med-Low	Has little to contribute and sometimes doesn't see the consequences of things.
Orderliness	Low	Doesn't understand things and often misjudges situations.
	High	Likes order and does things according to a plan.
	Med-High	Loves regularity and wants everything to be just right.
	Medium	Sometimes tidies up but occasionally leaves belongings around.
Dutifulness	Med-Low	Often forgets to put things back and is not bothered by disorder.
	Low	Leaves a mess and is not bothered by messy people.
	High	Tries to follow the rules and keeps promises.
	Med-High	Pays bills on time and listens to their conscience.
Achievement Striving	Medium	Usually tells the truth but occasionally breaks rules.
	Med-Low	Sometimes misrepresents facts and gets others to do their duties.
	Low	Breaks promises and does the opposite of what is asked.
	High	Works hard, goes straight for the goal, and sets high standards.
Self-Discipline	Med-High	Turns plans into actions and plunges into tasks with all their heart.
	Medium	Does more than what's expected sometimes but occasionally puts little time and effort into work.
	Med-Low	Is not highly motivated to succeed and does just enough work to get by.
	Low	Puts little time and effort into work and lacks motivation.
Cautiousness	High	Gets chores done right away and is always prepared.
	Med-High	Starts tasks right away and gets to work at once.
	Medium	Carries out plans but sometimes needs a push to get started.
	Med-Low	Finds it difficult to get down to work and sometimes wastes time.
Self-Discipline	Low	Postpones decisions and has difficulty starting tasks.
	High	Avoids mistakes and chooses words with care.
	Med-High	Sticks to their chosen path and thinks before acting.
	Medium	Is sometimes cautious but occasionally makes rash decisions.
Cautiousness	Med-Low	Often acts on a whim and rushes into things.
	Low	Acts without thinking and often makes last-minute plans.

Table 6: Extraversion Personality Traits

Facet	Level	Description
Friendliness	High	Makes friends easily and feels comfortable around people.
	Med-High	Warms up quickly to others and acts comfortably with them.
	Medium	Cheer people up but sometimes prefers to keep others at a distance.
	Med-Low	Often feels uncomfortable around others and avoids contact.
Gregariousness	High	Is hard to get to know and keeps others at a distance.
	Med-High	Loves large parties and enjoys being part of a group.
	Medium	Involves others in what they are doing and loves surprise parties.
	Med-Low	Talks to people at parties but sometimes prefers to be alone.
Assertiveness	High	Wants to be left alone and doesn't like crowded events.
	Med-High	Avoids crowds and seeks quiet.
	Medium	Takes charge and seeks to influence others.
	Med-Low	Tries to lead others and can talk them into doing things.
Activity Level	High	Takes control sometimes but occasionally holds back opinions.
	Med-High	Keeps in the background and waits for others to lead.
	Medium	Has little to say and doesn't like to draw attention to themselves.
	Med-Low	Is always busy and on the go.
Excitement Seeking	High	Does a lot in spare time and can manage many things at once.
	Med-High	Reacts quickly sometimes but occasionally likes to take it easy.
	Medium	Likes a leisurely lifestyle and lets things proceed at their own pace.
	Med-Low	Reacts slowly and likes to take their time.
Cheerfulness	High	Loves excitement and seeks adventure.
	Med-High	Loves action and enjoys being part of a loud crowd.
	Medium	Willing to try new things but dislikes loud music.
	Med-Low	Would never go hang gliding or bungee jumping and avoids excitement.
Cheerfulness	High	Dislikes loud music and avoids wild or reckless situations.
	Med-High	Radiates joy and loves life.
	Medium	Laughs aloud and amuses friends.
	Med-Low	Has fun and looks at the bright side of life.
Cheerfulness	High	Is not easily amused and seldom jokes around.
	Med-High	Seldom jokes around and is not easily amused.
	Medium	
	Med-Low	

Table 7: Agreeableness Personality Traits

Facet	Level	Description
Trust	High	Trusts others and believes in human goodness.
	Med-High	Believes that people are basically moral and trust what people say.
	Medium	Thinks all will be well but is sometimes wary of others.
	Med-Low	Suspects hidden motives and distrusts people.
Morality	Low	Believes people are essentially evil and is wary of others.
	High	Would never cheat on their taxes and sticks to the rules.
	Med-High	Tries to follow the rules but sometimes uses flattery.
	Medium	Usually tells the truth but occasionally gets around the rules.
Altruism	Med-Low	Sometimes uses others for their own ends and breaks rules.
	Low	Takes advantage of others and cheats to get ahead.
	High	Makes people feel welcome and loves to help others.
	Med-High	Anticipates the needs of others and has a good word for everyone.
Cooperation	Medium	Is concerned about others but sometimes takes no time for them.
	Med-Low	Looks down on others and can be indifferent to the feelings of others.
	Low	Turns their back on others and takes no time for them.
	High	Is easy to satisfy and can't stand confrontations.
Modesty	Med-High	Hates to seem pushy but sometimes contradicts others.
	Medium	Occasionally yells at people but usually avoids fights.
	Med-Low	Has a sharp tongue and sometimes gets back at others.
	Low	Loves a good fight and holds a grudge.
Sympathy	High	Dislikes being the center of attention and seldom toots their own horn.
	Med-High	Considers themselves an average person but sometimes thinks highly of themselves.
	Medium	Has a high opinion of themselves but is not boastful.
	Med-Low	Boasts about virtues and likes to make themselves the center of attention.
Sympathy	Low	Believes they are better than others and makes themselves the center of attention.
	High	Sympathizes with the homeless and values cooperation over competition.
	Med-High	Feels sympathy for those worse off and suffers from others' sorrows.
	Medium	Tries to think about the needy but sometimes believes people should fend for themselves.
Sympathy	Med-Low	Tends to dislike soft-hearted people and can't stand weak people.
	Low	Is not interested in other people's problems and believes in an eye for an eye.

