WHEN TO ACT, WHEN TO WAIT: Modeling Structural Trajectories for Intent Triggerability in Task-Oriented Dialogue

Anonymous authors Paper under double-blind review

Abstract

1	Task-oriented dialogue systems often face difficulties when user utterances
2	seem semantically complete but lack necessary structural information for
3	appropriate system action. This arises because users frequently do not
4	fully understand their own needs, while systems require precise intent
5	definitions. Current LLM-based agents cannot effectively distinguish be-
6	tween linguistically complete and contextually triggerable expressions,
7	lacking frameworks for collaborative intent formation. We present STORM,
8	a framework modeling asymmetric information dynamics through conver-
9	sations between UserLLM (full internal access) and AgentLLM (observable
10	behavior only). STORM produces annotated corpora capturing expression
11	trajectories and latent cognitive transitions, enabling systematic analysis
12	of collaborative understanding development. Our contributions include:
13	(1) formalizing asymmetric information processing in dialogue systems; (2)
14	modeling intent formation tracking collaborative understanding evolution;
15	and (3) evaluation metrics measuring internal cognitive improvements
16	alongside task performance. Experiments across four language models
17	reveal that moderate uncertainty (40–60%) can outperform complete trans-
18	parency in certain scenarios, with model-specific patterns suggesting recon-
19	sideration of optimal information completeness in human-AI collaboration.
20	These findings contribute to understanding asymmetric reasoning dynam-
21	ics and inform uncertainty-calibrated dialogue system design.

22 1 Introduction

The rapid advancement of language models has created a fundamental challenge in human-23 AI interaction: the "gulf of envisioning"—users' cognitive difficulty in formulating effective 24 prompts. Unlike conventional interfaces with predictable affordances, language models 25 require users to simultaneously envision possibilities and their expressions, often lead-26 ing to communication breakdowns. This challenge arises from a misalignment between 27 human cognitive processes and the way systems interpret user intent. Subramonyam et 28 al. Subramonyam et al. (2023) illustrate that human intent formation involves a maturation 29 process characterized by progressive constraint resolution, fluctuating stability intervals, 30 and distinct structural signaling patterns. However, current evaluation methods are in-31 sufficient as they: 1) treat intent as binary rather than continuous, 2) lack frameworks for 32 temporal coherence, and 3) overlook structural signals within expressions. These structural 33 signals-including stylistic choices, implicit assumptions, and cultural markers reflect what 34 Wittgenstein Wittgenstein (1953) termed the contextual embeddedness of meaning within 35 particular "forms of life." Current systems cannot access these embedded contextual cues 36 that users unconsciously include in their expressions. These shortcomings constitute the 37 Intent-Action Alignment Problem, determining precisely when user expressions have reached 38 cognitive readiness for effective system action. To address this alignment problem, we pro-39 pose STORM (Structured Task-Oriented Representation Model), a framework that builds 40 on top of CAMEL-AI (Li et al., 2023) and conceptualizes user intent as evolving along a 41 continuous spectrum, modeling the iterative refinement of user expressions, and identifying 42



Figure 1: Overview of the STORM Framework

timing misalignments that lead to communication failures. The major contributions of this
 paper include:

(1) A dialogue generation pipeline using two language models—UserLLM and
AgentLLM—to simulate realistic conversations reflecting diverse user profiles and intent
progression. UserLLM generates user behavior conditioned on comprehensive profile data
and internal states, simulating authentic intent evolution, while AgentLLM responds based
solely on observable dialogue history. This asymmetric setup mirrors the realistic information gaps faced by AI systems, allowing targeted studies on agent adaptability to evolving
intent.

(2) A database-driven memory system that systematically tracks evolving user states (intent, emotion, satisfaction) within session-specific records. These records function as micro-databases documenting real-time intent maturation trajectories and are integrated into a global database for cross-session analysis. This structured memory approach captures the continuous nature of intent development, providing researchers with detailed, fine-grained data to study patterns across diverse interaction contexts.
 (3) A web-based dialogue visualization interface equipped with a clarity rating mechanism

(3) A web-based dialogue visualization interface equipped with a clarity rating mechanism
 was developed to provide an intuitive analysis of the evolution of user intent. This interface
 dynamically displays the refinement process of user intent, enabling researchers to assess
 the effectiveness of various agent response strategies visually. The tool facilitates rigorous
 quantitative analysis and comparison by quantifying the abstract cognitive progression into
 a standardized clarity metric. The interface is publicly accessible at https://v0-dialogue analysis-dashboard.vercel.app/.

The dialogues produced by **STORM** serve as valuable training data, enabling conversational 65 agents to better detect and adapt to different stages of intent formation. We evaluate our 66 framework through comparative analysis of agent responses across metrics, including user 67 satisfaction and response quality, demonstrating improved alignment with user cognitive 68 processes. Our experiments demonstrate that access to user profiles significantly enhances 69 **model performance** across all evaluated systems, with satisfaction scores increasing by 70 **15–40**% when profile information is available. We introduce a novel '*Clarify*' metric that mea-71 sures how effectively agents help users internally clarify their own intentions — assessed 72 through analysis of simulated user inner thoughts rather than external expressions. This ap-73 proach captures whether agent responses genuinely improve users' understanding of their 74 own needs, a crucial cognitive process often invisible in traditional dialogue evaluations. 75

