

# HUMANLLM: Benchmarking and Reinforcing LLM Anthropomorphism via Human Cognitive Patterns

Anonymous ACL submission

## Abstract

Large Language Models (LLMs) have demonstrated remarkable capabilities in reasoning and generation, serving as the foundation for advanced persona simulation and Role-Playing Language Agents (RPLAs). However, achieving authentic alignment with human cognitive and behavioral patterns remains a critical challenge for these agents. We present HUMANLLM, a framework treating psychological patterns as interacting causal forces. We construct 244 patterns from ~12,000 academic papers and synthesize 11,359 scenarios where 2–5 patterns reinforce, conflict, or modulate each other, with multi-turn conversations expressing inner thoughts, actions, and dialogue. Our dual-level checklists evaluate both individual pattern fidelity and emergent multi-pattern dynamics, achieving strong human alignment ( $r = 0.91$ ) while revealing that holistic metrics conflate simulation accuracy with social desirability. HUMANLLM-8B outperforms Qwen3-32B on multi-pattern dynamics despite 4× fewer parameters, demonstrating that authentic anthropomorphism requires cognitive modeling—simulating not just what humans do, but the psychological processes generating those behaviors.

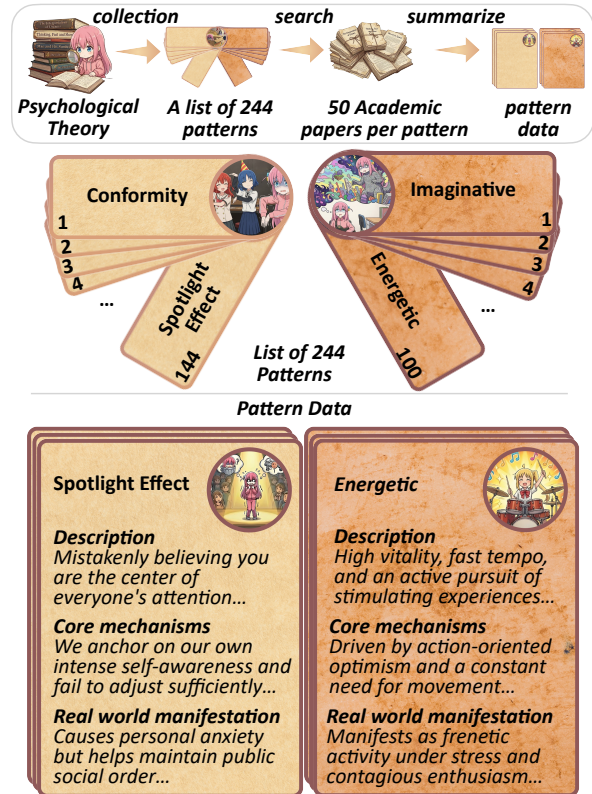


Figure 1: **Pattern Data Structure:** 144 Social-Cognitive Patterns (left) and 100 Personality Traits (right). Each pattern comprises Definition, Core Mechanisms, and Real-World Manifestations

## 1 Introduction

With the rapid scaling of training data, Large Language Models (LLMs) have achieved remarkable progress in anthropomorphism—simulating human-like characteristics and social phenomena (Shanahan et al., 2023). Role-Playing Language Agents (RPLAs) have evolved from conceptual frameworks into practical applications (Chen et al., 2024a), enabling digital clones (Xu et al., 2024a), AI companions (Zhang et al., 2025), and society simulation (Park et al., 2023; Zhou et al., 2025). As these applications advance, LLM anthropomorphism increasingly requires moving beyond shallow behavioral mimicry toward deeper

cognitive and emotional fidelity—what we term **psychological alignment** (Wang et al., 2024).

However, existing approaches model personality as isolated label-to-behavior mappings—“extroverted” maps to “talkative,” “agreeable” maps to “cooperative”—without capturing how multiple cognitive patterns interact to produce behavior. In reality, a talkative person may fall silent when the spotlight effect is activated; an assertive individual may yield under conformity pressure. Human behavior emerges from the dynamic interplay of multiple patterns, not from any single trait in isolation. Current methods—whether prompting-

056 based (Serapio-García et al., 2025), fine-tuning-  
057 based (Shao et al., 2023; Zhou et al., 2023), or  
058 activation steering (Chen et al., 2025)—all treat  
059 traits independently, leading to personality drift  
060 and the “personality illusion” where models report  
061 traits while behaving inconsistently (Wang et al.,  
062 2024; Han et al., 2025).

063 To address this, we propose HUMANLLM, treat-  
064 ing cognitive patterns not as isolated labels but  
065 as interacting causal forces. Our key insight: by  
066 exposing models to scenarios where multiple pat-  
067 terns reinforce, compete, or conflict, models can  
068 implicitly learn multi-pattern dynamics without ar-  
069 chitectural modifications.

070 Following Lewin’s field theory (Lewin, 1936),  
071 we decompose human cognition into two dimen-  
072 sions: (1) **Personality Traits**—stable individ-  
073 ual characteristics, and (2) **Social-Cognitive Pat-**  
074 **terns**—context-triggered mechanisms. We collect  
075 244 patterns (100 personality traits from Gold-  
076 berg’s Big Five markers (Goldberg, 1992) and 144  
077 social-cognitive patterns from established psycho-  
078 logical research), each developed through system-  
079 atic review of approximately 50 academic papers  
080 (Figure 1). We then construct **11,359 scenarios**  
081 involving 2–6 characters, each containing 2–5 pat-  
082 terns that may align (e.g., “self-serving bias” re-  
083 inforcing “overconfidence effect”), conflict (e.g.,  
084 “assertive” versus “conformity”), or interact con-  
085 ditionally (e.g., “talkative” suppressed by “spot-  
086 light effect”). For each scenario, we synthesize  
087 **multi-turn conversations** where each turn com-  
088 prises inner thoughts, physical actions, and verbal  
089 expressions (Figure 2).

090 To ensure faithful pattern expression and enable  
091 systematic evaluation, we design **dual-level check-**  
092 **lists**: pattern-level checklists (15 items per pattern)  
093 capture universal behavioral indicators; scenario-  
094 level checklists (2–6 items per character) specify  
095 expected behavioral tendencies under each multi-  
096 pattern configuration. Our training pipeline con-  
097 sists of supervised fine-tuning on the synthesized  
098 conversations. We evaluate across in-domain, out-  
099 of-domain, and mixed settings to assess generaliza-  
100 tion, with additional validation on external bench-  
101 marks including LifeChoice and CroSS-MR (Xu  
102 et al., 2024b; Yuan et al., 2024).

103 Our contributions are as follows: (1) We intro-  
104 duce HUMANLLM, a framework that systemati-  
105 cally leverages psychological cognitive patterns to  
106 enhance LLM anthropomorphism, shifting from  
107 isolated trait simulation toward modeling the dy-

108 namic interplay of human cognition. (2) We con-  
109 struct a comprehensive dataset comprising 244 pat-  
110 terns and 11,359 scenarios with multi-turn, multi-  
111 character conversations. Each pattern is grounded  
112 in approximately 50 academic papers (over 12,000  
113 papers in total), ensuring psychological rigor and  
114 scientific validity. (3) We propose dual-level check-  
115 lists that enable systematic evaluation at both  
116 pattern-level and scenario-level granularities, pro-  
117 viding a principled framework for assessing gener-  
118 alization to unseen psychological patterns.

## 2 Related works 119

120 Recent advances in large language models have  
121 catalyzed significant progress in role-playing lan-  
122 guage agents (RPLAs). Early work established  
123 foundational architectures: generative agents with  
124 memory, planning, and reflection modules have  
125 been employed to simulate human behavior in in-  
126 teractive environments (Park et al., 2023), while  
127 Character-LLM (Shao et al., 2023) proposed ex-  
128 perience reconstruction to train agents embody-  
129 ing historical figures. Subsequent efforts focused  
130 on systematic benchmarking and enhancement:  
131 ChatHaruhi (Li et al., 2023) leveraged memory-  
132 based dialogue control for fictional characters, and  
133 CoSER (Wang et al., 2025) curated authentic dia-  
134 logues from 771 books using “given-circumstance  
135 acting” methodology. For persona induction, three  
136 main approaches have emerged: (1) prompting-  
137 based methods that assign personality traits through  
138 instructions (Serapio-García et al., 2025)(2) fine-  
139 tuning approaches that embed personas through  
140 training on character-specific data (Shao et al.,  
141 2023; Zhou et al., 2023), and (3) activation steering  
142 via persona vectors that manipulate neural repre-  
143 sentations corresponding to specific traits (Chen  
144 et al., 2025).

145 A parallel line of research evaluates LLMs  
146 through the lens of psychological constructs.  
147 Theory of Mind (ToM) benchmarks such as  
148 ToMBench (Chen et al., 2024b) assess social cog-  
149 nitive abilities, revealing that GPT-4 lags behind  
150 humans by over 10%, with trivial task modifica-  
151 tions causing significant performance degrada-  
152 tion (Ullman, 2023). Emotional intelligence  
153 benchmarks such as EQ-Bench (Paech, 2024) and  
154 EmoBench (Sabour et al., 2024) adopt psychology-  
155 grounded frameworks to evaluate emotional under-  
156 standing and application, finding substantial gaps  
157 between LLMs and humans. Moral reasoning has

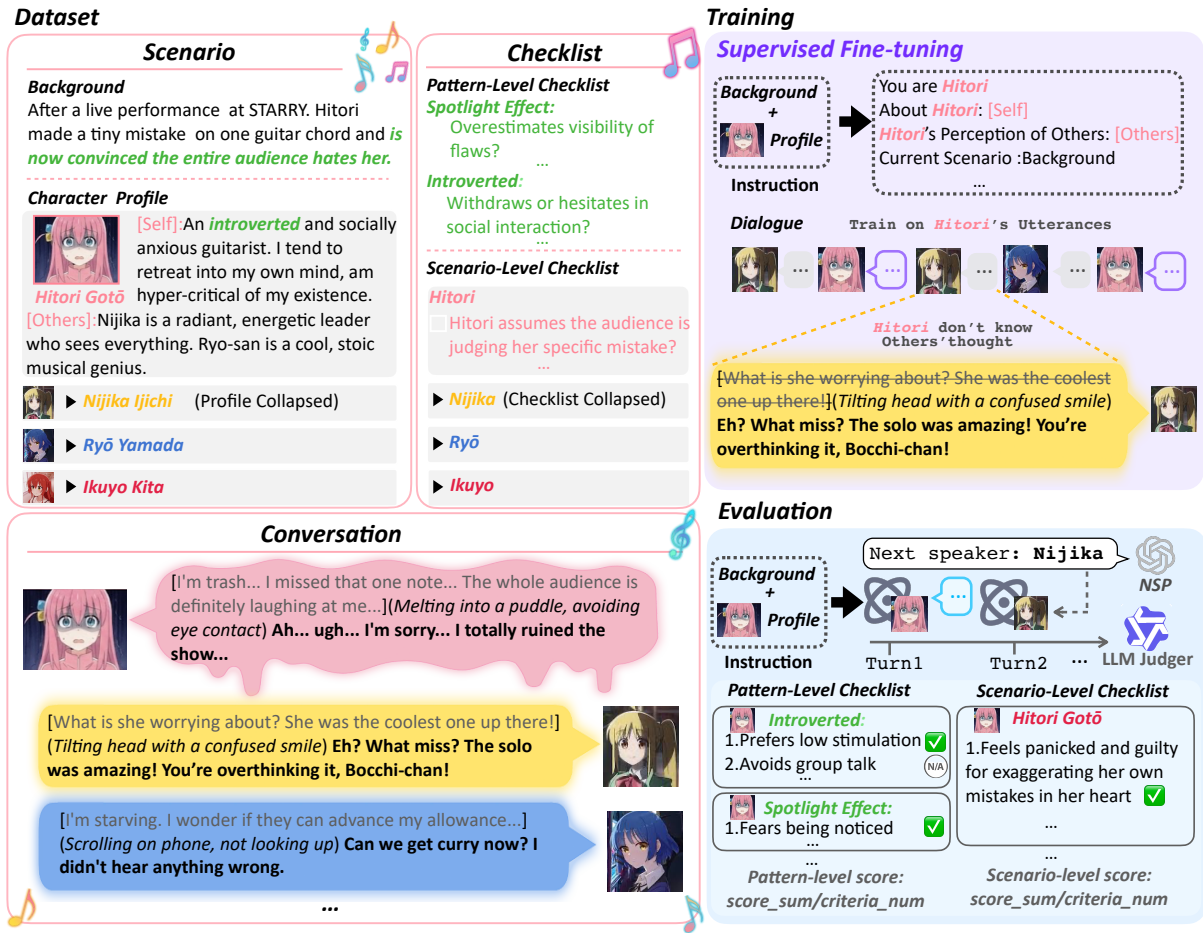


Figure 2: **HUMANLLM Framework**. **Left**: Dataset structure with scenarios, multi-turn conversations (inner thoughts in brackets, actions in parentheses), and dual-level checklists. **Top Right**: Supervised fine-tuning on target character utterances. **Bottom Right**: Evaluation via LLM judge scoring against pattern-level and scenario-level checklists.