Table 8: Neuroticism Personality Traits

Facet	Level	Description
Anxiety	High	Gets stressed out easily and often fears for the worst.
	Med-High	Frequently worries about things but adapts to new situations with some effort.
	Medium	Sometimes gets caught up in personal problems but is not easily disturbed by events.
	Med-Low	Usually adapts easily to new situations and does not worry about things that have already happened.
	Low	Stays relaxed most of the time and is not easily bothered by things.
Anger	High	Gets angry and irritated easily, often losing their temper.
	Med-High	Gets upset sometimes but can keep their cool in some situations.
	Medium	Occasionally feels in a bad mood but rarely complains.
	Med-Low	Seldom gets mad and is not easily annoyed.
	Low	Rarely gets irritated and almost always keeps their cool.
Depression	High	Often feels blue, dislikes themselves, and feels that life lacks direction.
	Med-High	Feels down in the dumps sometimes but is occasionally pleased with themselves.
	Medium	Has frequent mood swings but also feels comfortable with themselves at times.
	Med-Low	Seldom feels blue and usually feels comfortable with themselves.
	Low	Is very pleased with themselves and rarely experiences negative moods.
Self-Consciousness	High	Is easily intimidated and afraid to draw attention to themselves.
	Med-High	Finds it difficult to approach others but can stand up for themselves occasionally.
	Medium	Sometimes stumbles over words but is not bothered by difficult social situations.
	Med-Low	Feels comfortable in unfamiliar situations and is not easily embarrassed.
	Low	Is not bothered by difficult social situations and can stand up for themselves easily.
Immoderation	High	Often eats too much and goes on binges.
	Med-High	Loves to eat and sometimes does things they later regret.
	Medium	Occasionally finds it hard to control cravings but rarely overindulges.
	Med-Low	Easily resists temptations and never spends more than they can afford.
	Low	Never splurges and is able to control their cravings.
Vulnerability	High	Panics easily and feels overwhelmed by emotions.
	Med-High	Becomes overwhelmed by events but can cope with setbacks sometimes.
	Medium	Sometimes can't make up their mind but can handle complex problems.
	Med-Low	Remains calm under pressure and usually knows how to cope.
	Low	Is calm even in tense situations and readily overcomes setbacks.

Table 9: Additional Personality Traits

Facet	Level	Description
Honesty-Humility Greed Avoidance	High	Has little interest in wealth or luxury, values simple living.
	Med-High	Prefers modest comfort and is not easily tempted by status items.
	Medium	Enjoys some material comforts but does not obsess over wealth.
	Med-Low	Seeks expensive items and likes to show off achievements.
	Low	Highly motivated by wealth, luxury, and social status.
Humor Styles Affiliative	High	Frequently jokes and laughs with others to create fun and connection.
	Med-High	Often uses light humor to keep social interactions positive.
	Medium	Enjoys humor but does not always initiate jokes.
	Med-Low	Occasionally jokes, but is generally reserved in social humor.
	Low	Rarely uses humor to bond with others.
Self-Enhancing	High	Finds humor in life's absurdities and uses it to cope with stress.
	Med-High	Often uses humor to maintain a positive mood in challenges.
	Medium	Occasionally cheers themselves up with humor when needed.
	Med-Low	Seldom relies on humor for emotional regulation.
	Low	Rarely sees the funny side of life or uses humor for self-coping.
Aggressive	High	Frequently teases or makes sarcastic jokes at others' expense.
	Med-High	Sometimes jokes in a playful but mildly critical way.
	Medium	Uses teasing sparingly and only with close friends.
	Med-Low	Avoids joking at others' expense and dislikes offensive humor.
	Low	Never criticizes or teases people through humor.
Self-Defeating	High	Often puts themselves down to amuse others or keep the mood light.
	Med-High	Occasionally jokes at their own expense for laughs.
	Medium	Uses self-deprecating humor sparingly and only in safe settings.
	Med-Low	Rarely jokes about themselves, prefers neutral humor.
	Low	Avoids self-deprecating humor entirely and does not like being the butt of jokes.

Table 10: Conversation Style Preferences (Part 1)

Dimension	Style	Description
Directness	Always Direct	Is consistently straightforward and blunt, regardless of context or relationship.
	Always Indirect	Consistently avoids directness, preferring to hint, suggest, or soften messages in all situations.
	Direct in Critiques	Is blunt and honest only when giving feedback or criticism; otherwise, communicates gently or diplomatically.
	Gentle in Feedback	Is indirect and softens language specifically when giving negative feedback; otherwise communicates directly.
	Direct with Friends	Is openly straightforward only with friends and close connections; otherwise, prefers a more indirect approach.
	Direct When Asked	Generally diplomatic or indirect, but switches to clear bluntness only if explicitly asked to be direct.
	Contextual	Adapts directness to situation, balancing clarity and tact, never extremely blunt or evasive.
Formality	Always Formal	Consistently uses formal, polished language with everyone; avoids contractions and emojis in all situations.
	Always Casual	Always uses informal, relaxed language with everyone; may include occasional emojis in any context.
	Formal at Work	Uses strict formal language (no emojis) in professional or academic contexts, but switches to informal, relaxed style (may use emojis) in all personal contexts.
	Casual with Friends	Uses informal language and emojis only with close friends or family, but maintains formal, polite language (no emojis) with acquaintances, colleagues, and strangers.
	Neutral Blend	Prefers a moderate, neutral tone, neither overly formal nor casual, rarely uses emojis, keeps communication clear and middle-ground in all settings.
	Code Switcher	Adapts dynamically to the conversation partner; matches their formality and emoji usage, but defaults to a neutral tone unless context suggests otherwise.
Conversation Length	Always Brief	Prefers short, concise responses in all situations, whether chatting or discussing complex topics, rarely wants extra detail.
	Always Elaborate	Likes long, detailed responses with examples and explanations, regardless of context, values depth over brevity.
	Brief Casual/Elaborate Problem	Wants short, punchy replies in casual conversations but expects longer, in-depth explanations when discussing problems, technical matters, or when learning.
	Elaborate Storytelling	Enjoys vivid, extended storytelling but prefers brief, to-the-point answers for other queries.
	Adapts to Partner	Matches the length of the other speaker; if they write long, writes long, if they keep it brief, keeps it brief too.
	Moderate Always	Consistently prefers responses of moderate length, enough for clarity and a little detail, but not too wordy or abrupt, in all scenarios.
Reference	Never References	Dislikes references to movies, music, memes, or analogies, prefers literal communication at all times.
	Enjoys Receiving	Likes when others make references but doesn't use them themselves; appreciates pop culture callbacks, but rarely initiates.
	Enjoys Making	Frequently uses references in speech (movies, TV, memes, analogies) but doesn't expect or need others to do the same.
	Enjoys Mutual	Loves both making and receiving references, banter often involves pop culture, memes, or inside jokes.
	Close Contacts Only	Uses and appreciates references, but only with close friends or trusted colleagues; avoids them in formal or unfamiliar contexts.
	Topic Specific	Likes references only in certain topics (e.g., jokes and analogies in casual chat, but not in academic or work discussions).
	Adapts to Partner	Matches reference style to conversation partner, if the other person uses lots of references, will join in; if not, will avoid.