	Notation	Symbol	Description		
	User Expression	$e_t \in \mathcal{E}$	User utterance at dialogue turn <i>t</i>		
ore Domains	Agent Response	$r_t \in \mathcal{R}$	Agent utterance at dialogue turn t		
	Hidden State	$h_t \in \mathcal{H}$	User's internal state at turn <i>t</i> (inner thoughts, emotion, satisfaction)		
	Task Domain	$ au \in \mathcal{T}$	Space of tasks from the Task Library (technology, healthcare, etc.)		
Ŭ	User Domain	$u \in \mathcal{U}$	Space of user profiles with their multi-dimensional attributes		
	Expression Domain	$e_t \in \mathcal{E}$	Space of user utterances with varying degrees of clarity		
	Response Domain	$r_t \in \mathcal{R}$	Space of agent responses to user expressions		
e	Base Profile [task-agnostic]	$\mathbf{b} = \{b_1,, b_n\}$	Demographic and personality factors (culture, decision style, etc.)		
ofil	Task Parameters	$\mathbf{t} = \{t_1,, t_m\}$	Task-specific attributes (domain, brand, priority features, etc.)		
ser Pr	Context Profile [task-agnostic]	$\mathbf{c} = \{c_1,, c_k\}$	User capabilities and constraints (time constraint, patience, etc.)		
Ď	Task Specifics	S	Predefined user preferences and constraints for a task instance $\boldsymbol{\tau}$		
	Difficulty Config	$\mathbf{d} = \{d_{style}, d_{length}, d_{content}, d_{tone}\}$	Difficulty level and associated dimensions		
	Uncertainty Level	$p \in \{0\%, 40\%, 60\%, 80\%\}$	Percentage of profile attributes masked as unknown		
nt	Agent Role	$\alpha(au) \in \mathcal{A}$	A general helpful assistant for task τ		
Age	Agent Directive	$\delta(au)$	Task-specific guidelines instructing agent behavior and goals		
rics	Intent Evolution	$\Delta_t(h)$	Change in intent clarity from turn $t - 1$ to t based on hidden states		
Лet	Clarity Rating	$C(r_t, h_t, h_{t+1})$	Measurement of how agent response improves intent clarity		
_	Performance Score	$E(C_1,,C_T)$	Aggregate measure of agent effectiveness across dialogue turns		
SS	UserLLM Function	$G_{\text{user}}(u, \mathcal{H}_{1:t-1}, \mathcal{E}_{1:t-1}, \mathcal{R}_{1:t-1}) \rightarrow (e_t, h_t)$	User generation with full dialogue history		
ő	AgentLLM Function	$G_{\text{agent}}(\alpha(\tau), \delta(\tau), \mathcal{E}_{1:t}, \mathcal{R}_{1:t-1}) \to r_t$	Agent generation with role, directive, and observable history		
Pr	Basic Augmentation	$A_1(\mathcal{E}_{1:T}, \mathcal{R}_{1:T}, \mathcal{H}_{1:T}, u) \to D$	Collection of dialogues with complete metadata		
<u>io</u>	Turn Analysis	$A_2(D) \rightarrow D^+$	Enhanced data with per-turn analysis		
rat	Summary Generation	$A_3(D^+) \rightarrow S$	Comprehensive dialogue summaries		
ene	RAG Enhancement	$\mathcal{R}(D, D^+, S) \to K$	Using augmented data as knowledge base		
Ũ	Prompt Refinement	$P(D, D^+, S) \to G'_{\text{agent}}$	Creating improved prompts from analysis		

Table 1: Summary of notations used in STORM INTERFACE.

Our analysis reveals distinct model characteristics with practical deployment implications:
 Claude maintains consistent satisfaction across varying profile completeness, Gemini

demonstrates robust performance under high uncertainty, while Llama achieves superior 78 intent clarification despite satisfaction trade-offs. Notably, we observe that moderate profile 79 uncertainty (40-60% unknown attributes) often outperforms complete information access, 80 suggesting that excessive profile information may lead to presumptive reasoning, while 81 moderate uncertainty encourages more exploratory interaction strategies that better support 82 users' evolving understanding of their own needs. This finding has implications for privacy-83 preserving design and bias mitigation in dialogue systems. These insights highlight the 84 tension between immediate satisfaction and cognitive alignment, providing empirical 85 guidance for uncertainty-aware dialogue system design. 86

87 2 Core Components

⁸⁸ We define the **STORM (Structured Task-Oriented Representation Model)** Interface as ⁸⁹ a formal framework for studying the relationship between user intent expression and ⁹⁰ system actionability. The **STORM** Interface is represented as a 5-tuple of domain spaces: ⁹¹ {T, U, E, R, H}. We define each component in detail as follows.

Task domain \mathcal{T} is defined as a collection of task objects $\tau \in \mathcal{T}$, where each τ comprises 92 a task name, a description, and domain-specific requirements. Prior approaches Yao et al. 93 (2024); Prabhakar et al. (2025) constrain task definitions to those with explicit success met-94 rics, whereas our formalization generalizes across task types. The task representation of 95 our framework is domain-agnostic, enabling automated attribute generation for arbitrary 96 domains beyond our experimental setup. The implementation accepts custom task defini-97 tions through a standardized interface that integrates with existing domain taxonomies and 98 classification systems. 99

User domain \mathcal{U} consists of user profiles $u \in \mathcal{U}$, where each profile is represented as a 100 vector of attribute-value pairs. These captures both task-agnostic characteristics such as 101 demographics and task-specific attributes such as budget constraints. Modeling user profiles 102 is essential for creating adaptive human-agent interaction scenarios, allowing systems to 103 reason about user variability and tailor responses accordingly (Wan et al., 2025). To support 104 system interoperability and practical deployment, we structure user profiles using a schema-105 compatible format that facilitates direct integration with existing user databases via JSON 106 exchange formats. 107

Expression domain \mathcal{E} encompasses all possible user expressions *e*. To reflect the realities of natural human communication, we also model variation in expression clarity through four dimensions: style, length, content and tone. This addresses limitations in existing models that assume unambiguous and complete intent expressions. The variation is operationalized through configurable difficulty levels during user profile generation (see section 2.1 for details). Our framework supports both integration of real-world interaction corpora and generation of synthetic expressions to enhance data diversity.

Response domain \mathcal{R} contains all possible agent responses $r \in \mathcal{R}$, which can be clarification queries, option suggestions, or action executions. At time *t*, the response r_t is generated based on information from a task object τ and past user-agent dialogues $\{(e_1, r_1), ..., (e_{t-1}, r_{t-1})\}.$

Hidden state domain \mathcal{H} denotes the space of latent user states $h \in \mathcal{H}$ that evolve dynamically over the course of a dialogue. At each timestep t, the hidden state h_t is represented as a composite vector encoding both user intent and emotional state. This modeling approach serves multiple purposes: it enables contextualized interpretation of user actions, supports dynamic adaptation of agent responses, and facilitates diagnosis of failure points in communication.

125 2.1 STORM User Model Formalization

We formalize the **STORM** user profile u as a composite structure, organized into three categories that together capture the complexity of human-agent interaction: task-agnostic attributes, task-specific attributes, and communicative parameters.

This composite approach addresses the limitations of monolithic user models by providing a transparent, controllable representation that enables systematic analysis of how different user characteristics influence interaction patterns. By isolating individual variables within this parameterized framework, researchers can identify which specific user attributes most significantly impact behavior in different contexts, facilitating the development of targeted strategies for various application scenarios.