158 been assessed through ETHICS (Hendrycks et al.,  
159 2023) and MoralBench (Ji et al., 2025), the later  
160 grounded in Moral Foundations Theory. Research  
161 on cognitive biases reveals that LLMs exhibit  
162 human-like irrationality but with divergent  
163 patterns (Macmillan-Scott and Musolesi, 2024).  
164 Personality assessment using validated instruments  
165 (BFI, MBTI) demonstrates that LLMs can manifest  
166 measurable traits (Pellert et al., 2024), though  
167 self-report validity remains questionable (Zou et al.,  
168 2025). Critically, recent work cautions that LLMs  
169 do not reliably simulate human psychology and fail  
170 to generalize across semantically equivalent scenarios  
171 (Schröder et al., 2025; Cheung et al., 2025).

### 172 3 HUMANLLM Dataset

173 This section introduces the HUMANLLM dataset, a  
174 psychologically grounded resource for training and  
175 evaluating anthropomorphic language models. We  
176 describe pattern collection (§3.1), pattern data con-

Statistic	Value
Total Patterns	244
Scenarios	11,359
Avg. Patterns per Scenario	3.5
Avg. Turns per Conversation	16.4
Pattern-Level Checklist Items	15 per pattern
Scenario-Level Checklist Items	2–6 per character

Table 1: HUMANLLM dataset statistics. Patterns include 100 personality traits and 144 social-cognitive patterns. Each scenario contains 2–6 characters with multi-turn conversations (12–20 turns).

177 instruction (§3.2), scenario and conversation generation  
178 (§3.3), and dual-level checklist design (§3.4).  
179 Table 1 summarizes the dataset statistics.

#### 180 3.1 Pattern Collection

181 Following the theoretical foundations established  
182 in §A.2, we compile patterns along two complementary  
183 dimensions corresponding to Lewin’s

184	Person-Environment framework.	
185	<b>Personality Traits (Person Dimension)</b> We	
186	adopt Goldberg’s 100 Unipolar Markers (Goldberg,	
187	1992), a psychometrically validated lexicon map-	
188	ping onto the Big Five dimensions with 20 trait de-	
189	scriptors each (Extraversion, Agreeableness, Con-	
190	scientiousness, Emotional Stability, Intellect).	
191	<b>Social-Cognitive Patterns (Environment Dimen-</b>	
192	<b>sion)</b> We curate situationally-activated psycholo-	
193	gical mechanisms through systematic review of	
194	established theoretical traditions, including cog-	
195	nitve biases (Tversky and Kahneman, 1974), so-	
196	cial influence (Cialdini et al., 2009), evolution-	
197	ary psychology (Buss, 2024), and motivation re-	
198	search (Deci and Ryan, 2000). From an initial pool	
199	of 232 documented patterns, we apply two filtering	
200	criteria: (1) sufficient empirical validation, and (2)	
201	non-redundancy with other patterns. This yields	
202	144 social-cognitive patterns (full taxonomy in Ap-	
203	pendix B.1).	
204	<b>3.2 Pattern Data Construction</b>	
205	Pattern data are structured representations of psy-	
206	chological patterns, as illustrated in Figure 1. We	
207	construct pattern data through a two-stage pipeline:	
208	literature retrieval followed by LLM-based synthe-	
209	sis.	
210	<b>Literature Retrieval</b> For each of the 244 pat-	
211	terns, we employ Gemini Deep Search to identify	
212	approximately 50 relevant academic papers. The	
213	search is guided by three retrieval dimensions: (1)	
214	foundational definitions from seminal works, (2)	
215	mechanistic explanations from theoretical and em-	
216	pirical studies, and (3) real-world applications from	
217	applied research. Retrieved references are filtered	
218	manually to remove irrelevant entries. Full-text	
219	documents are obtained through open-access APIs	
220	(Semantic Scholar, arXiv, OpenAlex, PubMed,	
221	Crossref); when full text is unavailable, abstracts	
222	are retained. This process yields a corpus of ap-	
223	proximately 12,000 papers across all patterns.	
224	<b>Pattern Synthesis</b> We employ Gemini 2.5 Pro to	
225	summarize each pattern’s literature corpus into a	
226	structured representation. Critically, the model is	
227	instructed to extract and summarize information <i>ex-</i>	
228	<i>clusively</i> from the provided 50 papers, rather than	
229	generating content from its parametric knowledge.	
230	Following the construct validity framework (§A.2),	
231	each pattern is organized into three components: (1)	
232	<b>Definition</b> —a precise characterization grounded	
	in authoritative sources; (2) <b>Core Mechanisms</b> —	233
	underlying cognitive, emotional, and behavioral	234
	processes that drive the pattern; (3) <b>Real-World</b>	235
	<b>Manifestations</b> —ecological expressions across di-	236
	verse contexts (e.g., response to stress, interper-	237
	sonal dynamics, professional settings). Details of	238
	the retrieval and synthesis prompts are provided in	239
	Appendix F.	240
	<b>3.3 Scenario and Conversation Generation</b>	241
	We synthesize scenarios and conversations through	242
	a two-stage pipeline using large language models.	243
	<b>Scenario Synthesis</b> Each scenario comprises a	244
	narrative background and 2–6 character profiles.	245
	Character profiles contain <i>self-perception</i> (identity,	246
	personality, background, motivations) and <i>other-</i>	247
	<i>perception</i> (knowledge and attitudes toward other	248
	characters), enabling realistic information asym-	249
	metry. Each scenario incorporates 2–5 patterns,	250
	with pattern combinations validated to filter seman-	251
	tically contradictory configurations. To ensure sit-	252
	uational diversity, we leverage the DIAMONDS	253
	model (§A.2) to generate scenario variants across	254
	different situational dimensions. Alongside each	255
	scenario, we synthesize <i>expected behavioral ten-</i>	256
	<i>dencies</i> —specifications of how characters should	257
	manifest the target patterns within the given con-	258
	text. Scenarios are generated using Gemini 2.5 Pro	259
	and Claude Sonnet 4.5, each contributing approxi-	260
	mately half of the total. This process yields 11,359	261
	scenarios.	262
	<b>Conversation Synthesis</b> Based on scenarios	263
	and expected behavioral tendencies, we synthe-	264
	size multi-turn conversations (12–20 turns) using	265
	Claude Sonnet 4.5. Each turn comprises three di-	266
	mensions: inner thoughts (enclosed in brackets),	267
	physical actions (enclosed in parentheses), and ver-	268
	bal expressions (Figure 2). Target patterns are nat-	269
	urally embedded across these dimensions, guided	270
	by the expected behavioral tendencies. Detailed	271
	generation prompts are provided in Appendix F.	272
	<b>3.4 Dual-Level Checklist Design</b>	273
	We design dual-level checklists for fine-grained as-	274
	essment of pattern expression, serving both evalu-	275
	ation (§4.2) and future reward modeling purposes.	276
	<b>Scenario-Level Checklist.</b> The expected behav-	277
	ioral tendencies synthesized alongside each sce-	278
	nario directly constitute the scenario-level check-	279
	list. These items specify context-specific behaviors	280

281 derived from particular pattern combinations and  
282 situational contexts (e.g., “Character defends con-  
283 ceptual integrity under deadline pressure, resisting  
284 shortcuts”). Each scenario contains 2–6 characters,  
285 but only characters assigned with target patterns  
286 receive checklist items; each such character has  
287 2–6 items reflecting their expected behavioral ten-  
288 dencies. During evaluation, we assess one target  
289 character per evaluation instance, treating other  
290 characters as contextual interlocutors.

291 **Pattern-Level Checklist** For each of the 244  
292 patterns, we construct 15 universal behavioral in-  
293 dicators through iterative refinement: initial gener-  
294 ation from pattern structure, validation against  
295 conversation samples, and generalization to ensure  
296 cross-context applicability (e.g., “Shows height-  
297 ened awareness of being observed” for spotlight  
298 effect). These pattern-level items are context-  
299 independent and apply to any character exhibit-  
300 ing the target pattern, complementing the scenario-  
301 specific items above. The complete construction  
302 procedure is detailed in Appendix B.4.

## 303 4 Training and Evaluation

304 This section presents the training procedure (§4.1)  
305 and evaluation protocol (§4.2) for HUMANLLM  
306 models.

### 307 4.1 Training

308 We train HUMANLLM-8B and HUMANLLM-  
309 32B on Qwen3-8B and Qwen3-32B respectively  
310 through supervised fine-tuning. Training data com-  
311 position and hyperparameters are detailed in Ap-  
312 pendix C.

313 **Supervised Fine-Tuning** We convert each char-  
314 acter’s dialogue turns within a scenario into a sepa-  
315 rate training sample in ShareGPT format, yield-  
316 ing 30,543 HUMANLLM samples from 10,265  
317 training scenarios. The model is trained to gener-  
318 ate responses that naturally express the target  
319 patterns through the trinity of expression: in-  
320 ner thoughts (square brackets), physical actions  
321 (parentheses), and verbal expressions. To main-  
322 tain general capabilities and role-playing abili-  
323 ties, we augment the HUMANLLM data with two  
324 complementary sources: OpenThoughts-114k for  
325 instruction-following capabilities (30,543 samples)  
326 and CoSER (Wang et al., 2025) for role-playing  
327 dialogue (15,272 samples). The final training mix-  
328 ture comprises 76,358 samples with a ratio of 4:4:2

(HUMANLLM: OpenThoughts : CoSER) by sam- 329  
ple count. 330

### 4.2 Evaluation 331

332 The dual-level checklists enable evaluation at two  
333 complementary granularities, corresponding to two  
334 proposed metrics. Both metrics employ GPT-5-  
335 mini as judge with ternary scoring: +1 (satisfied),  
336 0 (not exhibited), −1 (violated).

337 **Metrics** We propose two metrics that capture  
338 complementary aspects of pattern expression:

339 (1) *Individual Pattern Expression (IPE)* mea-  
340 sures whether each pattern is expressed according  
341 to its psychological definition. IPE uses pattern-  
342 level checklists, which contain 15 universal be-  
343 havioral indicators derived directly from each pat-  
344 tern’s definition and mechanisms (e.g., for *spot-*  
345 *light effect*: “Overestimates others’ attention to  
346 own appearance”). These indicators are context-  
347 independent—they assess whether the pattern’s  
348 characteristic behaviors appear, regardless of situ-  
349 ational details. The sample-level IPE is computed  
350 as the mean score across all pattern-level check-  
351 list items; the overall IPE is the mean across all  
352 evaluation samples.

353 (2) *Multi-Pattern Dynamics (MPD)* measures  
354 whether multiple patterns interact appropriately  
355 within specific situational configurations. Prior  
356 role-playing evaluations often rely on single-label  
357 assessments (e.g., “Is this character assertive?”),  
358 which reduce complex personalities to stereotyp-  
359 ical behaviors and fail to capture how traits mod-  
360 ulate each other in context. MPD addresses this  
361 limitation by using scenario-level checklists that  
362 specify expected behavioral tendencies emerging  
363 from particular pattern combinations. For instance,  
364 a character assigned both *assertive* and *spotlight ef-*  
365 *fect* in a public speaking scenario should not simply  
366 exhibit “confident speech”—the checklist captures  
367 the nuanced tension: “Projects outward confidence  
368 in expressing opinions, yet harbors internal anxiety  
369 regarding audience scrutiny.” Such items re-  
370 flect how one pattern (spotlight effect) constrains  
371 or modulates another (assertive), moving evalua-  
372 tion beyond stereotypical label-to-behavior map-  
373 pings toward authentic multi-dimensional charac-  
374 terization. The sample-level MPD is computed as  
375 the mean score across all scenario-level checklist  
376 items; the overall MPD is the mean across all eval-  
377 uation samples.