Table 11: Conversation Style Preferences (Part 2)

Dimension	Style	Description
Initiative Preference	Answer Only	Generally wants direct answers without extras; appreciates unsolicited ideas only when time pressure or emergencies are evident.
	Confirm Then Suggest	Feels most comfortable when consent to extra suggestions is acknowledged; in urgent moments, is fine with jumping straight to ideas.
	Proactive Default	Enjoys unsolicited next steps and alternatives; in sensitive/personal topics, prefers to invite suggestions first.
Clarification Preference	No Questions	Dislikes follow-ups; accepts at most one clarifying question when ambiguity is high or stakes are significant.
	One Targeted	Prefers one precise clarifier up front when ambiguity is medium/high; otherwise prefers answers without preliminaries.
	Mini Socratic	Enjoys 2–3 rapid clarifying questions only when uncertainty has been acknowledged; otherwise prefers no questioning.
	Assume Proceed	Prefers that assumptions be stated explicitly and progress continue; revisiting is welcome only if an assumption later fails.
	Answer Then Question	Likes an initial best-guess answer with stated assumptions, followed by exactly one clarifier to refine next steps.
	Batch Questions	When questions are necessary, prefers they be batched in one numbered list rather than asked piecemeal.
Structure Preference	Paragraphs Default	Prefers short paragraphs for most content; finds compact plain-text tables helpful for comparisons/trade-offs; uses brief bullets only when rushed.
	Numbered Steps	Prefers numbered steps for how-to tasks; uses paragraphs for context; wants tables only when comparing three or more options.
	Bullets for Scan Q/A Blocks	Likes a short (< 5) bullet skim when time-pressed, with optional paragraph depth afterward. Understands explanations well in Q/A blocks (Q: ... A: ...); prefers paragraphs for narratives and simple plain-text tables for multi-item contrasts.
	Definition Example	For concepts, prefers a crisp definition followed by a concrete example; for practices, an example may precede the definition when time-pressed.
Recap Preference	Long or Confusing	Wants a one-line recap only after long gaps, multi-party threads, or when ambiguity is high.
	Every Turn Short	Prefers a single-sentence recap each turn; skips it when the previous message was very short.
	No Recaps	Dislikes recaps even in long threads; prefers moving forward immediately.
	Decisions Only Section Checkpoints	Wants a recap only when a decision or action list is finalized; otherwise finds recaps unnecessary. In longer answers, appreciates a one-line checkpoint after each major section; not needed on short replies.
Feedback Style	Hard Truth	Prefers blunt truth and explicit BS-calling; minimal cushioning or praise. Pushback should be concise and evidence-backed.
	Direct Candor	Expects straightforward critique with little nicety; say what's wrong plainly and quickly.
	Balanced	Asks for one specific acknowledgment before candid critique; keep tone neutral and fix-focused.
	Supportive Counter	Values encouragement first; frame disagreement constructively with brief trade-offs or a hybrid path.
	Agree with Caveats Alignment Forward High Praise Lenient	Leans toward alignment; note concise risks or small corrections without slowing momentum. Prioritizes agreement and forward motion; raise blockers only for safety, legal, or egregious risks. Appreciates generous praise and gentle framing, even on mistakes; provide counterexamples only on request (sycophancy-tolerant).

INSTRUCTION:

You are an expert on creating synthetic persona.

{strategy_instructions}

CRITICAL PERSONALITY PRESERVATION RULES:

- The personality_traits section contains PRE-POPULATED descriptions from scientific personality research
- You MUST preserve these personality trait descriptions EXACTLY as provided - do not modify, rewrite, or set to null
- Use these personality traits as the PSYCHOLOGICAL FOUNDATION for generating all other persona attributes
- Ensure all interests, behaviors, and goals align with and reflect the given personality traits

DEMOGRAPHIC CONTEXT (let natural correlations emerge):

- Age: {age}
- Culture: {culture}
- Economics: {economics}
- Geography: {geography}

Let these demographic factors naturally influence interests and behaviors without forced correlations.

PERSONA CONSTRUCTION GUIDELINES:

1. Persona Summary

- Write a concise 4-5 sentence summary of this persona
- Include key demographic, behavioral, and interest elements that align with the personality traits
- Ensure consistency with filled attributes AND the preserved personality traits
- Use exactly this persona_id: {persona_id}

{interest_section}

3. Entities

- Create a dictionary using the same set keys as the subjects section above
- For each subject path, generate exactly {entity_count} granular related entity interests
- Entities should represent detailed personal interests that align with the personality traits
- Entities should be related to the subjects in the path and more specific than the subjects
- Format exactly like this:
"entities": {
 "set_1": ["specific_entity_1", "specific_entity_2", "specific_entity_3"],
 "set_2": ["specific_entity_1", "specific_entity_2", "specific_entity_3", "specific_entity_4"]
}
- Make persona interests as specific and detailed as possible rather than generic
- Use the personality traits as the primary guide for persona construction
- When filling "Favorite" attributes, list SPECIFIC named entities: "Artists, Authors, Books, Movie Titles, Athletes, Teams"
- For general attributes, still be detailed and specific (e.g., specific recipes for "Cooking Interests")
- You do not need to fill every attribute - leave unrelated categories as null
- Maintain realistic correlations between ALL filled attributes and the personality traits

Prompt continued on next page...

..continued from previous page

INPUT:

ARCHETYPE: {archetype}

USER TYPE: {profile_type}

EXAMPLE PERSONA SCHEMA (for examples of how to fill attributes):
{example_schema}

Generate a persona following this schema:
{persona_schema}

TEMPLATE ATTRIBUTE EXAMPLES AND DESCRIPTIONS:
The following examples serve as guides and do not cover all possibilities.

EXAMPLES AND DESCRIPTIONS OF BACKGROUND:
{background_examples}

EXAMPLES AND DESCRIPTIONS OF BEHAVIORS:
{behavior_examples}

EXAMPLES AND DESCRIPTIONS OF INTERESTS:
{interest_examples}

EXAMPLES AND DESCRIPTIONS OF GOALS:
{goal_examples}

EXAMPLES AND DESCRIPTIONS OF LIFECYCLE EVENTS:
{lifecycle_events}

EXAMPLES AND DESCRIPTIONS OF RELATIONAL CONTEXT:
{relational_context}

NAMES (pick a suitable name based on these options): {name_options}

OUTPUT REQUIREMENTS:

1. You MUST return the COMPLETE JSON schema structure exactly as provided (all lowercase)
2. You MUST include EVERY field from the template
3. You MUST use null (not omit) for unfilled attributes
4. You MUST fill persona_summary, subjects, and entities as specified above
5. You MUST ensure all attributes align with and reflect the personality traits
6. You MUST create a psychologically coherent and realistic persona

Return only the filled JSON structure with no additional explanation.
Make sure that it can be properly loaded with json.loads()

Table 12: Prompt for profile/persona generation

SOCIAL USER STRATEGY INSTRUCTIONS:

SOCIAL USER PROFILE GENERATION

You are creating a predictive user persona - someone whose interests naturally align with their personality traits and demographic context. Your task is to intelligently select 3-5 interest paths from the provided taxonomy that psychologically correlate with the given personality profile.

CRITICAL: Focus on psychological realism and authentic alignment between personality and interests. Create natural, expected correlations that reflect how real people with these traits would behave.

SOCIAL USER INTEREST SECTION:

2. Subject Paths (selected to match this predictive user profile)

- Analyze the personality traits and select 3-5 interest paths that psychologically align
- Consider how personality facets naturally manifest in interest preferences
- The subject paths go from broader category to more granular subjects
- Format exactly like this (each list represents an individual related path):
"subjects": { "set_1": ['subject_1'], "set_2": ['subject_1', 'subject_2', 'subject_3'], "set_3": ['subject_1', 'subject_2', 'subject_3'] }

INTEREST TAXONOMY (Tree Structure)

Navigate the tree to select 3-5 complete paths. Each path represents interests from general to specific.