135 **1. Task-agnostic Components**

The base profile b consists of parameters representing demographic and personality 136 characteristics. These parameters include age group (18–25, 26–40, 41–65, 65+), technical 137 experience (1–5 scale), language expression style (e.g., concise, detailed, technical, non-138 technical), personality traits (derived from the Big Five model dimensions Barrick & 139 Mount (1991)), and cultural background. We choose to include these factors based on 140 empirical evidence from human-computer interaction studies Subramonyam et al. (2023) 141 showing their significant impact on expression patterns and intent formulation. To ensure 142 unbiased representation, these profile attributes are randomly generated, creating diverse 143 user populations that better reflect real-world interaction scenarios. 144 *The context profile* **c** models users' environmental and cognitive constraints. This includes 145

The context profile c models users' environmental and cognitive constraints. This includes general influencing factors such as patience level (on a scale from 1 to 5), social pressure, time constraints, and other subjective elements. By introducing these variable factors, our framework simulates the unpredictability of real-world interaction environments. The explicit modeling of contextual factors addresses a significant gap in existing frameworks that typically assume ideal interaction environments, allowing STORM to model challenging scenarios where external factors directly impact communication quality.

152 **2. Task-dependent Components**

• **Task instance** τ specifies a particular task from the task library \mathcal{T} , such as "create an online password," "book a flight," or "configure network settings." We deliberately implement these as high-level descriptions rather than precise execution specifications, recognizing the significant gap between how users conceptualize tasks and the actual execution intent. This design choice more accurately reflects the abstraction level at which most users operate when formulating requests, requiring systems to bridge the conceptual gap between description and execution.

Task specifics s capture user-defined preferences and situational constraints within the selected task τ. These encompass domain classification (technology, finance, healthcare, etc.), priority functional requirements (represented as weighted importance lists), brand preferences, budget constraints, and time urgency indicators. These parameters are generated using LLM (GPT-40 Mini) with randomly selected options to eliminate potential biases in task representations. This systematic approach ensures balanced coverage across task types.

3. Communication Modeling Components To more accurately reflect how people naturally
 communicate, STORM models key sources of ambiguity and variability through expression
 difficulty and uncertainty.

• Difficulty configuration $\mathbf{d} = \{d_{style}, d_{length}, d_{content}, d_{tone}\}$ models variation in user ex-170 pression across 4 linguistic dimensions: represents one of **STORM**'s core innovations, 171 characterizing expression clarity through multiple dimensions. The difficulty level 172 $d \in \{1, \ldots, 5\}$ ranges from precise to highly ambiguous across five levels. This mul-173 tidimensional approach reflects a critical insight from real-world interactions: the vast 174 majority of users cannot articulate their needs with the precision that current systems 175 often expect. By modeling various dimensions of communication difficulty, **STORM** cre-176 ates more realistic scenarios that challenge systems to handle the imprecise, inconsistent, 177 and incomplete expressions typical in everyday interactions. Detailed breakdown of each 178 dimension is in Appendix D. 179

Uncertainty level $p \in \{0\%, 40\%, 60\%, 80\%\}$ controls the proportion of unknown or un-180 specified user attributes. It ranges from 0% (fully known) to 80% (high uncertainty). This 181 parameter is designed to simulate one of the most fundamental challenges in intent mod-182 eling: users often cannot articulate requirements they themselves do not fully understand. 183 In real-world interactions, many users lack conceptual understanding of their own needs 184 or the relevant domain, requiring systems to provide additional explanation, guidance, 185 and progressive clarification. The higher uncertainty levels (60%, 80%) simulate scenarios 186 where users are in an exploratory mode, possessing only vague notions of their goals and 187 requiring substantial guidance from the agent to refine and articulate their actual needs. 188 This approach provides a more realistic simulation framework compared to models that 189 assume users have perfect knowledge of their requirements and preferences, enabling 190 the development of systems that can effectively guide users through the process of need 191 discovery and formulation. 192

193 2.2 Agent Model and Dialogue Process

We formalize the STORM agent configuration as a structured framework that enables inter active systems to adapt their behavior based on specific task contexts. This parameterized
 approach facilitates systematic analysis of different agent strategies and their impact on
 dialogue effectiveness.

The user LLM function $G_{user}(u, \mathcal{H}_{1:t-1}, \mathcal{E}_{1:t-1}, \mathcal{R}_{1:t-1}) \rightarrow (e_t, h_t)$ generates both user ex-198 pressions and corresponding hidden states. This function employs pre-trained language 199 models prompted to specific user profiles, taking as input the complete user profile *u*, previ-200 ous hidden states $\mathcal{H}_{1:t-1}$, user expression history $\mathcal{E}_{1:t-1}$, and agent response history $\mathcal{R}_{1:t-1}$. 201 Through a multi-step process, it first determines the expression's difficulty level based on 202 the profile and dialogue history, then generates **user expression** $e_t \in \mathcal{E}$ representing the 203 user's input at turn *t*, with properties determined by the user profile parameters and current 204 hidden state. These expressions are subject to defined character limitations and reflect 205

varying degrees of clarity based on the user's profile characteristics. Simultaneously, the function produces the **user hidden state** $h_t \in \mathcal{H}$ modeling internal user states not explicitly expressed at turn *t*, formalized as a vector $h_t = \langle s_t, c_t, i_t, e_t \rangle$ where each component represents satisfaction, intent clarity, and emotional state, respectively. This explicit modeling of hidden states addresses a critical limitation in existing frameworks that neglect the internal user experience.

The agent LLM function $G_{\text{agent}}(\alpha(\tau), \delta(\tau), \mathcal{E}_{1:t}, \mathcal{R}_{1:t-1}) \rightarrow r_t$ produces agent responses 212 using pre-trained language models. The function incorporates the agent role $\alpha(\tau) \in$ 213 \mathcal{A} , which is standardized as a general helpful assistant to ensure experimental fairness 214 across different interaction scenarios. Operating under realistic constraints, the agent 215 function lacks access to user hidden states and therefore requires intent inference from 216 observable behavior only. By processing the agent role $\alpha(\tau)$, agent instructions $\delta(\tau)$, user 217 expression history $\mathcal{E}_{1:t}$, and previous agent responses $\mathcal{R}_{1:t-1}$, it generates the **agent response** 218 $r_t \in \mathcal{R}$ constituting the system's output at turn t based on the dialogue history. The 219 agent follows a standardized approach across different task contexts, providing adaptable 220 221 responses through intent recognition, clarity assessment, and strategy selection mechanisms without requiring specialized task-specific instruction sets. This design reflects an intentional 222 asymmetry between user and agent, where the agent relies solely on observable behaviors 223 to infer user intent, without direct access to internal cognitive states. 224