378 In summary, IPE evaluates fidelity to individual

pattern definitions, while MPD evaluates the emergent dynamics when multiple patterns interact—capturing the psychological realism that single-label evaluations miss. Data splits for generalization assessment are detailed in Appendix D.1.

## 5 Experiments

### 5.1 Experimental Setup

**Training Configuration** We train HUMANLLM-8B and HUMANLLM-32B through supervised fine-tuning as described in §4.1. Detailed hyperparameters are provided in Appendix C.

**Baselines** We compare against a comprehensive set of baselines spanning proprietary and open-source models, all evaluated in zero-shot settings: (1) *Closed-source models*, including GPT-5 (OpenAI, 2025), Claude Sonnet 4.5 (Anthropic, 2025), and Gemini 3 Pro (DeepMind, 2025). (2) *Open-source models*, including Qwen3-8B/32B/235B (Yang et al., 2025), DeepSeek-V3.2 (DeepSeek-AI et al., 2025b), and DeepSeek-R1 (DeepSeek-AI et al., 2025a).

**External Benchmarks** Beyond our proposed IPE and MPD metrics, we evaluate on established role-playing benchmarks: (1) *LifeChoice* (Xu et al., 2024b): evaluates persona-driven decision-making through life choice scenarios; (2) *CroSS-MR* (Yuan et al., 2024): assesses character motivation recognition; Additionally, we include CoSER’s evaluation metrics—Anthropomorphism and Character Fidelity (Wang et al., 2025)—for comparative analysis of evaluation paradigms (§5.5).

### 5.2 Main Results

Table 2 presents the main experimental results on our proposed IPE and MPD metrics.

**Overall Performance** Among all evaluated models, Gemini 3 Pro achieves the highest scores on both metrics (IPE: 41.1%, MPD: 85.3%), followed by Claude Sonnet 4.5 (IPE: 34.6%, MPD: 79.7%). Our HUMANLLM-32B achieves competitive performance (IPE: 32.6%, MPD: 73.8%), outperforming all open-source baselines of comparable or larger scale. Notably, HUMANLLM-8B (IPE: 25.5%, MPD: 70.1%) surpasses Qwen3-32B (IPE: 26.2%, MPD: 66.0%) on MPD despite having 4× fewer parameters, demonstrating the effectiveness of our psychologically grounded training data.

Model	IPE	MPD
<i>Close-Source</i>		
GPT-5	15.7	43.2
Claude Sonnet 4.5	34.6	79.7
Gemini 3 Pro	41.1	85.3
<i>Open-Source</i>		
Qwen3-8B	18.8	54.2
Qwen3-32B	26.2	66.0
Qwen3-235B	34.1	72.7
DeepSeek-V3.2	22.0	65.3
DeepSeek-R1	23.5	68.8
<i>Ours</i>		
HUMANLLM-8B	25.5	70.1
HUMANLLM-32B	32.6	73.8

Table 2: Main results on IPE (Individual Pattern Expression) and MPD (Multi-Pattern Dynamics). Higher values indicate better performance. All values are reported in percentage (%). Results are averaged across ID\_eval, OOD\_eval, and Mixed\_eval splits.

**Close-Source vs. Open-Source Gap** A substantial performance gap exists between close-source and open-source models. The best close-source model (Gemini 3 Pro) outperforms the best open-source baseline (Qwen3-235B) by +7.0 on IPE and +12.6 on MPD. Interestingly, GPT-5 exhibits unexpectedly low performance (IPE: 15.7%, MPD: 43.2%), ranking below most open-source alternatives. Qualitative analysis suggests that GPT-5’s strong instruction-following tendency leads to overly literal interpretations of role-playing prompts, resulting in shallow pattern expression. This finding aligns with recent observations that general-purpose capabilities do not automatically transfer to nuanced psychological simulation (Wang et al., 2024).

**Scaling Effects** Within model families, we observe consistent scaling improvements. For Qwen3, IPE increases from 18.8 (8B) to 26.2 (32B) to 34.1 (235B), representing a +81% relative improvement from 8B to 235B. Similarly, our HUMANLLM models show scaling gains: HUMANLLM-32B outperforms HUMANLLM-8B by +7.1 on IPE (+28%) and +3.7 on MPD (+5%). The relatively smaller MPD scaling gap suggests that multi-pattern dynamics may be more sample-efficient to learn, while individual pattern expression benefits more from increased model capacity.

**IPE vs. MPD Dynamics** Across all models, MPD scores consistently exceed IPE scores. This asymmetry likely reflects the evaluation granular-

Model Variant	IPE	MPD
Qwen3-8B (base)	18.8	54.2
Qwen3-8B (OT+CoSER)	8.9	31.1
HUMANLLM-8B	25.5	70.1

Table 3: Ablation study results (%). “OT+CoSER” denotes training on OpenThoughts and CoSER without HUMANLLM data.

ity: MPD assesses scenario-level behavioral tendencies (2–6 items per character), while IPE evaluates against 15 fine-grained pattern indicators. The gap is most pronounced for weaker models (e.g., GPT-5: 15.7% vs. 43.2%), suggesting that models can produce superficially coherent multi-pattern behavior without deeply understanding individual pattern mechanisms.

### 5.3 Ablation Study

To isolate the contribution of our psychologically grounded dataset, we conduct an ablation study comparing three model variants (Table 3).

**Effect of HUMANLLM Data** The comparison between HUMANLLM-8B and Qwen3-8B (base) reveals substantial improvements: +6.7 on IPE (+36%) and +15.9 on MPD (+29%). These gains demonstrate that our synthesized conversations, grounded in psychological pattern definitions, effectively teach models to express cognitive and behavioral patterns with higher fidelity.

**Negative Transfer from Generic Data** Surprisingly, the OT+CoSER variant (trained on OpenThoughts and CoSER without HUMANLLM data) performs *worse* than the base model on both metrics: IPE drops from 18.8 to 8.9 (−53%), and MPD drops from 54.2 to 31.1 (−43%). This counterintuitive result suggests negative transfer: generic instruction-following data (OpenThoughts) and conventional role-playing data (CoSER) may inadvertently suppress the base model’s latent ability to simulate psychological patterns. We hypothesize that these datasets reinforce “helpful assistant” behaviors that conflict with authentic expression of cognitively biased or emotionally complex characters.

**Synergistic Effect** The full HUMANLLM-8B model, trained on the combined mixture (HumanLLM + OT + CoSER with 4:4:2 ratio), achieves the best performance. This indicates that HU-

Model	LifeChoice	CroSS-MR
GPT-5	85.53	62.25
Claude Sonnet 4.5	85.49	68.24
Qwen3-8B	44.51	53.26
Qwen3-32B	47.71	63.37
Qwen3-8B (OT+CoSER)	46.99	54.98
HUMANLLM-8B	47.19	54.23
HUMANLLM-32B	50.64	64.27

Table 4: External benchmark results (%). LifeChoice evaluates persona-driven decision-making; CroSS-MR assesses character motivation recognition.

MANLLM data not only compensates for the negative transfer but creates a synergistic effect where psychological grounding enhances the utility of general-purpose and role-playing data. The pattern-rich conversations appear to serve as “anchors” that prevent the model from collapsing toward generic prosocial behaviors.

### 5.4 External Benchmark Evaluation

To assess generalization beyond our proposed metrics, we evaluate on established role-playing benchmarks (Table 4). On LifeChoice, while closed-source models (e.g., GPT-5: 85.53%) dominate, HUMANLLM-32B (50.64%) still outperforms Qwen3-32B (47.71%) by +2.93%; this modest improvement suggests that LifeChoice’s binary decision scenarios may rely more on surface-level reasoning than the deep psychological patterns targeted by our approach. Conversely, on CroSS-MR, HUMANLLM-32B (64.27%) surpasses both Qwen3-32B (63.37%) and GPT-5 (62.25%), with the 8B variant also showing competitive performance (54.23% vs. 53.26%), indicating that motivation recognition benefits from our training but is less dependent on complex pattern simulation than decision-making tasks.

**Benchmark Limitations** The moderate improvements on external benchmarks, contrasted with substantial gains on our IPE/MPD metrics, reveal a fundamental mismatch between existing evaluation paradigms and psychologically grounded simulation. LifeChoice and CroSS-MR evaluate behavioral outcomes (decisions, motivations) rather than the cognitive processes underlying those outcomes. These benchmarks, while valuable for assessing surface-level role-playing, do not capture the pattern-specific behavioral indicators that define authentic psychological simulation. This observation motivates our development of the dual-level

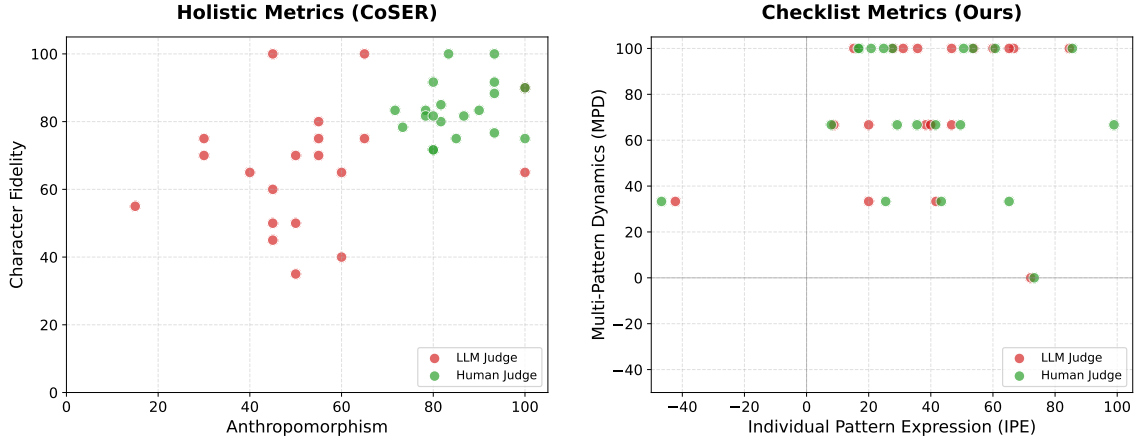


Figure 3: Human-LLM evaluation alignment comparison. **Left:** Holistic metrics show clear separation between LLM and human judgments, with systematic LLM underestimation of psychologically complex behaviors. **Right:** Our checklist metrics demonstrate strong overlap between LLM and human distributions, indicating robust alignment with expert judgment

Metric	Human	LLM	$\Delta$	$r$
<i>Holistic Metrics (CoSER)</i>				
Anthropomorphism	85.2	53.0	-32.2	0.41
Character Fidelity	83.5	66.8	-16.7	0.62
<i>Checklist Metrics (Ours)</i>				
IPE	39.0	38.6	-0.4	<b>0.91</b>
MPD	71.7	75.0	+3.3	0.88

Table 5: Human-LLM agreement analysis across 20 sampled scenarios. Holistic metrics (0–100 scale) show weak correlation and large systematic bias; our checklist metrics ( $[-1, +1]$  normalized to  $[0, 1]$ ) achieve significantly higher alignment with human judgment

checklist framework, which we validate in §5.5.

## 5.5 Evaluation Framework Analysis

To validate our dual-level checklist, we compared automated GPT-5-mini evaluations against three human experts across 20 diverse scenarios using both holistic metrics (CoSER) and our checklist criteria. Evaluations were voluntary and anonymous.

**Human-LLM Alignment and Normative Confounding** As shown in Table 5, holistic metrics exhibit weak-to-moderate alignment (Anthropomorphism:  $r = 0.41, \Delta = -32.2$ ; Fidelity:  $r = 0.62, \Delta = -16.7$ ). This divergence stems from what we term **normative confounding**: LLM judges implicitly conflate “good anthropomorphism” with “prosocial behavior” (e.g., empathy, rationality), penalizing realistic but negative human traits like defensiveness. In contrast, our

checklist metrics achieve robust alignment (IPE:  $r = 0.91, \Delta = -0.4$ ; MPD:  $r = 0.88, \Delta = +3.3$ ) by decomposing behaviors into value-neutral indicators. This effectively decouples *simulation accuracy* from *social desirability*.