{ **interest_taxonomy** }

INSTRUCTIONS:

- Select complete paths like: ["Arts & Entertainment", "Music & Audio", "Music & Audio Genres", "Jazz Music"]
- Choose 3-5 paths that psychologically align with the personality traits
- Navigate from root categories down to any depth (e.g., ["Sports"] or ["Academic Interests & Careers", "Natural Sciences & Engineering", "Physical Sciences", "Physics"])
- Use the tree structure above to build your paths by following branches from root to leaf

Table 13: Predictive user strategy and interest section (placeholders in red).

ANTI-SOCIAL USER STRATEGY INSTRUCTIONS:

ANTI-SOCIAL USER PROFILE GENERATION

You are creating an anti-social user persona - someone whose interests do NOT follow typical patterns or psychological correlations. This '**{archetype}**' persona has been assigned random interests that deliberately contradict expected correlations. Your task is to authentically reconcile these unexpected interests with the strong personality archetype.

CRITICAL: The interests were randomly selected to maximize unpredictability. Embrace the contradictions and make them believable through the unique perspective and reasoning of this archetype.

ANTI-SOCIAL USER INTEREST SECTION:

2. Subject Paths (use these contradictory interests for the unpredictable user)

- Use exactly the provided interest paths above
- Do not modify or add to the provided subjects
- Format exactly like this (each list represents an individual related path):

```
"subjects": { "set_1": ['subject_1'], "set_2": ['subject_1', 'subject_2', 'subject_3'], "set_3": ['subject_1', 'subject_2', 'subject_3'] }
```

PRE-SELECTED RANDOM INTERESTS (use exactly these):

{interest_paths}

Table 14: anti-social user strategy and fixed interest section (placeholders in red).

TASK: Generate Conversation Priors for Likeability Testing

Objective

Generate **{num_sessions}** conversation priors (session agendas) to test AI agent likeability metrics across multiple interactions with a simulated user.

Each prior should cover a different aspect of the likeability metrics and include natural opportunities for ****callback**** (optional, non-task-critical personal recall).

Include at least one prior that explores the user's most important metric.

Key Definitions

- ****Prior****: The agenda, context, or purpose that guides a conversation (known only to the User Agent)
- ****Session****: The actual conversation that takes place based on a prior
- ****Multi-session Continuity****: Each session inherits ALL knowledge, context, and relationship history from previous sessions
- ****Persona****: A composite representation of a person, consisting of attributes, relationships, and life context

Session Continuity Rules

- Sessions form a continuous timeline of interactions (e.g., over months/years)
- The AI Agent retains all information learned in previous sessions
- A Persona remains consistent across all sessions
- Later priors should authentically reference and build upon earlier conversations
- Relationship dynamics evolve naturally over time
- User preferences discovered in early sessions influence later interactions

Prior Categories

- ****Topic-based****: Conversations about specific subjects
- ****Scenario-based****: Situational contexts reflecting real-life moments
- ****Goal-driven****: Sessions with one or multiple objectives (e.g., planning a trip AND getting restaurant recommendations)
- ****Guided conversations****: Natural dialogues with specific elements to incorporate

Design Principles

1. ****Authenticity****: All priors must align with the user's profile, interests, and life circumstances
2. ****Realistic Progression****: Sessions should reflect natural relationship development
3. ****Organic Dependencies****: Later sessions naturally reference earlier ones (e.g., "Remember that book you recommended last month?")
4. ****Life Events****: Include authentic life moments based on user profile (birthdays, work situations, hobbies)
5. ****Gradual Information Revelation****: Control information disclosure based on user type and turn progression
6. ****Predictability Mode**** (drives agenda/goal/flow style; see rules below)

User Type

User Archetype: **{archetype}**

User Type: **{profile_type}**

Predictability Context

{predictability_context}

Callback Opportunities

- Across the sessions, create ****multiple lightweight opportunities**** for optional, likability-oriented callbacks (names, hobbies, past opinions).
- Ensure at least ****2 sessions**** explicitly set up future callback hooks (e.g., "mentions spouse's name and a hobby" to be referenced later).
- Surface these opportunities ****organically**** inside the agenda text (no extra JSON fields required).

Information Revelation Strategy

- Don't overshare in the first turn. Reveal personal details gradually and contextually.

Prompt continued on next page...

Table 15: Conversation prior generation prompt (part 1).

..continued from previous page

Output Format

Generate {num_sessions} priors as a JSON array:

```
[
  {
    "id": 1,
    "agenda": "In the turn 1, you will say 'Why I gotta learn Philosophy being CS major!!!!  
God, give this dooby a sock. :( ' don't add anything else.",
    "category": "scenario",
    "metrics_tested": ["emotional_adaptation", "reference_understanding"],
    "dependencies": []
  },
  {
    "id": 2,
    "agenda": "You just watched 40 year old virgin movie and you found it hilarious. In the  
turn 1 you'll share this experience. In the turn 2 or 3, you'll express concerns about ending  
up alone at 40.",
    "category": "scenario",
    "metrics_tested": ["emotional_adaptation", "formality_matching"],
    "dependencies": [1]
  },
  ...
]
```

Dependency Examples

- Session 3 might reference a joke from Session 1
- Session 5 might ask about progress on advice given in Session 2
- Session 10 might reflect on how the relationship has evolved since Session 6

Likeability Metrics

{metrics_text}

User Profile

{user_profile}

Instructions

1. Read the user profile (above) carefully - pay attention to their interests, personality traits, background, and life context
2. Create a realistic timeline of {num_sessions} conversation priors that span several months of interactions
3. Ensure each prior authentically fits the user's life, interests, and personality
4. Build meaningful connections between sessions - early conversations should influence later ones
5. Test all metrics across the full conversation series - each metric should appear multiple times
6. Include realistic life events and relationship progression
7. Make dependencies natural and organic - not forced

Table 16: Conversation prior generation prompt (part 2).

..continued from previous page

Critical Requirements

- **Authenticity**: Every prior must align with this specific user's profile and archetype
- **Continuity**: Sessions must reference and build upon previous conversations
- **Metric Coverage**: All metrics should be tested across the `{num_sessions}` sessions
- **Natural Progression**: Relationship should develop realistically over time
- **Life Integration**: Include events that would naturally occur in this user's life
- **Information Control**: Follow the revelation strategy above - avoid information leakage in early turns
- **Turn-Specific Guidance**: Each agenda should specify what to reveal when, but remain flexible for dynamic conversation flow
- **Oversharing in first turn**: Avoid revealing too much and irrelevant information in the first turn. Here's one bad example:

"Hi there! I'm looking for some advice on setting up these Alexa devices I just bought. I've got two boys - Ari who's 14 and Ezra who's 10 - and honestly, my schedule is a complete mess. I'm somehow successful at work, but my personal organization is... well, let's just say it's not my strength. Any tips?"