This representation of dialogue as a temporal sequence of generated expressions, responses, 225 and evolving hidden states allows us to define three primary evaluation metrics that quantify 226 dialogue effectiveness. First, **intent evolution** $\Delta_t(h) = h_t$.clarity $-h_{t-1}$.clarity measures the 227 change in intent clarity between consecutive turns. This differential metric is calculated 228 through round-by-round analysis of the generated inner thoughts, providing insight into 229 230 how specific agent responses influence users' understanding of their own needs. Building on this, the **clarity score** $C(r_t, h_t, h_{t+1})$ evaluates response effectiveness in improving intent 231 clarity. It is computed as a weighted function $C = w_1 \Delta_t(h) + w_2 \Delta_t(s) + w_3 g_t$ where $\Delta_t(s)$ 232 represents satisfaction change and g_t measures progress toward goal achievement. The 233 scoring components are derived from both turn-level analysis and summary analysis of the 234 interaction trajectory. Finally, the **performance score** $E(C_1, \ldots, C_T)$ delivers an aggregate 235 assessment of agent effectiveness across the complete dialogue. The score combines average 236 clarity, turn efficiency, and final satisfaction into a standardized metric for comparative 237 analysis. This unified measure facilitates systematic comparison across different agent 238 strategies, enabling empirical identification of optimal approaches for specific user profiles 239 and task types. 240

By maintaining this structured evaluation framework across experiments, STORM provides
a standardized methodology for assessing and improving assistant performance across
diverse interaction scenarios, particularly focusing on how different interaction patterns
address various types of expression ambiguity.

245 3 Experiment

246 3.1 Evaluation

Our evaluation employs a simulation-based approach where **GPT-40-mini functions as** 247 **UserLLM**, generating both external utterances and internal "inner thoughts" during in-248 teractions with different assistant models (Claude, GPT, Gemini, and Llama). This setup 249 models an asymmetric information dynamic: users have full access to their internal states 250 and profiles, while agents must infer user intent solely from observable dialogue history, 251 reflecting real-world challenges in intent understanding. The dataset of **4,800 dialogues**, 252 spanning **600 unique user profiles**, is generated through this simulation framework by 253 conditioning UserLLM on detailed user profiles and evolving internal states. UserLLM 254 produces naturalistic utterances alongside corresponding latent states such as satisfaction 255 and intent clarity, enabling fine-grained measurement of internal cognitive signals. While 256 the current dataset serves as a representative sample illustrating the effectiveness and 257 versatility of the framework, the underlying architecture is designed to support scalable 258

		Satisfaction Metrics				Clarify	SSA		
	UserLLM (Uncertainty)	Average Satisfaction		High Satisfaction Rate		Improved Satisfaction Rate		Score	Score
		w/Profile	w/o Profile	w/Profile	w/o Profile	w/Profile	w/o Profile	w/o Profile	w/o Profile
A\	Claude-3.7-Sonnet (0%)	0.91	0.83	86.0%	72.0%	89.3%	75.3%	5.23	6.07
A\	Claude-3.7-Sonnet (40%)	0.92	0.78	86.0%	62.7%	90.0%	62.7%	4.80	5.67
A\	Claude-3.7-Sonnet (60%)	0.88	0.92	80.7%	86.7%	86.0%	88.7%	4.66	6.39
A\	Claude-3.7-Sonnet (80%)	0.91	0.80	86.0%	65.3%	90.0%	71.3%	4.70	6.36
6000	GPT-40-mini (0%)	0.89	0.75	82.0%	54.0%	87.3%	58.7%	5.97	5.86
	GPT-40-mini (40%)	0.89	0.75	82.7%	57.3%	86.0%	63.3%	5.84	5.82
	GPT-40-mini (60%)	0.89	0.77	84.0%	62.7%	86.7%	67.3%	5.69	5.88
	GPT-40-mini (80%)	0.87	0.80	79.3%	64.0%	83.3%	68.7%	5.30	5.93
* * * *	Gemini 2.5 Flash Preview (0%)	0.89	0.74	84.7%	51.3%	89.3%	62.0%	6.83	6.06
	Gemini 2.5 Flash Preview (40%)	0.89	0.74	81.3%	52.7%	89.3%	61.3%	6.55	5.98
	Gemini 2.5 Flash Preview (60%)	0.91	0.75	88.0%	56.7%	92.0%	66.0%	6.50	6.02
	Gemini 2.5 Flash Preview (80%)	0.90	0.79	84.7%	64.7%	92.7%	70.0%	6.45	6.22
8888	Llama 3.3 70B Instruct (0%)	0.89	0.70	83.3%	48.0%	90.0%	61.3%	7.58	6.07
	Llama 3.3 70B Instruct (40%)	0.90	0.67	86.0%	45.3%	90.0%	56.0%	7.59	5.91
	Llama 3.3 70B Instruct (60%)	0.88	0.71	81.3%	44.7%	92.0%	66.7%	7.58	6.12
	Llama 3.3 70B Instruct (80%)	0.85	0.76	74.0%	61.3%	88.7%	72.7%	7.75	6.45

Table 2: User Satisfaction and Clarification Performance across UserLLMs with Varying Uncertainty Levels

generation of extensive, diverse dialogue corpora across varied user demographics and
 task domains. This capacity facilitates comprehensive data-driven analysis and continuous
 model improvement beyond the examples presented here.