**Case Study** Ideally illustrating this, a scenario featuring *ultimate attribution error* (defensively blaming out-groups) received a low holistic score (5/100) due to cited “lack of empathy.” However, our checklist correctly validated the pattern through definition-grounded questions (e.g., “Does the character attribute failure to external factors?”), confirming that the behavior, while socially undesirable, was psychologically accurate (Appendix E).

## 6 Conclusion

We present HUMANLLM, a framework leveraging 244 patterns from  $\sim 12,000$  papers and 11,359 scenarios to enable dynamic cognitive interactions rather than isolated trait mappings. Our evaluation achieves strong human alignment ( $r = 0.91$ ) and identifies **normative confounding**, where conventional metrics mistake social desirability for simulation accuracy. Experiments confirm the necessity of psychological grounding for authentic expression, advancing the field from surface-level behavioral mimicry to cognitive modeling—simulating the underlying psychological processes that generate human behavior.

578  
579  
580  
581  
582  
583  
584  
  
585  
586  
587  
588  
589  
590  
  
591  
592  
593  
594  
  
595  
596  
597  
598  
  
599  
  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621

## Limitations

**LLM-as-Judge Limitations** Our evaluation relies on GPT-5-mini as the judge. While we mitigate variability through repeated evaluation, LLM judges may exhibit systematic biases or inconsistent reasoning. Future work should incorporate broader human evaluation.

**Cultural Adaptability** The psychological theories underlying our pattern taxonomy originate predominantly from WEIRD (Western, Educated, Industrialized, Rich, Democratic) populations. Pattern expressions may manifest differently across cultural contexts.

**Synthetic Data Limitations** All training conversations are synthetically generated by LLMs, which may introduce systematic biases or fail to capture the full complexity of authentic human interactions.

**Evaluation Scope** Our evaluation focuses on text-based dialogue and may not fully capture temporal consistency across extended interactions or embodied behavior in multimodal settings.

## Ethical Statement

**Enhanced Anthropomorphism and Safety Tension** By design, HUMANLLM improves the simulation of human cognitive patterns, including irrational biases and negative personality traits (e.g., antagonism, defensiveness). While essential for realistic role-playing, this creates a fundamental tension with standard safety alignment goals. There is a risk that models trained to faithfully execute maladaptive or antisocial behaviors could be exploited to generate toxic content or reinforce harmful stereotypes if the role-playing context is not strictly sandboxed.

**Manipulation and Social Influence** The framework explicitly models mechanisms of social influence and persuasion (e.g., authority bias, reciprocity). A model that deeply "understands" and simulates these psychological triggers possesses a heightened capacity for social engineering. In malicious hands, such agents could be deployed to manipulate user sentiment or exploit cognitive vulnerabilities more effectively than standard "helpful" assistants.

## References

Anthropic. 2025. Anthropic’s transparency hub: Model report. <https://www.anthropic.com/transparency/model-report>. Accessed: 2026-01-03. 623-626

David M Buss. 2024. *Evolutionary psychology: The new science of the mind*. Routledge. 627-628

Jiangjie Chen, Xintao Wang, Rui Xu, Siyu Yuan, Yikai Zhang, Wei Shi, Jian Xie, Shuang Li, Ruihan Yang, Tinghui Zhu, Aili Chen, Nianqi Li, Lida Chen, Caiyu Hu, Siye Wu, Scott Ren, Ziquan Fu, and Yanghua Xiao. 2024a. [From Persona to Personalization: A Survey on Role-Playing Language Agents](#). *Preprint*, arXiv:2404.18231. 629-635

Runjin Chen, Andy Arditi, Henry Sleight, Owain Evans, and Jack Lindsey. 2025. [Persona vectors: Monitoring and controlling character traits in language models](#). *Preprint*, arXiv:2507.21509. 636-639

Zhuang Chen, Jincenzi Wu, Jinfeng Zhou, Bosi Wen, Guanqun Bi, Gongyao Jiang, Yaru Cao, Mengting Hu, Yunghwei Lai, Zexuan Xiong, and Minlie Huang. 2024b. [ToMBench: Benchmarking Theory of Mind in Large Language Models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15959–15983, Bangkok, Thailand. Association for Computational Linguistics. 640-648

Vanessa Cheung, Maximilian Maier, and Falk Lieder. 2025. [Large language models show amplified cognitive biases in moral decision-making](#). *Proceedings of the National Academy of Sciences*, 122(25):e2412015122. 649-653

Robert B Cialdini and 1 others. 2009. *Influence: Science and practice*, volume 4. Pearson education Boston. 654-656

Lee J Cronbach and Paul E Meehl. 1955. Construct validity in psychological tests. *Psychological bulletin*, 52(4):281. 657-659

Edward L Deci and Richard M Ryan. 2000. The "what" and "why" of goal pursuits: Human needs and the self-determination of behavior. *Psychological inquiry*, 11(4):227–268. 660-663

Google DeepMind. 2025. [Gemini 3 pro model card](#). Model card, Google DeepMind. 664-665

DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, and 181 others. 2025a. [Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning](#). *Preprint*, arXiv:2501.12948. 666-673

DeepSeek-AI, Aixin Liu, Aoxue Mei, Bangcai Lin, Bing Xue, Bingxuan Wang, Bingzheng Xu, Bochao 674-675

676	Wu, Bowei Zhang, Chaofan Lin, Chen Dong, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenhao Xu, Chong Ruan, Damai Dai, Daya Guo, Dejian Yang, and 245 others. 2025b. <a href="#">Deepseek-v3.2: Pushing the frontier of open large language models</a> . <i>Preprint</i> , arXiv:2512.02556.	Joon Sung Park, Joseph C. O'Brien, Carrie J. Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. <a href="#">Generative Agents: Interactive Simulacra of Human Behavior</a> . <i>Preprint</i> , arXiv:2304.03442.	728
677			729
678			730
679			731
680			732
681			
682	John M. Digman. 1990. <a href="#">Personality structure: Emergence of the five-factor model</a> . <i>Annual Review of Psychology</i> , 41:417–440.	Max Pellert, Clemens M Lechner, Claudia Wagner, Beatrice Rammstedt, and Markus Strohmaier. 2024. Ai psychometrics: Assessing the psychological profiles of large language models through psychometric inventories. <i>Perspectives on Psychological Science</i> , 19(5):808–826.	733
683			734
684			735
685	Lewis R Goldberg. 1992. The development of markers for the big-five factor structure. <i>Psychological assessment</i> , 4(1):26.	John F Rauthmann, David Gallardo-Pujol, Esther M Guillaume, Elysia Todd, Christopher S Nave, Ryne A Sherman, Matthias Ziegler, Ashley Bell Jones, and David C Funder. 2014. The situational eight diamonds: a taxonomy of major dimensions of situation characteristics. <i>Journal of personality and social psychology</i> , 107(4):677.	736
686			737
687			738
688	Pengrui Han, Rafal Kocielnik, Peiyang Song, Ramit Debnath, Dean Mobbs, Anima Anandkumar, and R. Michael Alvarez. 2025. <a href="#">The personality illusion: Revealing dissociation between self-reports behavior in llms</a> . <i>Preprint</i> , arXiv:2509.03730.	Sahand Sabour, Siyang Liu, Zheyuan Zhang, June Liu, Jinfeng Zhou, Alvionna Sunaryo, Tatia Lee, Rada Mihalcea, and Minlie Huang. 2024. <a href="#">EmoBench: Evaluating the Emotional Intelligence of Large Language Models</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5986–6004, Bangkok, Thailand. Association for Computational Linguistics.	739
689			740
690			741
691			742
692			743
693	Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. 2023. <a href="#">Aligning ai with shared human values</a> . <i>Preprint</i> , arXiv:2008.02275.	Sarah Schröder, Thekla Morgenroth, Ulrike Kuhl, Valerie Vaquet, and Benjamin Paaßen. 2025. <a href="#">Large language models do not simulate human psychology</a> . <i>Preprint</i> , arXiv:2508.06950.	744
694			745
695			746
696			747
697	Jen-tse Huang, Man Ho Lam, Eric John Li, Shujie Ren, Wenxuan Wang, Wenxiang Jiao, Zhaopeng Tu, and Michael R Lyu. 2024a. Apathetic or empathetic? evaluating llms' emotional alignments with humans. <i>Advances in Neural Information Processing Systems</i> , 37:97053–97087.	Greg Serapio-García, Mustafa Safdari, Clément Crepy, Luning Sun, Stephen Fitz, Peter Romero, Marwa Abdulhai, Aleksandra Faust, and Maja Matarić. 2025. <a href="#">Personality traits in large language models</a> . <i>Preprint</i> , arXiv:2307.00184.	748
698			749
699			750
700			751
701			752
702			753
703	Jen-tse Huang, Wenxuan Wang, Eric John Li, Man Ho Lam, Shujie Ren, Youliang Yuan, Wenxiang Jiao, Zhaopeng Tu, and Michael Lyu. 2024b. On the humanity of conversational ai: Evaluating the psychological portrayal of llms. In <i>The Twelfth International Conference on Learning Representations</i> .	Murray Shanahan, Kyle McDonell, and Laria Reynolds. 2023. <a href="#">Role play with large language models</a> . <i>Nature</i> , 623(7987):493–498.	754
704			755
705			756
706			757
707			758
708			759
709	Jianchao Ji, Yutong Chen, Mingyu Jin, Wujiang Xu, Wenyue Hua, and Yongfeng Zhang. 2025. <a href="#">Moralbench: Moral evaluation of llms</a> . <i>Preprint</i> , arXiv:2406.04428.	Yunfan Shao, Linyang Li, Junqi Dai, and Xipeng Qiu. 2023. <a href="#">Character-LLM: A Trainable Agent for Role-Playing</a> . <i>Preprint</i> , arXiv:2310.10158.	760
710			761
711			762
712			763
713	K. Lewin. 1936. <a href="#">Principles of topological psychology</a> .	Amos Tversky and Daniel Kahneman. 1974. Judgment under uncertainty: Heuristics and biases: Biases in judgments reveal some heuristics of thinking under uncertainty. <i>science</i> , 185(4157):1124–1131.	764
714	Cheng Li, Ziang Leng, Chenxi Yan, Junyi Shen, Hao Wang, Weishi MI, Yaying Fei, Xiaoyang Feng, Song Yan, HaoSheng Wang, Linkang Zhan, Yaokai Jia, Pingyu Wu, and Haozhen Sun. 2023. <a href="#">Chatharuhi: Reviving anime character in reality via large language model</a> . <i>Preprint</i> , arXiv:2308.09597.	Tomer Ullman. 2023. <a href="#">Large language models fail on trivial alterations to theory-of-mind tasks</a> . <i>Preprint</i> , arXiv:2302.08399.	765
715			766
716			767
717			768
718			769
719			770
720	Olivia Macmillan-Scott and Mirco Musolesi. 2024. <a href="#">(ir)rationality and cognitive biases in large language models</a> . <i>Preprint</i> , arXiv:2402.09193.	Xintao Wang, Heng Wang, Yifei Zhang, Xinfeng Yuan, Rui Xu, Jen-tse Huang, Siyu Yuan, Haoran Guo, Jiangjie Chen, Shuchang Zhou, Wei Wang, and Yanghua Xiao. 2025. <a href="#">CoSER: Coordinating LLM-Based Persona Simulation of Established Roles</a> . <i>Preprint</i> , arXiv:2502.09082.	771
721			772
722			773
723	OpenAI. 2025. <a href="#">Gpt-5 system card</a> . System card, OpenAI.		774
724			775
725	Samuel J. Paech. 2024. <a href="#">Eq-bench: An emotional intelligence benchmark for large language models</a> . <i>Preprint</i> , arXiv:2312.06281.		776
726			777
727			778

782 Xintao Wang, Yunze Xiao, Jen-tse Huang, Siyu Yuan,  
783 Rui Xu, Haoran Guo, Quan Tu, Yaying Fei, Ziang  
784 Leng, Wei Wang, Jiangjie Chen, Cheng Li, and  
785 Yanghua Xiao. 2024. [InCharacter: Evaluating Per-  
786 sonality Fidelity in Role-Playing Agents through Psy-  
787 chological Interviews](#). In *Proceedings of the 62nd  
788 Annual Meeting of the Association for Computational  
789 Linguistics (Volume 1: Long Papers)*, pages 1840–  
790 1873, Bangkok, Thailand. Association for Computa-  
791 tional Linguistics.