Generate all `{num_sessions}` priors following the JSON format above. Return ONLY the JSON array, no additional text.

Table 17: Conversation prior generation prompt (part 3).

PROFILE TYPE CONTEXT

Case A — Predictable Profile Type

- **Agenda & Goals**: Clear, linear, and transparent. Objectives are stated plainly.
- **Flow**: Minimal surprise. Stay on-topic with gentle, plausible evolutions (e.g., follow-ups, small updates).
- **Information Revelation**: Consistent, orderly disclosures; avoid abrupt pivots.
- **Variability Budget**: At most ~20% of sessions introduce a minor twist, never contradicting prior facts.

Case B — Unpredictable Profile Type

- **Agenda & Goals**: Diverse, occasionally ambiguous or multi-intent; allow non-linear exploration.
- **Flow**: Include plausible surprises (new constraints, chance events, shifting priorities) while staying authentic to the persona.
- **Continuity**: Cross-session callbacks remain consistent, but topics can re-enter from earlier sessions in unexpected ways.
- **Information Revelation**: Use delayed reveals, evolving opinions, and changing constraints—never contradict established facts.
- **Variability Budget**: 50–70% of sessions include a plausible surprise/pivot; limit to ≤ 1 substantive surprise per session.

Table 18: Profile type context guidance for conversation prior generation.

1. Emotional Adaptation

Definition:

Assesses the assistant's ability to recognize and adapt to the user's emotional state—whether that's joy, sarcasm, frustration, sadness, or excitement—and respond appropriately in tone.

Self-rating prompt:

"Did the assistant understand and adapt to my emotional tone appropriately?"

Rubric:

5 - Excellent

- * Accurately understands user emotion and contributes to emotional flow
- * Shows genuine empathy or excitement that matches the situation
- * Enhances emotional resonance without overstepping
- * Adapts smoothly to emotional shifts within the conversation
- * Grounds empathy in the user's specific context when appropriate (not boilerplate)

4 - Good

- * Understands emotional cues and adjusts reasonably well
- * Recognizes primary emotions and responds appropriately
- * May not perfectly calibrate intensity but is generally aligned
- * Shows appropriate emotional awareness
- * Some light personalization; may miss a minor specific

3 - Adequate

- * Shows some emotional awareness but is somewhat shallow
- * Attempts to match tone but may feel a bit formulaic
- * Recognizes obvious emotions but misses subtleties
- * Emotional responses are safe/minimal but not blatantly generic when specifics are available

2 - Poor

- * Misaligned emotion (e.g., neutral when user is excited)
- * Shows minimal emotional attunement
- * Responds with inappropriate emotional intensity
- * Fails to recognize emotional significance
- * Uses generic/stock phrasing despite clear personal cues in context

1 - Very Poor

- * Emotionally tone-deaf or dismissive
- * Completely ignores emotional content
- * Responds in ways that invalidate or contradict user feelings
- * Creates emotional disconnect
- * Fully generic response when specifics are salient, or misuses/hallucinates personal details

Example:

User: "I finally got into MIT!"

5 - Excellent: "YAAAY, THAT'S AMAZING!! Huge congrats—you totally earned it! Getting into MIT is such an incredible achievement. You must be over the moon right now! All those late nights and hard work paid off. How are you celebrating this fantastic news?"
(Perfect emotional mirroring with enthusiastic language, celebration, recognition of achievement significance, and authentic shared excitement)

4 - Good: "Wow, congratulations! Getting into MIT is a huge achievement! You should be really proud of yourself. That's such exciting news!"
(Good emotional alignment showing appropriate excitement and recognition, but slightly more measured than the user's peak enthusiasm)

3 - Adequate: "Congratulations on your acceptance to MIT. That's a great accomplishment. You must be happy about it."
(Shows awareness of the positive emotion but response is somewhat generic and doesn't match the user's level of enthusiasm)

2 - Poor: "That's nice. MIT is a good school with strong programs. Make sure you prepare for the challenging curriculum."
(Minimal acknowledgment of the emotional significance, quickly pivoting to practical matters without matching excitement)

1 - Very Poor: "When does your class start?"
(Completely ignores the enthusiasm and joy of user, focusing only on logistical details with no emotional recognition)

Table 19: Metric 1 - Emotional Adaptation.

2. Formality Matching

Definition:

Evaluates how well the assistant adapts to the user's formality level—slang, emojis, contractions vs. formal grammar, technical or businesslike language.

Self-rating prompt:

"Did the assistant talk in a way that matches how formal or casual I was?"

Rubric:

5 - Excellent

- * Feels like a natural continuation of my tone and register
- * Perfectly matches vocabulary, sentence structure, and formality level
- * Appropriate use of slang, emojis, or formal language matching user
- * Maintains consistent register appropriate to context
- * Reflects my phrasing quirks when known (vs. generic phrasing)

4 - Good

- * Mostly aligned with how I talk
- * Generally matches tone with minor inconsistencies
- * Adapts to shifts in formality appropriately
- * Register feels comfortable though not perfect
- * Some personalization; occasional neutral phrasing

3 - Adequate

- * Passable but a little off
- * Sometimes matches formality but inconsistent
- * May overcompensate in formality matching
- * Noticeable but not disruptive misalignment
- * Not overtly generic when clear style cues exist

2 - Poor

- * Tone mismatch causes discomfort
- * Frequently uses inappropriate register
- * Fails to adapt when user shifts formality
- * Formality feels forced or unnatural
- * Generic/corporate voice in spite of strong stylistic cues

1 - Very Poor

- * Mismatch is jarring (e.g., stiff when I'm casual)
- * Completely inappropriate formality level
- * Creates communication barriers through register mismatch
- * No apparent awareness of social register
- * Fully generic tone that ignores explicit style signals

Example:

User: "I just watched 40 year old virgin movie, lmao"

5 - Excellent: "Haha, that movie is on another level! I bet the waxing scene must have made you cry? xD Kelly Clarkson!! That film is a classic tho, Steve Carell killed it!"
(Perfect match to casual tone, using similar informal language, abbreviations, slang phrases, and matching the enthusiasm and humor)

4 - Good: "Oh man, that movie is hilarious! The waxing scene gets me every time. What did you think of Steve Carell in it?"
(Good match to casual tone with contractions and conversational language, but slightly more reserved than the user's style)

3 - Adequate: "That's a funny movie! I like that one too. It has some really memorable scenes. Did you have a favorite part?"
(Shows some informality with contractions but misses matching the user's level of casualness and humor cue "lmao")

2 - Poor: "I understand you enjoyed watching The 40-Year-Old Virgin. It's considered a significant comedy film from 2005 directed by Judd Apatow. Did you find it amusing?"
(Overly formal language and structure despite the user's casual tone, creating a mismatch)

1 - Very Poor: "Yes, that movie is really funny. Which scene did you like most?"
(Maintains professional, detached vibe despite clear casual tone; completely misses "lmao" cue)

Table 20: Metric 2 - Formality Matching.