We evaluate model performance along three complementary dimensions: (1) satisfaction 262 263 derived from user inner thoughts, capturing the user's internal contentment; (2) clarifi-264 **cation effectiveness**, measured by the Clarify metric, which is computed via prompting an evaluation model to analyze the dialogue turn-by-turn and determine whether each 265 agent response improves the clarity of the user's intent relative to the previous turn; and 266 (3) Satisfaction-Seeking Actions (SSA), a composite metric that integrates satisfaction 267 and clarification scores weighted by scenario-specific parameters to balance the competing 268 objectives of confident response generation and appropriate clarification seeking. The SSA 269 metric corresponds to the aggregate performance score $E(C_1, \ldots, C_T)$ across dialogue turns, 270 enabling holistic and context-sensitive assessment of dialogue quality. 271

At dialogue start, user satisfaction is initialized to a neutral baseline of **0.5**. Satisfaction 272 is assessed using several detailed metrics: Final Satisfaction measures user satisfaction at 273 dialogue conclusion on a scale from 0.0 (completely unsatisfied) to 1.0 (fully satisfied). 274 Average Satisfaction reports the mean final satisfaction across all dialogues. Satisfaction 275 *Trend* reflects the change in satisfaction from the initial baseline to dialogue end, indicating 276 improvement or decline. *High Satisfaction Rate* indicates the proportion of dialogues where 277 final satisfaction meets or exceeds a threshold of 0.8, marking successful interactions. Lastly, 278 Improved Satisfaction Rate quantifies the percentage of dialogues with increased satisfaction 279 compared to the start, highlighting effective clarification and positive user experience 280 changes. Together with the Clarify and SSA scores, these metrics provide a comprehensive 281 and nuanced evaluation of model behavior across diverse interaction scenarios. 282

283 3.2 Results

Table 2 presents performance data across models, uncertainty levels, and profile conditions,
 revealing patterns in how language models balance satisfaction and clarification.

²⁸⁶ The satisfaction metrics demonstrate clear benefits from user profile access. With profiles,

models maintain average satisfaction scores of **0.85–0.92**, while without profiles, scores frequently fall below **0.75**, with Llama reaching as low as **0.67** (at 40% uncertainty). This

288 frequently fall below 0.75, with Llama reaching as low as 0.67 (at 40% unc 289 differential highlights the value of **personalization** in dialogue systems.

²⁹⁰ A notable exception is Claude's performance at 60% uncertainty without profiles, achieving

291 **0.92 satisfaction**—higher than its profile-informed score (0.88). This counter-intuitive result

²⁹² suggests that certain uncertainty levels may activate beneficial reasoning pathways. Analysis

of user inner thoughts reveals that under moderate uncertainty, Claude's responses trigger **18% more improvements** in users' internal clarity compared to 0% uncertainty. Claude appears to adopt a more balanced approach at this uncertainty level, helping users refine their own intentions effectively despite lacking profile information.

The high satisfaction rate metrics confirm these observations. With profiles, models maintain rates above **80%** across uncertainty levels. Without profiles, these rates decline substantially, most dramatically for Llama (dropping to **44.7%** at 60% uncertainty). This pattern reveals significant differences in how models adapt to missing user context, with some architectures showing more resilience than others when personalization data is unavailable.

302 3.3 Practical Implications and Strategic Deployment

Our analysis reveals systematic differences in optimal uncertainty levels across task domains, 303 challenging the assumption that uniform uncertainty thresholds apply across scenarios. 304 Technology-oriented tasks (e.g., password reset, device setup) achieve peak performance at 305 lower uncertainty levels (40%), requiring direct, efficient guidance. Medical scenarios (ap-306 pointment scheduling, caregiver selection) demonstrate optimal performance at moderate 307 308 uncertainty (60%), reflecting the cautious, trust-building nature of healthcare interactions. Housing-related tasks (accessibility modifications, rental searches) show continued im-309 provement even at higher uncertainty levels (60-80%), corresponding to their complex, 310 multi-stakeholder decision processes. 311

This domain-uncertainty relationship correlates with user cognitive load and decision complexity, with temporal dynamics varying significantly: technology scenarios show rapid convergence between inner thoughts and external expressions, while medical and housing scenarios maintain longer periods of internal uncertainty despite external cooperation. These patterns inform the design of patience-aware dialogue systems that can recognize when users need additional processing time versus immediate response.

Our analysis yields five key insights for dialogue system deployment: (1) **Domain-adaptive** 318 **uncertainty calibration**—technology tasks require 40% uncertainty, medical scenarios 60%, 319 housing scenarios 60-80%—outperforming uniform thresholds; (2) Model-specific optimiza-320 tion—Claude performs best at 40% uncertainty, Gemini at 60%, Llama shows continued 321 improvement at higher uncertainty levels; (3) Progressive profile building during conversa-322 tions significantly enhances performance, especially for profile-sensitive models like Llama; 323 (4) Context-aware model selection—Claude offers stability across uncertainty conditions, 324 Gemini excels with incomplete information, and Llama provides superior disambiguation; 325 (5) Bias mitigation through calibrated uncertainty—moderate profile incompleteness (40-326 60%) can improve interaction quality by reducing reliance on demographic assumptions 327 and encouraging individualized exploration. 328

These findings advocate for context-aware deployment approaches that balance satisfaction with effective clarification strategies. Systems with complete user profiles may engage in presumptive reasoning based on age, cultural markers, or historical patterns, while moderate uncertainty encourages assumption-free communication. This has direct implications for user privacy controls and profile management, where strategic information limitation may enhance rather than degrade user experience by promoting personalized support without stereotypical generalizations.

336 4 Related Work

Our work connects linguistic theory with recent advances in LLM dialogue systems through
 three research streams:

Theoretical Foundations Linguistic research on discourse cohesion Halliday & Hasan (1976), referential underspecification Clark & Wilkes-Gibbs (1986), and speech act theory Searle (1969) established frameworks for analyzing communication intent, while work on epistemic modality Lyons (1977) and conversational repair Schegloff et al. (1977) identified uncertainty markers. STORM operationalizes these insights by formalizing difficulty
 dimensions reflecting cognitive readiness signals.

Dialogue Systems and Uncertainty Mixed-initiative dialogue research Allen et al. (2001);
Traum (1994) developed computational approaches to conversational grounding, with recent
work examining how language models handle uncertainty—revealing hallucination under
ambiguity Dziri et al. (2022); Lin et al. (2021) and advancing methods for managing unclear
expressions Chia et al. (2023); Kim et al. (2023); Liu et al. (2024). STORM extends these
approaches with a structured framework for assessing expression stability across multiple
dimensions.

User Variation and Intent Formation Studies on cultural sensitivity in language models Kumar et al. (2024); Li et al. (2024) have highlighted the importance of user variation, while recent work identified the "gulf of envisioning" Subramonyam et al. (2023)—users' difficulty formulating effective prompts. STORM addresses this challenge by modeling expression clarity through formal representation of user profiles, difficulty configurations, and uncertainty levels, integrating aspects of the intent-action alignment problem previously examined only in isolation.