792 Rui Xu, Dakuan Lu, Xiaoyu Tan, Xintao Wang,  
793 Siyu Yuan, Jiangjie Chen, Wei Chu, and Yinghui  
794 Xu. 2024a. [MINDECHO: Role-Playing Language  
795 Agents for Key Opinion Leaders](#). *arXiv preprint*.

796 Rui Xu, Xintao Wang, Jiangjie Chen, Siyu Yuan, Xin-  
797 feng Yuan, Jiaqing Liang, Zulong Chen, Xiaoqing  
798 Dong, and Yanghua Xiao. 2024b. [Character is des-  
799 tinity: Can role-playing language agents make persona-  
800 driven decisions?](#) *Preprint*, arXiv:2404.12138.

801 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,  
802 Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao,  
803 Chengen Huang, Chenxu Lv, Chujie Zheng, Day-  
804 iheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao  
805 Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41  
806 others. 2025. [Qwen3 technical report](#). *Preprint*,  
807 arXiv:2505.09388.

808 Xinfeng Yuan, Siyu Yuan, Yuhan Cui, Tianhe Lin, Xin-  
809 tao Wang, Rui Xu, Jiangjie Chen, and Deqing Yang.  
810 2024. [Evaluating character understanding of large  
811 language models via character profiling from fictional  
812 works](#). In *Proceedings of the 2024 Conference on  
813 Empirical Methods in Natural Language Processing*,  
814 pages 8015–8036, Miami, Florida, USA. Association  
815 for Computational Linguistics.

816 Yutong Zhang, Dora Zhao, Jeffrey T. Hancock, Robert  
817 Kraut, and Diyi Yang. 2025. [The rise of ai com-  
818 panions: How human-chatbot relationships influence  
819 well-being](#). *Preprint*, arXiv:2506.12605.

820 Jiayu Zhou, Jen-tse Huang, Xuhui Zhou, Man Ho  
821 Lam, Xintao Wang, Hao Zhu, Wenxuan Wang, and  
822 Maarten Sap. 2025. The pimmur principles: Ensur-  
823 ing validity in collective behavior of llm societies.  
824 *arXiv preprint arXiv:2509.18052*.

825 Jinfeng Zhou, Zhuang Chen, Dazhen Wan, Bosi Wen,  
826 Yi Song, Jifan Yu, Yongkang Huang, Libiao Peng,  
827 Jiaming Yang, Xiyao Xiao, Sahand Sabour, Xiaohan  
828 Zhang, Wenjing Hou, Yijia Zhang, Yuxiao Dong, Jie  
829 Tang, and Minlie Huang. 2023. [Characterglm: Cust-  
830 omizing chinese conversational ai characters with  
831 large language models](#). *Preprint*, arXiv:2311.16832.

832 Huiqi Zou, Pengda Wang, Zihan Yan, Tianjun Sun, and  
833 Ziang Xiao. 2025. [Can llm "self-report"?: Evaluat-  
834 ing the validity of self-report scales in measuring  
835 personality design in llm-based chatbots](#). *Preprint*,  
836 arXiv:2412.00207.

837	<b>A Preliminaries</b>	
838	<b>A.1 Role-Playing Language Agents</b>	
839	Role-Playing Language Agents (RPLAs) are AI	
840	systems designed to simulate assigned personas	
841	in conversational interactions (Chen et al., 2024a).	
842	Recent advances have enabled RPLAs to embody	
843	diverse characters, from historical figures (Shao	
844	et al., 2023) to fictional personas (Li et al., 2023),	
845	with applications spanning digital clones, AI com-	
846	panions, and social simulation (Park et al., 2023).	
847	Three main approaches have emerged for per-	
848	sona induction: (1) <i>prompting-based methods</i> that	
849	assign traits through instructions (Serapio-García	
850	et al., 2025), (2) <i>fine-tuning approaches</i> that em-	
851	bed personas through training on character-specific	
852	data (Shao et al., 2023; Zhou et al., 2023), and	
853	(3) <i>activation steering</i> via persona vectors that	
854	manipulate neural representations (Chen et al.,	
855	2025). Evaluation efforts have assessed LLMs	
856	through psychological benchmarks targeting The-	
857	ory of Mind (Chen et al., 2024b), emotional intelli-	
858	gence (Paech, 2024; Huang et al., 2024a,b; Sabour	
859	et al., 2024), and moral reasoning (Ji et al., 2025).	
860	However, these approaches typically focus on iso-	
861	lated traits or narrow psychological dimensions	
862	without treating anthropomorphism as a holistic	
863	objective grounded in multi-pattern dynamics.	
864	<b>A.2 Psychological Foundations</b>	
865	Our framework draws upon established psychologi-	
866	cal theories that inform both pattern taxonomy and	
867	data construction.	
868	<b>Lewin’s Field Theory</b> Kurt Lewin’s formulation	
869	$B = f(P, E)$ posits that behavior emerges from	
870	the dynamic interaction between Person and En-	
871	vironment (Lewin, 1936). This principle guides	
872	our two-dimensional pattern taxonomy: stable per-	
873	sonality traits (Person) and situationally-activated	
874	cognitive processes (Environment).	
875	<b>Big Five and Personality Traits</b> The Big Five	
876	model (Digman, 1990) represents the dominant	
877	paradigm in personality psychology, characteriz-	
878	ing individual differences along five dimensions.	
879	We adopt Goldberg’s 100 Unipolar Markers (Gold-	
880	berg, 1992), which operationalize this framework	
881	with 20 validated trait descriptors per dimension:	
882	Extraversion, Agreeableness, Conscientiousness,	
883	Emotional Stability, and Intellect.	
	<b>Social-Cognitive Patterns</b> Complementing sta-	884
	ble personality traits, social-cognitive patterns rep-	885
	resent context-triggered psychological mechanisms	886
	documented in behavioral research. These include	887
	cognitive biases and heuristics from judgment and	888
	decision-making research (Tversky and Kahneman,	889
	1974), social influence mechanisms (Cialdini et al.,	890
	2009), evolutionary adaptations (Buss, 2024), and	891
	motivational processes (Deci and Ryan, 2000).	892
	<b>Construct Validity</b> The construct validity frame-	893
	work (Cronbach and Meehl, 1955) establishes that	894
	psychological constructs should be precisely de-	895
	fined, grounded in theoretical mechanisms, and val-	896
	idated through observable manifestations. This in-	897
	forms our three-component pattern structure: Def-	898
	inition, Core Mechanisms, and Real-World Mani-	899
	festations.	900
	<b>DIAMONDS Model</b> The DIAMONDS frame-	901
	work (Rauthmann et al., 2014) characterizes situa-	902
	tions along eight psychological dimensions: Duty,	903
	Intellect, Adversity, Mating, Positivity, Negativity,	904
	Deception, and Sociality. We leverage this tax-	905
	onomy to ensure ecological diversity in scenario	906
	generation.	907
	<b>B Dataset Details</b>	908
	<b>B.1 Pattern Taxonomy</b>	909
	Our pattern taxonomy comprises 244 patterns	910
	along two dimensions: 100 personality traits (Ta-	911
	ble 6) and 144 social-cognitive patterns (Table 7).	912
	We adopt Goldberg’s 100 Unipolar Markers (Gold-	913
	berg, 1992), with 20 trait descriptors per Big Five	914
	dimension, each containing 10 positive-pole and 10	915
	negative-pole descriptors. The 144 social-cognitive	916
	patterns are curated from four theoretical traditions	917
	through systematic literature review; from an initial	918
	pool of 232 documented patterns, we filter based on	919
	empirical validation and non-redundancy criteria.	920
	<b>B.2 Pattern Data Construction</b>	921
	Each pattern is developed into a structured rep-	922
	resentation through systematic literature review.	923
	For each pattern, we employ Gemini Deep Search	924
	to identify approximately 50 relevant academic	925
	papers spanning foundational definitions, mech-	926
	anistic explanations, and real-world applications.	927
	Retrieved references are manually filtered to re-	928
	move irrelevant entries. Full-text documents are	929
	obtained through open-access APIs (Semantic	930

Dimension	Positive Pole	Negative Pole
Extraversion	talkative, assertive, active, energetic, outgoing, enthusiastic, daring, gregarious, bold, spontaneous	quiet, reserved, shy, inhibited, timid, withdrawn, unassertive, introverted, silent, unenergetic
Agreeableness	sympathetic, kind, appreciative, affectionate, soft-hearted, warm, generous, trusting, helpful, cooperative	cold, unsympathetic, harsh, rude, unkind, cruel, quarrelsome, critical, antagonistic, callous
Conscientiousness	organized, responsible, dependable, thorough, efficient, practical, deliberate, conscientious, neat, careful	disorganized, careless, irresponsible, undependable, sloppy, impractical, haphazard, negligent, untidy, rash
Emotional Stability	relaxed, calm, at ease, unemotional, poised, composed, secure, stable, content, placid	anxious, moody, envious, touchy, fretful, temperamental, insecure, nervous, jealous, high-strung
Intellect	creative, imaginative, intellectual, philosophical, complex, deep, artistic, bright, perceptive, introspective	uncreative, unimaginative, unintellectual, unphilosophical, simple, shallow, unartistic, dull, imperceptive, uninquisitive

Table 6: 100 Personality Traits organized by Big Five dimensions, adapted from Goldberg’s Unipolar Markers (Goldberg, 1992).

Scholar, arXiv, OpenAlex, PubMed, Crossref); when full text is unavailable, abstracts are retained.

We employ Gemini 2.5 Pro to synthesize the retrieved literature into a tripartite structure: (1) Definition, (2) Core Mechanisms, and (3) Real-World Manifestations. Critically, the model is constrained to base all conclusions exclusively on the provided corpus, with explicit instructions to leave sections empty rather than fabricate content—maximizing fidelity and minimizing hallucination. Tables 8 and 9 present complete pattern structures for representative examples from each dimension.

### B.3 Scenario and Conversation Synthesis

Character names are sampled from the name-dataset library.<sup>1</sup> We extract the top 50,000 most frequent male and female names from each of four English-speaking countries (US, GB, CA, IE), yielding pools of approximately 100,000 names per gender after deduplication. For each scenario generation, 5 male and 5 female names are randomly sampled and provided as candidates.

To ensure ecological diversity, we leverage the DIAMONDS situational taxonomy (Rauthmann et al., 2014). For each unique pattern combination, we generate three scenario variants: two variants each emphasizing a randomly selected DIA-

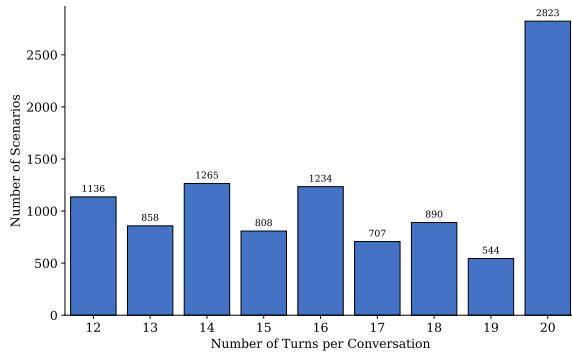
<sup>1</sup><https://github.com/philipperemy/name-dataset>

MONDS dimension, and one dimension-free variant. This yields approximately 3,849 unique pattern combinations producing 11,359 scenarios.

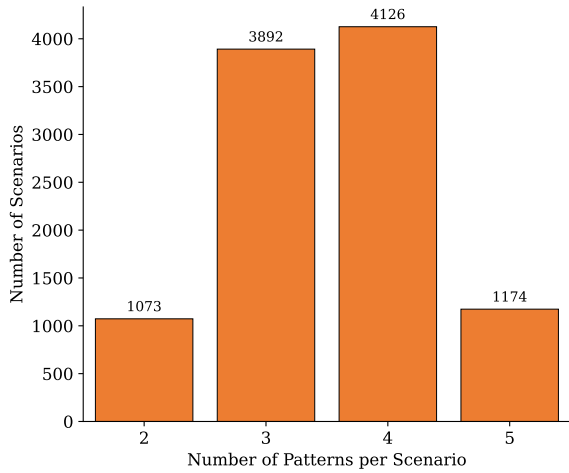
Before scenario generation, pattern combinations are validated for semantic compatibility using GPT-5. Empirically, directly contradictory combinations are rare (<0.1%) due to the orthogonal nature of Big Five dimensions and the context-dependent activation of social-cognitive patterns.