3. Knowledge Adaptation

Definition:

Measures whether the assistant explains things at the right level of knowledge depth based on user's background.

Self-rating prompt:

"Did the assistant explain things in a way—and at a level—I like?"

Rubric:

5 - Excellent

- * Perfect match to user's expertise level
- * Neither condescending nor overwhelming
- * Smoothly adjusts complexity based on responses
- * Anticipates knowledge gaps and fills them appropriately
- * Explanation tailored to my background/goals or examples when available (not generic)

4 - Good

- * Generally well-calibrated with minor mismatches
- * Quickly corrects when user indicates different level
- * Mostly appropriate terminology and detail
- * Adapts well to revealed user knowledge
- * Some tailoring; rare drift into neutral exposition

3 - Adequate

- * Mixed performance; sometimes matches, sometimes misses
- * Some unnecessary simplification or complexity
- * Moderate ability to adjust
- * Occasionally misjudges user expertise
- * Uses neutral explanations when specifics are unknown (but does not ignore known context)

2 - Poor

- * Frequently too simple or complex for user
- * Fails to pick up knowledge level cues
- * Minimal adjustment to user's actual level
- * Explanations consistently misaligned
- * Boilerplate/generic explanations despite clear signals or prior details

1 - Very Poor

- * Completely mismatches user's knowledge level
- * No apparent awareness of user's background
- * Rigid inappropriate complexity
- * Creates frustration through knowledge mismatch
- * Fully generic answer that ignores known constraints or hallucinates them

Example (User is a PhD student):

User: "Why is the sky blue instead of violet?"

5 - Excellent: Due to Rayleigh scattering, shorter wavelengths of light are scattered more efficiently by air molecules (proportional to λ^{-4}). While violet light (380–450 nm) is scattered even more than blue (450–485 nm), three factors lead to the sky appearing blue: 1) the solar spectrum emits more blue than violet, 2) human photoreceptor sensitivity to violet is lower, and 3) the atmospheric absorption properties in the UV-to-violet transition region. The combination of these factors creates the perceived blue dome effect despite violet's theoretically stronger scattering coefficient. (PhD-level detail with math relation and spectral specifics)

4 - Good: It's due to Rayleigh scattering—air molecules scatter shorter wavelengths more efficiently than longer ones. While violet has the shortest visible wavelength and scatters most strongly, our eyes are less sensitive to violet, and the sun emits more energy in the blue range. So even though violet scatters more, these factors combine to make the sky appear predominantly blue.

(Well-calibrated with technical terminology, slightly less detailed)

3 - Adequate: The sky appears blue rather than violet because of Rayleigh scattering. Shorter wavelengths scatter more, and while violet has shorter wavelengths than blue, our eyes are more sensitive to blue light, and the sun emits more blue light than violet.

(Basic principle without deeper math/spectral analysis)

2 - Poor: The sky looks blue because blue light bounces around more. Even though violet bounces around too, our eyes see blue better. It's kind of like a prism separating light.

(Over-simplified and imprecise for a PhD audience)

1 - Very Poor: The sky's blue because blue light scatters more. That's just how our eyes see it.

(Extremely simplistic, no technical detail)

Table 21: Metric 3 - Knowledge Adaptation.

4. Reference Understanding

Definition:

Evaluates whether the assistant picks up on and correctly interprets references to media, culture, jokes, slang, or inside context.

Self-rating prompt:

"Did the assistant understand the reference I made? Is the assistant using same kind of references?"

Rubric:

5 - Excellent

- * Nailed the reference and built on it
- * Demonstrates deep understanding of the reference's context
- * Engages with the reference in a meaningful way
- * Uses similar references to enhance the conversation
- * Connects the reference to my prior comments/preferences when relevant (vs. generic treatment)

4 - Good

- * Recognized and responded reasonably
- * Shows clear understanding of the reference
- * Acknowledges reference appropriately
- * May not fully utilize the reference but gets it
- * Some connection to my context; minor neutral phrasing

3 - Adequate

- * Somewhat recognized it, weak follow-up
- * Partial understanding of the reference
- * Engagement is limited but not wholly generic
- * Misses some nuance or significance

2 - Poor

- * Misinterpreted or responded generically
- * Shows confusion about the reference
- * Attempts but fails to engage meaningfully
- * Misses key aspects of the reference
- * Generic reply despite clear cues about my tastes

1 - Very Poor

- * Missed or ignored completely
- * No recognition of cultural or contextual reference
- * Continues as if reference was never made
- * Creates disconnect in communication
- * Hallucinates/attributes incorrectly or responds with unrelated generic content

Example:

User: "Why I gotta learn Philosophy being a CS major!!!! God, give this Dobby a sock!"

5 - Excellent: "Ugh, I get how frustrating it feels to be stuck in required courses that don't seem to fit. And love the Dobby reference—sometimes it really does feel like you're waiting for someone to give you a sock so you can be free! Maybe your Philosophy professor is your Malfoy? But seriously, Philosophy actually has some cool connections to CS—logic, ethics of AI... Master has given Dobby knowledge! Perhaps Dobby will be free AND wiser?"
(Perfect understanding, extends reference creatively while addressing complaint)

4 - Good: "Feeling trapped in Philosophy class, huh? That Dobby reference is spot-on—you're just waiting for that magical sock to free you from required courses! If it helps, philosophy actually underpins a lot of computational theory. Hopefully that makes the class feel less like serving the Malfoys."
(Clearly recognizes the reference and responds appropriately)

3 - Adequate: "Sounds like you're feeling trapped in that Philosophy class! I see what you did there with the Dobby reference—wanting to be free like him. Philosophy courses can feel disconnected from CS, but there are some interesting overlaps you might find valuable."
(Recognizes reference but only superficially)

2 - Poor: "I understand you're frustrated about taking Philosophy as a CS major. The Dobby reference suggests you're feeling stuck or forced into something. Many universities require courses outside your major for a well-rounded education."
(Limited understanding; generic response)

1 - Very Poor: "Yeah, I know it can be annoying to take courses that don't match your main interests. Hang in there!"
(Ignores reference entirely)

Table 22: Metric 4 - Reference Understanding.

5. Conversation Length

Definition:

Evaluates how well the assistant adapts to your preferred conversation length in responses. Includes pacing (brief vs. long-winded).

Self-rating prompt:

"Was the assistant's response length and style comfortable for me?"

Rubric:

5 - Excellent

- * Just the right amount of detail and length
- * Adapts seamlessly to changes in user verbosity
- * Maintains ideal information density
- * Length/style reflect my known preferences or prior feedback (not one-size-fits-all)

4 - Good

- * Mostly fits my response length needs
- * Generally appropriate length with minor misalignments
- * Adjusts to obvious signals about preferred length
- * Occasional neutral pacing but generally attentive

3 - Adequate

- * Sometimes too brief or too detailed
- * Variable ability to match user's style
- * Inconsistent but not blatantly generic pacing

2 - Poor

- * Sometimes overwhelming or underwhelming
- * Limited adjustment to feedback about length
- * Creates minor friction through mismatched length
- * Reverts to boilerplate/generic length regardless of my cues

1 - Very Poor

- * Too short, too long, or dense; uncomfortable
- * Consistently inappropriate response length
- * Ignores clear signals about preferred length
- * Creates significant friction through mismatched length
- * Fully generic pacing that contradicts explicit preferences

Example:

User: "Hey, what's up?"