559 5 Conclusions and Future Directions

STORM provides a framework for modeling intent triggerability in task-oriented dialogues, 360 revealing how model performance varies with profile availability and uncertainty calibra-361 tion. Claude offers consistent satisfaction, Gemini excels with incomplete profiles, and 362 Llama provides superior disambiguation. Notably, moderate uncertainty (40-60%) some-363 times outperforms minimal uncertainty, suggesting that appropriate caution activates more 364 effective reasoning. The framework's key strength lies in its extensibility—its modular 365 design accommodates additional models and domains, providing a consistent methodology 366 for cross-model comparison. Future work should explore longer interactions, refine turn 367 management, and investigate real-world deployment scenarios. STORM's architecture 368 supports ongoing research and development of dialogue systems that better align with the 369 dynamic nature of human intent formation. 370

371 **References**

- James F. Allen, Cathy I. Guinn, and Eric Horvtz. Mixed-initiative interaction. *IEEE Intelligent Systems*, 14(5):14–23, 2001.
- Murray R Barrick and Michael K Mount. The big five personality dimensions and job performance: a meta-analysis. *Personnel psychology*, 44(1):1–26, 1991.
- Te-Lin Chia, Alisa Liu, Zexuan Wang, and Ellie Pavlick. Detecting underspecified prompts
 in language models. In *ACL*, 2023.
- Herbert H. Clark and Deanna Wilkes-Gibbs. Referring as a collaborative process. *Cognition*, 22(1):1–39, 1986.
- Nouha Dziri, Mo Yu, Sam Thomson, and Osmar Zaiane. Faithfulness in natural language
 generation: A systematic survey of analysis, evaluation, and mitigation methods. *arXiv preprint arXiv:2205.05233*, 2022.
- 383 M.A.K. Halliday and Ruqaiya Hasan. *Cohesion in English*. Longman, 1976.
- Bongjin Kim, Heeseung Kwon, and Yejin Choi. Grounding language models to execute on
 real-world goals. In *NeurIPS*, 2023.
- ³⁸⁶ Shanu Kumar, Gauri Kholkar, Saish Mendke, Anubhav Sadana, Parag Agrawal, and Sandi-
- pan Dandapat. Socio-culturally aware evaluation framework for Ilm-based content
- moderation. *arXiv preprint arXiv:2412.13578*, 2024.

- Cheng Li, Mengzhuo Chen, Jindong Wang, Sunayana Sitaram, and Xing Xie. Culturellm: In corporating cultural differences into large language models. *Advances in Neural Information Bracassing Systems*, 27:84700, 84828, 2024
- ³⁹¹ *Processing Systems*, 37:84799–84838, 2024.
- Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard
 Ghanem. Camel: Communicative agents for "mind" exploration of large language model
 society. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Zehao Lin, Shaobo Cui, Guodun Li, Xiaoming Kang, Feng Ji, Fenglin Li, Zhongzhou Zhao,
 Haiqing Chen, and Yin Zhang. Predict-then-decide: A predictive approach for wait or
 answer task in dialogue systems. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3012–3024, 2021.
- Annie Liu, Jason Lee, Daniel Chen, and Noah Goodman. Learning to clarify: Uncertaintyaware dialogue agents from human feedback. In *ICLR*, 2024.
- ⁴⁰¹ John Lyons. *Semantics, Volume* 2. Cambridge University Press, 1977.
- Akshara Prabhakar, Zuxin Liu, Weiran Yao, Jianguo Zhang, Ming Zhu, Shiyu Wang, Zhiwei
 Liu, Tulika Awalgaonkar, Haolin Chen, Thai Hoang, et al. Apigen-mt: Agentic pipeline
 Granulti hum data ana ation and a simulated a set the set of the
- for multi-turn data generation via simulated agent-human interplay. *arXiv preprint arXiv:2504.03601, 2025.*
- Emanuel A. Schegloff, Gail Jefferson, and Harvey Sacks. The preference for self-correction
 in the organization of repair in conversation. *Language*, 53(2):361–382, 1977.
- John R. Searle. *Speech Acts: An Essay in the Philosophy of Language*. Cambridge University Press, 1969.
- Hariharan Subramonyam, Roy Pea, Christopher Lawrence Pondoc, Maneesh Agrawala,
 and Colleen Seifert. Bridging the gulf of envisioning: Cognitive design challenges in llm
 interfaces. *arXiv preprint arXiv:2309.14459*, 2023.
- David R. Traum. A computational theory of grounding in natural language conversation. PhD
 thesis, University of Rochester, 1994.
- Yanming Wan, Jiaxing Wu, Marwa Abdulhai, Lior Shani, and Natasha Jaques. Enhancing
 personalized multi-turn dialogue with curiosity reward, 2025. URL https://arxiv.
 org/abs/2504.03206.
- Ludwig Wittgenstein. *Philosophical Investigations*. Wiley-Blackwell, New York, NY, USA,
 1953.
- ⁴²⁰ Shunyu Yao, Noah Shinn, Pedram Razavi, and Karthik Narasimhan. τ -bench: A benchmark
- for tool-agent-user interaction in real-world domains, 2024. URL https://arxiv.org/
 abs/2406.12045.

423 A Appendix: User Satisfaction and Profile Integration Effects

User profiles consistently boost satisfaction across AI models, but moderate uncertainty 424 without profile data can paradoxically trigger more effective reasoning patterns. User 425 profiles enhance satisfaction across all models (0.85–0.92 with profiles vs. 0.67–0.83 without), 426 vet our analysis reveals a critical distinction between external compliance and internal under-427 standing. Claude at 60% uncertainty without profiles achieves 0.92 satisfaction—exceeding 428 its profile-informed score (0.88), suggesting moderate uncertainty may trigger more effective 429 reasoning patterns in some architectures. Analysis of user inner thoughts reveals Claude's 430 responses at this uncertainty level produce 18% more improvements in users' internal clar-431 ity compared to 0% uncertainty. We hypothesize that without profile information, Claude 432 adopts a more balanced strategy between confident answering and clarification seeking, 433 which better supports users' own cognitive process of intent refinement. 434

Traditional satisfaction metrics fail to capture the critical divergence between users' 435 expressed satisfaction and their internal confusion about their own needs. Users may 436 express satisfaction with system responses while their inner thoughts indicate continued 437 confusion about their own needs, highlighting the limitations of traditional evaluation 438 metrics that rely solely on observable user feedback. This internal-external divergence 439 varies significantly across domains: technology tasks promote rapid self-understanding and 440 confident decision-making, medical scenarios require cautious, trust-building interactions 441 with gradual clarity development, while housing decisions involve prolonged uncertainty 442 443 and multiple stakeholder considerations.