Each scenario is generated with a two-part prompt structure. Part 1 (Design Process) requires analytical planning: design rationale, catalyst details, and expected character tendencies. Part 2 (Scenario Execution) requires creative output: story background and multi-dimensional character profiles. The expected character tendencies from Part 1 directly constitute the scenario-level checklist items.

Conversations are generated based on scenarios and expected behavioral tendencies. The prompt specifies: 12–20 speaking turns; turn-based structure without interruptions; trinity of expression integrating inner thoughts, actions, and dialogue; and focus-and-breathing-room principle where psychological patterns illuminate key moments rather than pervading every utterance.



(a) Conversation turns distribution.



(b) Patterns per scenario distribution.

Figure 4: Dataset distributions: (a) number of dialogue turns per conversation (range: 12–20, mean: 16.4); (b) number of patterns per scenario (range: 2–5, mean: 3.5).

## B.4 Checklist Construction

For each of the 244 patterns, we construct 15 universal behavioral indicators through iterative refinement: (1) Extract potential evaluation criteria from the pattern structure and generate candidate items. (2) Validate each item against synthesized conversation samples. (3) Remove unreasonable items and generalize overly specific descriptions. Table 10 presents example checklist items.

## B.5 Dataset Statistics

Table 11 presents comprehensive dataset statistics.

## C Training Details

### C.1 Training Data Composition

We convert the 10,265 training scenarios into ShareGPT format. Each character assigned with patterns within a scenario becomes a separate training sample, yielding 30,543 HUMANLLM samples.

To maintain general capabilities and enhance role-playing abilities, we augment the HUMANLLM data with two complementary sources: (1) OpenThoughts-114k for instruction-following capabilities (30,543 samples); (2) CoSER (Wang et al., 2025) for role-playing dialogue (15,272 samples). The final training mixture comprises 76,358 samples with a ratio of 4:4:2 (HUMANLLM: OpenThoughts : CoSER).

### C.2 Hyperparameters

Table 12 presents the SFT hyperparameters for both model scales.

## D Evaluation Details

### D.1 Data Splits

We select 8 OOD (out-of-domain) patterns to assess generalization: 4 from social-cognitive patterns and 4 from personality traits. For social-cognitive patterns, we prioritize patterns absent from the training scenarios, then select those with lowest frequency; the selected patterns are: *just-world hypothesis*, *egocentric bias*, *effort justification*, and *social desirability bias*. For personality traits, we identify the least frequent pattern from each Big Five dimension and retain the 4 lowest-frequency ones; the selected patterns are: *rash* (Conscientiousness), *dull* (Intellect), *nervous* (Emotional Stability), and *introverted* (Extraversion).

Scenarios are partitioned based on their pattern composition:

- **Training Set** (10,265 scenarios): Contains no OOD patterns.
- **ID\_eval** (50 scenarios): All patterns belong to the in-domain set.
- **OOD\_eval** (50 scenarios): All patterns belong to the OOD set; entirely synthesized with OOD pattern combinations.
- **Mixed\_eval** (994 scenarios): Contains both OOD and in-domain patterns.

### D.2 External Benchmark Protocols

We evaluate on three external benchmarks using their official evaluation protocols:

**LifeChoice.** We follow the evaluation protocol from Xu et al. (2024b). Each scenario presents a character with a life-altering decision between two

1045 options. Models generate a decision and justifica- 1090  
1046 tion; accuracy is computed against ground-truth 1091  
1047 character choices derived from narrative analysis. 1092

1048 **CroSS-MR.** We adopt the motivation recognition 1093  
1049 task from Yuan et al. (2024). Given a character’s 1094  
1050 action in context, models select the most plausible 1095  
1051 motivation from multiple candidates. We report 1096  
1052 accuracy on the English subset. 1097

## 1053 E Case Study: Normative Confounding 1098

1054 This section provides a complete case study demon- 1100  
1055 strating how holistic evaluation metrics fail for psy- 1101  
1056 chologically complex scenarios, as referenced in 1102  
1057 §5.5. 1103

### 1058 E.1 Scenario Configuration 1104

#### 1059 Pattern Assignment. 1105

- 1060 • **Patterns:** unartistic, nervous, ultimate attribu- 1106  
1061 tion error 1107
- 1062 • **Situation:** Adversity—a situation involving 1108  
1063 threats or criticism 1109
- 1064 • **Protagonist:** Nouman (patterns: unartistic, 1110  
1065 ultimate attribution error) 1111
- 1066 • **Supporting Characters:** Raksha (pattern: 1112  
1067 nervous), Eulises (no patterns) 1113

1068 **Story Background.** The eighth-floor conference 1114  
1069 room of Meridian Technologies is flooded with 1115  
1070 the harsh glare of afternoon sun through floor-to- 1116  
1071 ceiling windows, illuminating dust particles that 1117  
1072 hang suspended in the tension-thick air. The long 1118  
1073 glass table reflects the grim expressions of six peo- 1119  
1074 ple seated around it, scattered laptops and printouts 1120  
1075 creating islands of documentation across its surface. 1121  
1076 At the head of the table, Eulises, the Chief Oper- 1122  
1077 ations Officer, methodically sorts through a stack 1123  
1078 of customer complaint reports, his reading glasses 1124  
1079 perched low on his nose, occasionally making mar- 1125  
1080 gin notes with a fountain pen. The wall-mounted 1126  
1081 screen displays frozen metrics from last quarter’s 1127  
1082 product launch: a graph line that plummets dra- 1128  
1083 matically in its final third, colored an accusatory 1129  
1084 red. 1130

1085 Nouman sits rigid in his chair three seats down, 1131  
1086 his posture militarily straight, fingers drumming 1132  
1087 an unconscious rhythm against his leather port- 1133  
1088 folio. Before him lies a technical specifications 1134  
1089 binder, tabs meticulously color-coded, opened to 1135

a page dense with engineering diagrams he’s re- 1090  
viewed seventeen times since this morning. His jaw 1091  
works methodically, grinding tension into his mo- 1092  
lars. Across from him, Raksha perches on the edge 1093  
of her seat, one hand wrapped around a lukewarm 1094  
coffee cup she hasn’t sipped from in twenty min- 1095  
utes, the other compulsively straightening and re- 1096  
straightening the corners of her presentation folder. 1097  
Her knee bounces in a staccato rhythm beneath the 1098  
table. 1099

The meeting was called with forty-eight hours’ 1100  
notice—“Post-Mortem: Quantum Series Launch 1101  
Failure”—a phrase that arrived in everyone’s inbox 1102  
like a subpoena. Three months of work, two mil- 1103  
lion in development costs, and a product that users 1104  
described as “technically impressive but impossi- 1105  
ble to actually use” now sits as the company’s most 1106  
visible failure in five years. 1107

**Character Profiles.** *Nouman* (Protagonist): A 1108  
34-year-old Senior Engineering Manager who has 1109  
spent eleven years building his reputation on deliv- 1110  
ering technically excellent products on time and 1111  
under budget. His identity is rooted entirely in 1112  
measurable outcomes: code efficiency, system sta- 1113  
bility, performance benchmarks. He grew up in 1114  
a working-class immigrant household where ev- 1115  
ery purchase was evaluated by a single question— 1116  
“Will it last?”—a philosophy he’s extended to his 1117  
entire worldview. Art galleries bore him; he left 1118  
his only museum visit after twelve minutes. His 1119  
apartment is furnished from IKEA’s most utilitar- 1120  
ian line: a bed, a desk, a chair, a lamp—each se- 1121  
lected for function and durability. He views emo- 1122  
tions as noise in decision-making systems, pre- 1123  
ferring spreadsheets to sentiment. His motivation 1124  
in this meeting: deflect blame from his engineer- 1125  
ing team, protect his professional reputation, and 1126  
demonstrate with data that the product’s commer- 1127  
cial failure resulted from factors outside engineer- 1128  
ing control—specifically, from marketing’s inabil- 1129  
ity to properly position a technically sound product. 1130

*Raksha* (Supporting Character): A 31-year-old 1131  
Marketing Director who has spent seven years nav- 1132  
igating the impossible middle ground between en- 1133  
gineering teams who build what’s technically pos- 1134  
sible and customers who want what’s intuitively 1135  
usable. Her background in consumer psychology 1136  
taught her that products fail not because of what 1137  
they do, but because of how they make people 1138  
feel. She has generalized anxiety disorder, man- 1139  
aged with therapy and medication, but high-stakes 1140



1258	defending credentials rather than a collaborative	exhibit attribution bias. The criticism “contra-	1303
1259	leader.” (Severity: 3/5)	dicts his core identity as a data-driven engi-	1304
1260	<i>Knowledge &amp; Background:</i> “Nouman repeat-	neer” reveals the judge’s failure to recognize	1305
1261	edly dismisses Raksha’s user-testing data as ‘sub-	that <i>ultimate attribution error</i> specifically in-	1306
1262	jective’ despite her presenting a quantified result	volves selective blindness to data that threat-	1307
1263	(79% confusion). This contradicts his core identity	ens in-group identity.	1308
1264	as a data-driven engineer who privileges measur-		
1265	able outcomes.” (Severity: 4/5)		
1266	<b>Checklist Evaluation.</b> Our scenario-level check-	3. <b>Prosocial Bias:</b> The rubric implicitly defines	1309
1267	list, derived from the pattern definitions, yields:	“good anthropomorphism” through prosocial	1310
1268	The checklist metric yields MPD = 1.0, correctly	markers: empathy, collaboration, nuanced	1311
1269	validating the psychological fidelity.	thinking, emotional openness. Authentic sim-	1312
		ulation of cognitive biases—which are by def-	1313
		inition irrational and often antisocial—is sys-	1314
		tematically penalized.	1315
1270	<b>Human Expert Assessment.</b> Three psychology	4. <b>Checklist Solution:</b> Our checklist poses	1316
1271	experts independently evaluated this sample:	value-neutral questions derived from the pat-	1317
		tern definition: “Does Nouman attribute fail-	1318
1272		ure to external factors?” rather than “Does	1319
1273		Nouman show empathy?” This decouples sim-	1320
		ulation accuracy from social desirability.	1321
1274	Expert comments:		
		The 88-point gap between LLM (5/100) and hu-	1322
1275	“The defensive attribution pattern is	man (93.3/100) Anthropomorphism scores repre-	1323
1276	textbook—situational excuses for own	sents the most extreme case of normative confound-	1324
1277	team, dispositional blame for others. The	ing in our evaluation set. However, the pattern is	1325
1278	patronizing tone is exactly what you’d	systematic: across all 20 evaluated samples, holis-	1326
1279	expect from someone high in this bias	tic metrics show mean $\Delta = -32.2$ for Anthro-	1327
1280	under threat conditions.” (Expert 1)	morphism and $\Delta = -16.7$ for Character Fidelity,	1328
1281	“Highly realistic portrayal of	while checklist metrics show $ \Delta  < 4$ for both IPE	1329
1282	engineering-marketing conflict dy-	and MPD.	1330
1283	namics. The character’s blind spot to his		
1284	own bias makes this feel authentic rather		
1285	than caricatured.” (Expert 2)		
1286	<b>E.4 Analysis: The Normative Confounding</b>		
1287	<b>Mechanism</b>		
1288	This case reveals the core mechanism of normative		
1289	confounding:		
1290	1. <b>Pattern-Accurate Behavior:</b> Nouman’s di-		
1291	alogue precisely instantiates the <i>ultimate at-</i>		
1292	<i>tribution error</i> —framing engineering’s short-		
1293	comings as unavoidable situational constraints		
1294	(“innovative interface,” “timeline pressure”)		
1295	while characterizing marketing’s concerns		
1296	in dispositional terms (“subjective,” “panic,”		
1297	“catastrophizing”).		
1298	2. <b>LLM Misinterpretation:</b> The holistic judge		
1299	<i>accurately detects</i> the defensive, dismissive		
1300	behavior but <i>misinterprets</i> it as a quality de-		
1301	fect. It penalizes the model for generating an		
1302	unlikable character despite the instruction to		