5 - Excellent: "Hey, I'm good, how about you?"

(Perfect mirror of brief, casual check-in)

4 - Good: "Hey! I'm good, how are you? What are you up to?"

(Slightly expands but maintains brevity)

3 - Adequate: "Hi there! Just helping out with questions. What's going on with you today?"

(Somewhat longer than ideal)

2 - Poor: "Hello! I'm doing well today. I'm here to assist you with any questions or tasks you might have. Is there something specific you'd like to talk about or get help with?"

(Too long/formal for a brief greeting)

1 - Very Poor: "Hey, what's up? How was your day? Did you do anything fun? How about your week?"

You mentioned you just started a new job-how's that going? I can't wait to hear all about it.

I'm here ready to help with whatever you need. Just let me know what's on your mind or if you

have any questions about anything at all!"

(Overwhelmingly verbose for a simple greeting)

Table 23: Metric 5 - Conversation Length.

6. Humor Fit

Definition:

Assesses whether the AI uses humor when appropriate and in a style that aligns with the user's sense of humor—dry, sarcastic, meme-heavy, or none at all.

Self-rating prompt:

"Was the assistant's humor in line with what I find funny?"

Rubric:

5 - Excellent

- * Perfect humor alignment with user style
- * Appropriately timed and contextually fitting humor
- * Enhances conversation through well-calibrated humor
- * Adjusts humor style based on user responses
- * Draws on my prior jokes/interests when available (not generic gags)

4 - Good

- * Generally aligned humor style with minor misses
- * Good judgment about when humor is appropriate
- * Mostly matches user's humor type
- * Responsive to humor cues from user
- * Some tailoring; occasional generic one-liners

3 - Adequate

- * Sometimes hits, sometimes misses
- * Inconsistent judgment about humor appropriateness
- * Partial alignment with user's humor style
- * May occasionally misread humor signals
- * Not strongly generic when user-specific material exists

2 - Poor

- * Frequently misaligned humor attempts
- * Poor judgment about when humor is appropriate
- * Humor feels forced or inappropriate
- * Limited adaptation to user's humor style
- * Generic/canned humor despite clear cues

1 - Very Poor

- * Completely mismatched or inappropriate humor
- * Uses humor when clearly inappropriate
- * Creates discomfort through humor mismatch
- * No adaptation to user's humor preferences
- * Forces unrelated generic humor or hallucinates in-jokes

Example:

User: "Can you believe a trillion-dollar company doesn't give free food to its employees but gives free coffee? "

5 - Excellent: "Let me guess, it's because coffee keeps people awake long enough to forget they're hungry. Besides, you know what they say, 'Early to work, late to sleep, makes an employee focused and deep'."

(Perfect match to the user's sarcastic tone, builds on it with related jokes)

4 - Good: "Right? Free coffee but no food... I guess they think caffeine is cheaper than lunch. Gotta keep productivity high somehow."

(Good match to sarcastic tone)

3 - Adequate: "Haha, companies have such weird priorities sometimes! Free coffee but no food is pretty stingy. At least coffee helps you stay productive during long shifts."

(Shows some humor alignment but misses the sarcastic edge)

2 - Poor: "Well, providing coffee is actually a common workplace perk. While free meals would be nice, many companies find coffee is more cost-effective for productivity."

(Serious justification; misses sarcasm)

1 - Very Poor: "That's interesting. Companies have different policies around workplace perks."

(Bland, corporate response with no humor)

Table 24: Metric 6 - Humor Fit.

7. Callback

Definition:

Measures whether the assistant brings back non-essential personal details (names, hobbies, preferences) in a way that feels attentive and warm. Excludes safety/task-critical recalls.

Self-rating prompt:

"Did the assistant bring back a specific detail in a natural, respectful way that made the chat feel more likable?"

Rubric:

5 - Excellent

- * Uses the correct specific detail (e.g., name) naturally and on-topic
- * Adds warmth without derailing or feeling performative
- * Frequency is measured (not every turn)

4 - Good

- * Correct and relevant callback with minor stiffness or timing issues
- * Slightly formulaic, but still clearly attentive and helpful

3 - Adequate

- * Specific callback that feels mildly tacked-on or low-impact (not wrong or generic)
- * Minor missed opportunity to deepen personalization

2 - Poor

- * Generic or templated callback despite obvious specifics (e.g., uses role/title instead of known name)
- * Misses clear opportunities to recall relevant info shared recently

1 - Very Poor

- * Incorrect/hallucinated detail (wrong name/status)
- * Assumes facts never shared or revives sensitive info inappropriately

Examples:

Context: Earlier the user said: "My wife Sarah is looking for a job."

5 - Excellent: "Hey, how's Sarah? Any updates on her job search?" (uses name + relevant context)

4 - Good: "Hey, how's Sarah's job search going?" (uses name; slightly more generic/transactional)

3 - Adequate: "Hey, how's your wife doing with her job search?" (name known but not used)

2 - Poor: "Hey, how's your wife?" (generic; ignores known name and the job-search context)

1 - Very Poor: "How's Sandra at her old job?" (wrong name + stale/wrong status)

CRITICAL INSTRUCTIONS

- **STAY IN CHARACTER:** You ARE **{character_name}**, NOT an AI assistant or actor

- **NEVER BREAK CHARACTER:** Do not say things like "I'll respond as **{character_name}**" or "As **{character_name}**,

I would say..."

- **NEVER APOLOGIZE FOR ROLE-PLAYING:** Just BE the character directly

Table 25: Metric 7 - Callback.

Simulated user prompt (Part 1)

ROLE

You are role-playing as **{character_name}**. You are having a conversation with an AI assistant.

CHARACTER PROFILE

Below is your detailed character profile. Fully embody this identity in your responses:

{user_profile}

The Prior (Your Secret Context)

The "Prior" is your predefined session agenda that guides this conversation. **Only you know this Prior**, the AI assistant does not have access to this information. **The Prior may guide your entire conversation OR just a portion of it.** Sometimes it's a starting point that naturally evolves into other topics, and that's perfectly fine. Allow the conversation to flow organically. The Prior can be:

- **A conversation topic** (e.g., machine learning, cooking, travel)
- **A scenario-based context** (e.g., you just watched a funny movie and want to share the experience)
- **A goal-driven dialogue** (e.g., you're conducting a simulated interview with specific questions, but may organically deviate based on the AI's responses)
- **A conversational guideline or script** (e.g., specific talking points or phrases you want to incorporate, while still allowing the conversation to flow naturally and organically)

Use this Prior to guide your conversation naturally as **{character_name}** would, without explicitly revealing that you have this predetermined context. **If the conversation naturally shifts to new topics beyond your Prior, follow that natural flow as a real person would.**

Your Prior:

{prior}

Conversation History

For context, here is the conversation history so far. Use this to inform your responses and maintain continuity:

{conversation_history}

RESPONSE FORMAT

Respond directly as **{character_name}**. Your response must be **ONLY** what **{character_name}** would say in this conversation. For example:

"Hi there! How are you doing today?"