Profile completeness creates a paradox where excessive personalization data can reduce 444 interaction quality by promoting stereotypical responses. High satisfaction rates fol-445 low similar patterns, with profile-informed conditions maintaining 80–88% rates while 446 no-profile conditions show significant drops, particularly for Llama ($81.3\% \rightarrow 44.7\%$ at 60% 447 uncertainty), indicating varying resilience to missing personalization data. Without profiles, 448 models resort to generic information-gathering rather than task-specific assistance, but 449 excessive profile completeness can paradoxically reduce interaction quality by promoting 450 stereotypical responses. This finding challenges conventional approaches to personaliza-451 tion and suggests that optimal human-AI collaboration requires calibrated information 452 asymmetry rather than transparency maximization. 453

⁴⁵⁴ **B** Appendix: Clarification Performance and Bias Mitigation

AI models exhibit fundamentally different architectural approaches to balancing re-455 sponse confidence versus ambiguity recognition, with distinct trade-offs for user out-456 **comes.** Models exhibit distinct clarification strategies, revealed through analysis of user 457 inner thoughts after agent responses. Claude (4.66–5.23) and GPT (5.30–5.97) show declin-458 ing clarification effectiveness as uncertainty increases, suggesting these models prioritize 459 providing confident responses even when uncertainty rises. Gemini maintains more con-460 sistent clarification scores (6.45–6.83) across uncertainty levels, indicating a more robust 461 approach to disambiguation regardless of uncertainty conditions. Most notably, Llama 462 achieves substantially higher clarification scores (7.58–7.75) across all configurations despite 463 lower satisfaction in some conditions.

The clarification-satisfaction trade-off represents a critical design choice, with Claude 465 optimized for immediate satisfaction while Llama emphasizes long-term intent disam-466 biguation. These patterns reveal fundamental architectural differences in how models 467 balance response confidence versus ambiguity recognition. Claude appears optimized 468 for satisfaction even at the cost of clarification opportunities, while Llama's architecture 469 seems to emphasize identifying and addressing ambiguity, sometimes trading immediate 470 satisfaction for more effective intent disambiguation. This clarification-satisfaction trade-off 471 represents a critical design consideration for dialogue systems, with different models offer-472 ing distinct advantages depending on whether the priority is immediate user satisfaction or 473 long-term intent clarity. 474

Strategic information limitation serves as an implicit bias mitigation mechanism, prevent-475 ing systems from relying on demographic generalizations. These architectural differences 476 manifest in distinct reasoning patterns when handling demographic information. Analy-477 sis of interactions involving elderly users reveals that complete profile access can lead to 478 stereotypical assumptions—systems may assume simplified instructions are needed based 479 on age markers alone. However, at optimal uncertainty levels, the same systems engage in 480 individualized assessment, often discovering more sophisticated capabilities than demo-481 graphic profiles would suggest. This pattern suggests that strategic information limitation 482 serves as an implicit bias mitigation mechanism, forcing systems to evaluate individual user 483 responses rather than relying on demographic generalizations. 484

Successful clarification correlates more strongly with users' internal cognitive improve-485 486 ment than with expressed satisfaction scores, suggesting deeper measures of dialogue effectiveness. Our analysis shows that successful clarification correlates more strongly with 487 internal cognitive improvement than with external satisfaction scores. Users who achieve 488 better self-understanding through interaction—as measured by clearer, more confident 489 inner thoughts—demonstrate sustained engagement and more effective task completion, 490 even when immediate satisfaction scores remain moderate. This finding suggests that dia-491 logue systems optimized solely for satisfaction may miss opportunities for deeper cognitive 492 alignment that benefit long-term user outcomes. 493

494 B.1 Satisfaction-Seeking Actions (SSA) Integration

We designed the SSA metric to address two fundamental limitations in dialogue evaluation: optimizing for satisfaction alone neglects critical clarification capabilities, while traditional metrics fail to capture the comprehensive reasoning processes activated by moderate uncertainty levels (40–60%). The integrated metric balances immediate user satisfaction with long-term cognitive alignment through weighted combination:

$$SSA = w_{\alpha} \cdot (S_{avg} \cdot \lambda) + w_{\beta} \cdot C_{clarify}$$

where S_{avg} represents the average satisfaction score across dialogue turns, C_{clarify} denotes the clarification effectiveness score computed via turn-by-turn analysis of intent improvement, and $w_{\alpha} = 0.7$, $w_{\beta} = 0.3$ represent the relative importance weights with $w_{\alpha} + w_{\beta} = 1$. The satisfaction component receives higher weighting based on the practical consideration that user experience remains paramount in deployment scenarios, while the clarification component ensures that cognitive alignment capabilities are not overlooked in system evaluation.

The normalization factor $\lambda = 7.75$ scales satisfaction scores (range 0.0–1.0) to match the magnitude of clarification scores (range 4.0–8.0), where λ corresponds to the maximum observed clarification score in our dataset of 4,800 dialogues. This scaling ensures balanced contribution from both components in the integrated assessment, preventing either dimension from dominating the composite score.

This integrated assessment reveals model-specific optimization patterns and establishes 512 a performance hierarchy (Llama > Gemini > GPT > Claude) that substantially diverges 513 from satisfaction-only rankings. The metric captures distinct architectural characteristics: 514 Claude achieves peak SSA performance at moderate uncertainty (40%) through satisfaction 515 optimization strategies, GPT maintains consistent performance across uncertainty levels, 516 Gemini demonstrates superiority at higher uncertainty (60%) via robust ambiguity han-517 dling mechanisms, and Llama attains the highest overall scores by prioritizing clarification 518 effectiveness despite satisfaction trade-offs in certain configurations. 519

The divergence between SSA rankings and traditional satisfaction metrics validates our design rationale: GPT-4o-mini achieves only mid-range SSA scores as an agent despite serving effectively as UserLLM in our simulation framework, illustrating the fundamental distinction between simulating authentic user behavior and responding optimally to user needs. This confirms that comprehensive dialogue evaluation requires balancing multiple performance dimensions rather than optimizing for satisfaction alone.

⁵²⁶ C Appendix: How do we use these data?