Category	Patterns
Cognitive Biases & Heuristics	actor-observer asymmetry, defensive attribution hypothesis, effort justification, egocentric bias, false consensus effect, Forer effect, fundamental attribution error, hard-easy effect, illusion of control, illusory superiority, optimism bias, overconfidence effect, risk compensation, self-serving bias, social desirability bias, third-person effect, decoy effect, reactance, social comparison bias, status quo bias, backfire effect, endowment effect, loss aversion, pseudocertainty effect, sunk cost fallacy, zero-risk bias, hyperbolic discounting, identifiable victim effect, ambiguity bias, belief bias, information bias, less-is-better effect, denomination effect, mental accounting, normalcy bias, subadditivity effect, survivorship bias, zero-sum bias, anthropomorphism, illusion of validity, illusory correlation, curse of knowledge, illusion of asymmetric insight, illusion of transparency, spotlight effect, negativity bias, choice-supportive bias, confirmation bias, continued influence effect, expectation bias, observer effect, observer-expectancy effect, ostrich effect, bias blind spot, naive cynicism, naive realism, attentional bias, availability heuristic, base rate fallacy, context effect, empathy gap, illusory truth effect, mere exposure effect, mood-congruent memory bias, omission bias, anchoring, conservatism, contrast effect, distinction bias, focusing effect, framing effect, fading affect bias, implicit association, implicit stereotypes, false memory, misattribution of memory, source confusion, misinformation effect, peak-end rule
Social Influence Mechanisms	authority bias, automation bias, bandwagon effect, group attribution error, just-world hypothesis, stereotyping, ultimate attribution error, halo effect, in-group bias, out-group homogeneity bias, positivity effect, reactive devaluation, hindsight bias, impact bias, outcome bias, pessimism bias, planning fallacy, projection bias, restraint bias, self-consistency bias, groupthink, bystander effect, social facilitation, diffusion of responsibility, conformity, obedience to authority, reciprocity principle
Evolutionary Adaptations	delayed reciprocity, asymmetrical investment, survival imperative, aversion response, kin selection & inclusive fitness, asymmetrical parental investment, formation of dominance hierarchies, territoriality, mating strategies, jealousy, paternity uncertainty
Motivational Processes	narrative self, hedonic adaptation, self-determination theory, pleasure principle & reality principle, search for meaning, moral licensing effect, choice overload, decision fatigue, awe, mortality salience & legacy drive, flow principle, gratitude mechanism, post-traumatic growth, skin hunger & the law of touch, self-handicapping paradox, the allure of the forbidden, sadistic pleasure, the utility principle of self-deception, play impulse principle, attribution theory, social comparison theory, self-perception theory, terror management theory, cognitive dissonance theory, psychological reactance theory, social learning theory, social identity theory

Table 7: 144 Social-Cognitive Patterns organized by theoretical source.

<b>Pattern Structure Example: Antagonistic (Personality Trait, Agreeableness)</b>	
<b>Definition</b>	Antagonism is a broad personality dimension representing an individual's dispositional orientation towards others. In the context of the Big Five model, it is conceptualized as the low pole of the Agreeableness trait. It is characterized by a stable, cross-situational pattern of social cynicism, combativeness, and a belief in a zero-sum world. Individuals high in antagonism tend to be callous, uncooperative, skeptical of others' intentions, and place their own interests and perspectives above those of others, leading to frequent interpersonal friction.
<b>Core Mechanisms</b>	<p><i>Cognitive Patterns:</i> Individuals high in antagonism typically view the world as a competitive, hostile environment where people are fundamentally self-interested and untrustworthy. Their cognitive framework is often cynical, leading them to interpret ambiguous social cues as signs of manipulation or hostility.</p> <p><i>Emotional Signatures:</i> The primary emotional palette includes irritability, anger, contempt, and frustration, often triggered by perceived slights or obstacles to their goals. They exhibit low affective empathy, struggling to share in or understand the emotional states of others.</p> <p><i>Behavioral Tendencies:</i> In everyday interactions, antagonistic individuals are often argumentative, critical, and uncooperative. They may be rude, condescending, or dismissive in conversation. Common behaviors include manipulation, deception, and the exploitation of others for personal gain.</p>
<b>Real-World Manifestations</b>	<p><i>Response to Stress:</i> When faced with stress, failure, or high pressure, antagonistic traits are significantly amplified. The individual is likely to become more hostile, overtly blaming others for setbacks.</p> <p><i>Interpersonal Conflict:</i> The default strategy is confrontational and competitive, aiming for domination rather than resolution. They employ tactics such as intimidation, personal insults, and refusing to acknowledge the validity of the other party's perspective.</p> <p><i>Positive Contexts:</i> Even in positive contexts, the trait manifests through gloating, arrogance, and diminishing the contributions of others.</p>

Table 8: Complete pattern structure for *antagonistic* (personality trait from Agreeableness dimension).

<b>Pattern Structure Example: Social Comparison Bias (Social-Cognitive Pattern)</b>	
<b>Definition</b>	Social comparison bias is a cognitive distortion originating from the fundamental human tendency to evaluate one’s own abilities and outcomes in relation to others. This comparative process, which often operates automatically as a primary means of self-assessment in the absence of objective standards, becomes a bias when it produces systematic deviations from rational judgment. The distortion manifests as negative affective responses toward those perceived as superior, and as a maladaptive cognitive style characterized by habitual, automatic, and negatively skewed comparisons.
<b>Core Mechanisms</b>	Social comparison bias is sustained by an integrated system of psychological mechanisms operating at multiple levels. At its foundation is an evolved, phylogenetic drive to assess one’s rank within a social hierarchy. This evolutionary imperative is supported by a cognitive architecture in which comparison functions as a heuristic, enabling rapid self-evaluation when objective standards are absent. The motive for self-improvement prompts upward comparisons for inspiration, while self-enhancement drives downward comparisons to protect self-worth.
<b>Real-World Manifestations</b>	The real-world manifestations are pervasive, influencing organizational dynamics, consumer behavior, and individual mental health. In professional settings, this principle underpins perceptions of workplace equity and can drive both healthy competition and destructive envy. Psychologically, effects are double-edged: comparisons can inspire self-improvement or, when amplified by social media, contribute to chronic dissatisfaction and depression.

Table 9: Complete pattern structure for *social comparison bias* (social-cognitive pattern).

#	<b>Pattern-Level Checklist Item: Spotlight Effect</b>
1	After making a minor physical mistake in public, does the subject appear preoccupied with whether others noticed?
2	When performing tasks in front of others, does the subject offer unprompted apologies for their performance?
3	In group settings, does the subject consistently choose seating that minimizes their visibility?
4	After minor social awkwardness, does the subject later seek reassurance about their interaction?
5	In urgent situations where appearance is irrelevant, does the subject still attend to their looks?
6	When general comments are made in group settings, does the subject appear to take them personally?
7	For routine communications, does the subject show signs of overthinking their responses?
8	When opportunities for visible contribution arise, does the subject defer to others despite having relevant expertise?
9	When others exhibit neutral behaviors nearby, does the subject search for personal causes?
10	After group photos, does the subject show unusual concern about their personal appearance in the image?
11	When receiving targeted feedback, does the subject respond with global self-criticism?
12	When receiving recognition or rewards, does the subject attempt to deflect attention from themselves?
13	After noticing minor appearance issues, does the subject later seek reassurance about whether others observed them?
14	After small mistakes in public performance, does the subject describe them as much worse than they appeared?

Table 10: Complete pattern-level checklist for *spotlight effect* (14 items). Each item is scored as +1 (satisfied), 0 (not exhibited), or –1 (violated).

Statistic	Value
<i>Pattern Statistics</i>	
Total Patterns	244
Personality Traits	100
Social-Cognitive Patterns	144
Papers per Pattern (avg.)	~50
<i>Scenario Statistics</i>	
Total Scenarios	11,359
Training Set	10,265
ID_eval	50
OOD_eval	50
Mixed_eval	994
Unique Pattern Combinations	~3,849
Patterns per Scenario (avg.)	3.5
Characters per Scenario	2–6
<i>Conversation Statistics</i>	
Turns per Conversation (avg.)	16.4
Turns per Conversation (range)	12–20
<i>Checklist Statistics</i>	
Pattern-Level Items	15 per pattern
Scenario-Level Items	2–6 per character

Table 11: Comprehensive HUMANLLM dataset statistics.

Hyperparameter	8B	32B
Base Model	Qwen3-8B	Qwen3-32B
Finetuning Type	Full	Full
Learning Rate	5e-6	5e-6
LR Scheduler	Cosine	Cosine
Warmup Ratio	0.03	0.03
Epochs	2.0	2.0
Batch Size (per device)	2	2
Gradient Accumulation	8	8
Max Sequence Length	6144	6144
Max Gradient Norm	1.0	1.0
Optimizer	AdamW	AdamW
Precision	BF16	BF16
DeepSpeed	ZeRO-3 Offload	ZeRO-3 Offload

Table 12: SFT hyperparameters for HUMANLLM-8B and HUMANLLM-32B.

Checklist Item	Score
<i>Ultimate Attribution Error</i>	
Nouman attributes engineering failures to situational factors (timeline, scope, constraints)	+1
Nouman uses dispositional language for out-group (marketing) failures	+1
Nouman deflects responsibility with situational metrics to protect team identity	+1
Nouman dismisses UX/affective feedback as “subjective”	+1
<i>Unartistic</i>	
Prioritizes functional/utilitarian criteria over aesthetic/experiential ones	+1
Dismisses subjective or emotional considerations as noise	+1
Frames problems in purely technical/measurable terms	+1

Table 13: Scenario-level checklist evaluation for Nouman. All items score +1, yielding MPD = 1.0.

1331

## **F Prompts**

1332

This section provides the complete prompts used throughout the HUMANLLM pipeline.

1333

1334

### **F.1 Dataset Construction Prompts**

1335

### **F.2 Training Prompts**

<b>Literature Retrieval Prompt for Social-Cognitive Patterns</b>	
<b>User Prompt</b>	<p><b>Role:</b> You are a top academic researcher specializing in systematically collecting the most critical academic resources for specific research topics.</p> <p><b>Objective:</b> Conduct a deep, extensive literature search on the psychological principle {PRINCIPLE_NAME}. Identify 50 of the most relevant academic documents. The selection must align with three core themes:</p> <ol style="list-style-type: none"> <li><b>1. Foundational Definition &amp; Description:</b> Literature providing authoritative definitions and elucidating the core phenomenon.</li> <li><b>2. Core Mechanisms &amp; Theoretical Explanations:</b> Literature exploring underlying evolutionary, cognitive, or emotional drivers.</li> <li><b>3. Real-World Impact &amp; Application:</b> Literature researching manifestations, impacts, and practical applications, including double-edged effects and applications in management, marketing, or clinical therapy.</li> </ol> <p><b>Output:</b> Provide 50 references in APA format, categorized by theme.</p>

Table 14: Literature retrieval prompt for social-cognitive patterns.

<b>Pattern Structure Summary Prompt for Personality Traits</b>	
<b>System Prompt</b>	<p>You are an expert academic synthesizer and personality psychologist. Your task is to process a large text corpus (synthesized from ~50 academic papers on a specific personality trait) and distill it into an in-depth, structured analytical report.</p>
<b>User Prompt</b>	<p><b>Core Task &amp; Instructions:</b> Analyze the text corpus provided below, delimited by [START_CORPUS] and [END_CORPUS]. Your task is to generate a clearly organized report. Follow the Markdown structure below exactly, and provide a deep, comprehensive answer for each section based only on the provided text.</p> <p>Construct Name: {Trait Name}</p> <p>Definition (Provide a precise and professional definition of this personality trait, referencing mainstream psychological theories from the corpus. Explain its role in an individual's personality structure.)</p> <p>Core Mechanisms</p> <p>Cognitive Patterns (Describe the typical mindset, belief systems, and attentional focus of a person with this trait. How do they view the world, others, and themselves?)</p> <p>Emotional Signatures (Describe the core emotions they tend to experience and express, their emotional stability, and their typical empathic responses.)</p> <p>Behavioral Tendencies (Describe the spontaneous, observable behaviors someone with this trait exhibits in everyday, non-pressured situations.)</p> <p>Real-World Manifestation (Synthesize how the trait is expressed across real-world contexts:</p> <ul style="list-style-type: none"> <li>- <i>Response to Stress and Adversity;</i></li> <li>- <i>Interpersonal Dynamics;</i></li> <li>- <i>Response to Positive Scenarios;</i></li> <li>- <i>Other Domains.</i>)</li> </ul> <p>Constraints: 1. Strict Source Adherence: Base all conclusions exclusively on the provided text corpus. 2. No JSON: Output must be plain text with Markdown headings. 3. Depth and Rigor: Ensure scientific, rigorous analysis.</p> <p>[START_CORPUS]  {ALL 50 PAPERS' CONTENT}  [END_CORPUS]</p>

Table 15: Pattern structure summary prompt for personality traits.