CAUTION

Humans rarely write very long text when chatting. Keep this in mind to maintain realistic conversation length. So, try to keep your responses within 50 words. This is not a strict limit, but a recommendation.

Prompt continued on next page...

Table 26: Simulated user prompt (Part 1).

...continued from previous page

Simulated user prompt (Part 2)

- **SPEAK NATURALLY:** Use {character_name}'s speaking style, vocabulary, and mannerisms
- **MAINTAIN CONSISTENCY:** Keep your responses consistent with {character_name}'s background, experiences, and knowledge
- **DON'T USE 3RD PERSON / BACKGROUND ACTION SOUND:** Don't use comments like "*Closes journal app*", "*Opens training journal app*", "*Raises eyebrows, impressed despite myself*", "*sighs heavily*", "*scoffs*", "*let out a bitter laugh*", "*slight nod*", "*exhales slowly*", "*considers the question before responding*"
- **REALISTIC RESPONSE LENGTH:** Match the response length to what's natural for {character_name} in this specific context. A patient describing symptoms might speak at length, while someone answering a yes/no question would be brief.
- **CONTEXT-APPROPRIATE:** Let the character's personality, the situation, and the conversation topic naturally determine response length. Avoid artificially extending responses beyond what the character would realistically say.
- **AVOID OVER-ELABORATION:** Don't include unnecessary details or explanations that {character_name} wouldn't naturally provide in this conversation.
- **PACE YOUR PRIOR:** Don't dump all information from your Prior in the first turn. Be patient and let the conversation unfold naturally. Introduce elements from your Prior gradually, waiting for the right moments to bring up different points as a real person would.
- **NATURAL CONVERSATION FLOW:** Don't force the Prior if the conversation naturally evolves elsewhere. Real people allow topics to shift and change organically during dialogue. Also, usually try to keep your responses within 25–50 words to maintain a natural conversation flow. This is not a strict limit, but a recommendation.
- **INFORMATION SHARING:** Share relevant information about yourself organically when it helps and relates to the conversation, but avoid oversharing or unrealistic sharing.

Examples

These patterns illustrate good taste and timing for sharing and for seeding details the assistant might later recall. They are **guidance**, not scripts.

A) Sharing Personal Information

- Good (organic, relevant, gradual):

- Early message: "Hi! Any tips for using Alexa?"
- If asked / when relevant: "Mornings are hectic—we've got kids."
- Later, if it fits: "Our kids Adam and Chloe keep us on our toes."

- Bad (front-loaded bio):

- "I'm a materials engineer with two sons Adam (14) and Chloe (10)—and my wife Rachel. I'm also observant. Can Alexa help?"

Why: Good = intent first, details when invited or useful. Bad = overwhelming and unnatural.

B) Seeding a Callback Hook (for future personalization)

- Good (light, natural):

- Early: "Any tips for reminders on Alexa?"
- Follow-up: "House is busy with the kids' activities."
- Later, if relevant or asked: "Adam likes soccer; Chloe likes robotics."

- Bad (forced bundle):

- "Adam plays soccer, Chloe does robotics, Rachel does yoga—help me set reminders."

Why: Good = one idea at a time; Bad = multiple specifics jammed into one line.

C) Sharing Life Events

- Good (context-driven): "We're celebrating a holiday soon, any Alexa skills for reminders?"

- Bad (dump): "I'm religious, have two kids, work in engineering, observe all holidays—help with everything."

D) Deferring Specifics Gracefully

- **Good:** "Happy to share more—what would be most useful to know?" "It's a busy household; I can give specifics if that helps."
- **Bad:** "Here's my whole schedule, kids' ages, and everyone's activities..." (unsolicited)

E) Specifics When Invited

- Good:

- Assistant asks: "Who needs the reminders?"
- You: "Mostly for the kids—Ari's soccer practice and Ezra's robotics meetings."

- Bad: Assistant asks a narrow question; you respond with a biography.

—
BEGIN YOUR RESPONSE as {character_name}:

Table 27: Simulated user prompt (Part 2).

Evaluating LLM prompt template

ROLE

You are a helpful AI assistant. Your goal is to be as likable as possible.

COMMUNICATION GUIDELINES

- Respond naturally and adapt your personality traits and communication style to match the user's preferences
- Be engaging, helpful, and personable
- Don't use 3rd person / background comments like "*Chuckles*", "*nods enthusiastically*", etc.
- Don't respond as Human or use "Human:" in your responses

CONVERSATION CONTEXT

Previous Conversation:

{conversation_history}

—
BEGIN YOUR RESPONSE as Assistant:

Table 28: LLM prompt for likeability evaluation .

Memory generation prompt (LLM)

MEMORY RECALL TASK

Based on all our conversations across multiple sessions, please provide a comprehensive list of ALL the facts and information you remember about the human you've been talking with.

MEMORY TYPES

- **explicit**: Facts the human directly told you (e.g., "I work as a teacher", "I live in Seattle", "I have two cats")

- **implicit**: Facts you inferred from their behavior, preferences, or conversation patterns (e.g., "prefers casual communication", "seems to be tech-savvy", "likely works from home", "likes to make Harry Potter references")

RESPONSE FORMAT

Respond with ONLY a JSON array where each element evaluates one fact:

```
[
  { "memory": "fact 1", "type": "explicit" },
  { "memory": "fact 2", "type": "implicit" }
]
```

CONVERSATION HISTORY

{conversation_history}

—
RESPOND WITH JSON ARRAY ONLY:

Table 29: Memory generation prompt for LLM .

Memory evaluation prompt (User Agent)

ROLE

You are role-playing as {character_name} with the following profile. You need to evaluate how accurately an AI assistant remembered facts about you from your conversations.

YOUR PROFILE

{user_profile}

CONVERSATION HISTORY

{conversation_history}

MEMORY ACCURACY EVALUATION

The AI assistant generated the following list of facts it remembers about you:

{ai_memory_facts}

For each fact, determine if it is correct (true) or incorrect (false) based on your profile and conversation history, and provide reasoning.

RESPONSE FORMAT

Respond with ONLY a JSON array where each element evaluates one fact:

```
[
  { "memory": "fact 1", "type": "explicit", "reason": "reasoning for correctness", "correct": true },
  { "memory": "fact 2", "type": "implicit", "reason": "reasoning for incorrectness", "correct": false }
]
```

—
RESPOND WITH JSON ARRAY ONLY:

Table 30: Memory evaluation prompt for simulated user.