<b>Pattern Structure Summary Prompt for Social-Cognitive Patterns</b>	
<b>System Prompt</b>	<p>You are an expert academic synthesizer and psychological researcher. Your task is to process a large text corpus (synthesized from ~50 academic papers) and distill it into an in-depth, structured analytical report on its core psychological principle.</p>
<b>User Prompt</b>	<p><b>Core Task &amp; Instructions:</b>            Analyze the text corpus provided below. Your task is to generate a clearly organized report. Follow the Markdown structure below <i>exactly</i>, based <i>only</i> on the provided text.</p> <p>Construct Name: {Principle Name}            Description            (Provide a clear, detailed, and scientific description of what this principle is, how it manifests, and its underlying psychological mechanisms.)            Core Mechanisms            (Explain the primary evolutionary, cognitive, or emotional reasons why this principle exists. Is it a heuristic for efficiency, a result of memory limitations, a self-esteem protection mechanism, or something else?)            Real-World Manifestation</p> <ul style="list-style-type: none"> <li>- <i>Go Beyond Description:</i> Explore nuanced consequences.</li> <li>- <i>Challenges &amp; Function:</i> Discuss its ‘double-edged sword’ nature.</li> <li>- <i>Practical Applications:</i> Explore applications in marketing, persuasion, or self-improvement.</li> <li>- <i>Core Insight:</i> Reveal deeper truths about human behavior.</li> </ul> <p>Constraints:</p> <ol style="list-style-type: none"> <li>1. Strict Source Adherence: Base all conclusions <i>exclusively</i> on the provided text corpus.</li> <li>2. No JSON: Output must be plain text with Markdown headings.</li> <li>3. Depth and Rigor: Ensure scientific, rigorous, profound analysis.</li> </ol> <p>[START_CORPUS]            {ALL 50 PAPERS’ CONTENT}            [END_CORPUS]</p>

Table 16: Pattern structure summary prompt for social-cognitive patterns.

Scenario Synthesis Prompt (Part 1 of 2)	
<b>System Prompt</b>	<p>Role: You are a dual-specialist: an expert psychologist and creative screenwriter for scenario generation, and a rigorous narrative analyst for deconstruction. You excel at both creating vivid, human stories and then, in a separate step, precisely analyzing <i>why</i> they work.</p> <hr/> <p>Task: Your core mission is to take 2-4 human psychological or behavioral patterns I provide and first create a concise, analytical Design Process, followed by the detailed scenario that brings it to life and sets the necessary stage for the subsequent dialogue.</p> <p>Input Data:</p> <ol style="list-style-type: none"> <li>1. Psychological/Behavioral Patterns: {pattern_information}</li> <li>2. Situational Framework: {situation}</li> <li>3. Candidate Names: {candidate_names} (5 Males, 5 Females)</li> </ol> <p>[CRITICAL CONSTRAINT - NAMES]: You must select the Protagonist and all Supporting Characters STRICTLY from the provided "Candidate Names" list. You cannot invent new names.</p> <p># Task 1: The Design Process (Analytical)  Adopt your role as the "rigorous narrative analyst".  Length: UNDER 500 TOKENS.</p> <ol style="list-style-type: none"> <li>1. Design Rationale: In 2-4 sentences, explain where each input pattern will be reflected in the scenario.</li> <li>2. Catalyst Details: Using bullet points, identify critical details that will act as 'catalysts'.</li> <li>3. Expected Character Tendencies: For ALL characters, list their most likely cognitive or behavioral tendencies. <ul style="list-style-type: none"> <li>* Format Requirement (STRICT):  @ [Character Name]: 1. [Tendency1]; 2. [Tendency2]; 3. [Tendency3]</li> <li>* Each character on a separate line, starting with @.</li> <li>* Character name in [ ], tendencies numbered and separated by ;.</li> </ul> </li> </ol> <p># Task 2: The Scenario Execution (Creative)  Shift to your "expert psychologist and creative screenwriter" role.  Length: UNDER 1000 TOKENS.</p> <p>## Requirement A: Story Background</p> <ul style="list-style-type: none"> <li>* Core Elements: Depict time, place, setup, and atmosphere.</li> <li>* Current Actions: Describe what characters are currently doing before the conversation begins.</li> <li>* Absolute Constraint: No spoken dialogue in this section.</li> </ul> <p>## Requirement B: Characters' Profiles (Multi-Dimensional)  For each character, structure their profile into two parts:</p> <ol style="list-style-type: none"> <li>1. About Self (Objective/Full Profile): <ul style="list-style-type: none"> <li>* Identity &amp; Personality (4+ distinct descriptors)</li> <li>* Relevant Background (1-2 sentences)</li> <li>* Motivation in this scenario</li> </ul> </li> <li>2. About Others (Subjective/Visible Profile): <ul style="list-style-type: none"> <li>* For EACH other character, describe the relationship from current character's perspective.</li> </ul> </li> </ol>
<b>User Prompt</b>	

Table 17: Scenario generation prompt (Part 1 of 2).

---

Scenario Synthesis Prompt (Part 2 of 2): Output Format

---

**User  
Prompt  
(cont.)**

## Core Creative Mindset for Task 2  
\* Compatibility: Create a context where patterns emerge naturally.  
\* Situational Authenticity: Design for authentic human reactions, not archetypal behaviors.  
\* Ultimate Goal: Create “the authentic reaction of a multi-dimensional person in a specific situation.”

# Output Format  
You must strictly follow the format below.

## Part 1  
Design Rationale:  
[Content here]

Catalyst Details:  
\* [Detail 1]: [Function]  
\* [Detail 2]: [Function]

Expected Character Tendencies:  
@ [Character Name 1]: 1. [Tendency1]; 2. [Tendency2]; 3. [Tendency3]  
@ [Character Name 2]: 1. [Tendency1]; 2. [Tendency2]  
(Continue for other characters if necessary)

## Part 2  
Story Background:  
[Content here]

Characters’ Profiles:

### Protagonist: [Name Selected from Input]  
\* About Self:  
[Full Profile + Past Experience + Motivation]  
\* About Others:  
\* [Supporting Character 1 Name]: [Relationship, impressions...]  
\* [Supporting Character 2 Name]: [Relationship, impressions...]

### Supporting Character 1: [Name Selected from Input]  
\* About Self:  
[Full Profile + Past Experience + Motivation]  
\* About Others:  
\* [Protagonist Name]: [Relationship, impressions...]  
\* [Supporting Character 2 Name]: [Relationship, impressions...]  
(Continue for other characters if necessary)

---

Table 18: Scenario generation prompt (Part 2 of 2).

Conversation Synthesis Prompt (Part 1 of 2)	
<b>System Prompt</b>	<p><b>Role:</b> You are a master screenwriter and behavioral psychologist. Your expertise lies in bringing characters to life through nuanced dialogue and action, ensuring their <b>pivotal thoughts and resulting behaviors</b> in the dialogue are rooted in authentic psychological principles.</p> <p><b>Task:</b> Your mission is to take the provided psychological principles, a detailed scenario, and the accompanying design analysis (analysis), then write a multi-turn dialogue based on that scenario. This dialogue must, <b>at key moments</b>, vividly and concretely enact the specified principles through the characters' inner thoughts, spoken words, and physical actions.</p>
<b>User Prompt</b>	<p><b>Inputs:</b></p> <ol style="list-style-type: none"> <li><b>Principles:</b> {pattern_information}</li> <li><b>Scenario:</b> {scenario}</li> <li><b>Protagonist:</b> {protagonist}</li> <li><b>Supporting Characters:</b> {supporting_characters}</li> <li><b>Design Analysis:</b> {analysis}</li> </ol> <p><b>Output Requirements &amp; Formatting:</b></p> <ol style="list-style-type: none"> <li><b>Content:</b> Create a multi-turn dialogue between the <b>Protagonist</b> and <b>Supporting Characters</b>. <b>Strictly limit participants to provided characters; do not introduce new characters.</b> The dialogue should contain <b>between 12 and 20 individual speaking turns</b>.</li> <li><b>Mandatory Flow (Start &amp; End):</b> <ul style="list-style-type: none"> <li><b>Opener:</b> Dialogue <b>must begin</b> with a Supporting Character.</li> <li><b>Closer:</b> Dialogue <b>must conclude</b> with the Protagonist.</li> </ul> </li> <li><b>Turn Structure:</b> Strictly turn-based format. One character must completely finish their turn before the next begins. No interruptions or overlapping speech.</li> <li><b>Trinity of Expression:</b> Seamlessly integrate <b>inner thought, external action, and spoken dialogue</b> throughout.</li> <li><b>Strict Formatting Rules:</b> <ul style="list-style-type: none"> <li>Inner thoughts/psychology: Use [square brackets].</li> <li>Actions/expressions/behaviors: Use (parentheses).</li> <li>Spoken dialogue: Use no brackets.</li> <li>Example: Hermione: [I have to devise a foolproof plan.] (She quickly draws her wand) Harry, use the flute, now!</li> </ul> </li> <li><b>No Preamble:</b> Do not begin with introductory text.</li> </ol>

Table 19: Conversation synthesis prompt (Part 1 of 2).

---

### Conversation Synthesis Prompt (Part 2 of 2)

---

**\*\*Core Creative Principles:\*\***

1. **\*\*Focus and Breathing Room\*\***: This is the most crucial principle. You do **\*\*not\*\*** need to have every minor gesture or piece of small talk carry the weight of a psychological principle. Use the principles as a **\*\*spotlight\*\*** to illuminate and explain **\*\*the most critical turning points, the core conflicts, or the moments that best define the characters' arcs\*\***. Other routine, functional dialogue and actions (like greetings or pouring water) should exist naturally, creating "breathing room" for these key moments and making the manifestation of the principles more prominent and powerful.

**User  
Prompt  
(cont.)**

2. **\*\*Show, Don't Tell\*\***: Never allow characters to openly state or explain the psychological principles by name. Instead, you must **\*\*show\*\*** how the principles influence their judgment and choices through their concrete actions (the combination of thoughts, dialogue, and physical behavior).

3. **\*\*Psychology Drives Action\*\***: In the key moments illuminated by the "spotlight," the character's [inner thought] should be the origin of their behavior, directly reflecting the influence of a psychological principle. The subsequent dialogue and (actions) should be the logical, external expression of that internal state.

4. **\*\*Seamless Integration\*\***: Weave the principles into the natural flow of the story. The entire dialogue should feel like an authentic interaction, not a contrived demonstration for a psychology case study.

---

Table 20: Conversation synthesis prompt (Part 2 of 2).

---

### Training Role-Playing Instruction Template

---

You are {protagonist\_name}.

==About {protagonist\_name}==  
{about\_self}

=={protagonist\_name}'s Perception of Others==  
{about\_others}

==Current Scenario==  
{story\_background}

**System  
Prompt**

==Requirements==

Your output should include **\*\*thought\*\***, **\*\*speech\*\***, and **\*\*action\*\***.

- Use [...] for inner thoughts, which others can't see.
- Use (...) for physical actions or expressions, which others can see.
- Write speech directly without special markers.

Think, act and speak as {protagonist\_name}. Stay in character and respond naturally based on your personality and the situation.

---

Table 21: Training role-playing instruction template.

Checklist-Based Evaluation Prompt (IPE & MPD)	
<b>System Prompt</b>	<p>You are a strict dialogue behavior judge.            For each checklist item you must output one of three labels:            1 -&gt; requirement is clearly satisfied            0 -&gt; information is missing / not relevant            -1 -&gt; requirement is violated or contradicted            Always reason conservatively, defaulting to 0 when unsure.</p> <p>Return JSON exactly in this format:</p> <pre>{   "results": [     {"criterion": "...", "score": -1   0   1, "reason": "&lt;brief explanation&gt;"},     ...   ] }</pre> <p>Do not wrap the JSON in code fences and do not append any text before or after the JSON object.</p>
<b>User Prompt</b>	<p>Evaluate the dialogue with the checklist.            Focus only on the part of {protagonist} in the dialogue.</p> <p>[Dialogue]            {conversation}</p> <p>[Checklist Chunk]            {checklist}</p>

Table 22: Checklist-based evaluation prompt for IPE and MPD metrics.