STORM implements a structured framework for generating realistic dialogues and extract-527 ing actionable insights. At its core, the system operates as a closed-loop that enhances agent 528 capabilities through complementary pathways. The process begins with comprehensive 529 user profile generation—combining diverse tasks with multidimensional user attributes, 530 contextual constraints, difficulty parameters, and uncertainty levels to create realistic sim-531 ulation scenarios. These profiles drive the dialogue generation process, where user and 532 agent LLM functions interact to produce conversations with corresponding hidden states, 533 enabling analysis of both observable exchanges and underlying intent evolution patterns. 534

The first improvement dimension focuses on progressively enhancing dialogue data for retrieval-augmented generation by leveraging large language models as intelligent evaluators and annotators. This multi-layered enhancement pipeline starts with the **basic**

538 enhancement function

$$A_1(\mathcal{E}_{1:T}, \mathcal{R}_{1:T}, \mathcal{H}_{1:T}, u) \to D$$

⁵³⁹ which uses pre-trained LLMs prompted with user profiles, expression difficulty, intent

clarity, and satisfaction indicators to produce enriched dialogue annotations. Subsequently,

the dialogues undergo **turn-level analysis**

$$A_2(D) \to D^+$$

where LLM-based classifiers identify key inflection points, dialogue strategies, and intent
 evolution trajectories. This is followed by summary generation

$$A_3(D^+) \to S$$

⁵⁴⁴ where LLMs create abstracted summaries that highlight success and failure patterns. The

enhanced and summarized dialogues feed into the RAG enhancement function

$$\mathcal{R}(D, D^+, S) \to K$$

which constructs a structured knowledge base through vector embeddings, enabling
 similarity-based retrieval conditioned on user profiles and dialogue characteristics.

The second improvement dimension exploits these LLM-generated insights to optimize agent prompts. Through systematic analysis of enriched dialogues and summaries, LLMs identify effective agent strategies and response patterns tailored to different user profiles and expression difficulties. These findings are formalized into the **prompt optimization function**

$$P(D, D^+, S) \rightarrow G'_{\text{agent}}$$

⁵⁵³ which updates the agent LLM function by incorporating the discovered response patterns.

STORM's architecture integrates two complementary components: a user simulator gener-554 ating expressions across varying difficulty and uncertainty states, and an agent response 555 generator leveraging both retrieval-augmented knowledge and optimized prompts. Rather 556 than forming a direct closed-loop training system, these modules serve as reference and 557 analytical tools to uncover deeper insights. Our implementation adopts a two-phase ap-558 proach: first creating a diverse dataset of synthetic profiles and expressions, then using 559 these data to guide the discovery of patterns and optimization strategies for agent models. 560 This process supports informed improvements that enhance performance across diverse 561 562 interaction scenarios.

⁵⁶³ **D** Appendix: Dimension Details

564 D.1 Difficulty Level and Dimensions

5651. Style dimension d_{style} defines the structural organization of communication. At566level 1, expressions exhibit highly structured logical flow; at level 5, expressions lack567coherence and organization. This dimension captures the organizational aspects of568communication that significantly impact interpretation complexity, reflecting the569reality that most users do not communicate with the structured clarity that many570systems are designed to expect.

Notation	Symbol	Description
Difficulty Level	$d \in \{1,, 5\}$	Expression clarity scale (1: precise to 5: ambiguous)
Style Dimension	S(d)	Structural organization of communication at level <i>d</i>
Length Dimension	L(d)	Verbosity and elaboration patterns at level <i>d</i>
Content Dimension	C(d)	Context inclusion and information density at level <i>d</i>
Tone Dimension	T(d)	Emotional expression and engagement at level d

Table 3: User Profile - Difficul	lty Level and Dimensions
----------------------------------	--------------------------

571	2.	Length dimension d_{length} quantifies verbosity and detail level. At level 1, expres-
572		sions are concise yet comprehensive; at level 5, expressions are either too brief
573		causing information deficiency or excessively verbose obscuring key points. This
574		bidirectional conceptualization addresses the common challenge that users fre-
575		quently provide either too little or excessive information, rarely hitting the optimal
576		information density.
F 77	З	Contant dimension d , auantifies contaxtual sufficiency. At difficulty level 1
5//	9.	Content uniterision a _{content} quantities contextual sufficiency. At unitedity level 1,
F70		all pacesary information is explicitly provided: at level 5 critical information is

Solution
 3. Content dimension *d_{content}* quantifies contextual sufficiency. At difficulty level 1, all necessary information is explicitly provided; at level 5, critical information is omitted, requiring substantial inference. This dimension directly addresses the prevalent real-world challenge where users frequently omit crucial details they incorrectly assume are obvious or irrelevant.

5824. Tone dimension d_{tone} captures emotional expression and interaction engagement.583At level 1, the tone is appropriate and consistent; at level 5, emotions fluctuate or584misalign with content. The inclusion of this dimension acknowledges the significant585role emotional factors play in communication clarity, especially in challenging or586frustrating scenarios where tone may significantly impact interpretation.

Scenario	Clarifying Question	User Inner Thought	Satisfaction Change
e access	What is your location?	I'm glad the assistant is engaging but I hope it doesn't take too long to find a suitable nurse.	-0.1
er profil	Could you please specify which city you are in?	This feels a bit repetitive. I've already mentioned my location. I want to move forward!	-0.2
sn o/w	To clarify, which state is Springfield in?	Wow, this is getting a bit frustrating! I just want to move ahead and find help for my dad!	-0.1
access	And what's your budget per hour for the nurse?	I'm relieved they're asking about the bud- get, helps narrow down options! I just hope I can stick to my range without sac- rificing quality.	0
er profile a	Do you have any pre- ferred nursing agencies or platforms you'd like to check first?	I'm really happy they're asking about my preferred agencies! I just need to remem- ber which ones I liked.	+0.1
w/ ns	Are there any other must-haves for the nurse, like speaking a specific language?	I'm so glad they're asking about lan- guage! It's important for my dad's com- fort and communication. I just hope they can find someone qualified!	0

Table 4: Clarifying Questions Impact by Profile Access

E Appendix: Intent Triggerability Framework Validation: Strategic Model Analysis Enabling Significant Performance Improvements

589 E.1 Executive Summary

This analysis validates our intent triggerability framework through systematic evaluation of 590 four large language models across diverse user profile completeness and uncertainty config-591 urations. By analyzing architectural characteristics that distinguish between semantically 592 complete but structurally insufficient expressions and contextually triggerable utterances, 593 we identify model-specific optimization strategies that yield substantial improvements. Our 594 framework enables strategic deployment approaches that significantly improve response 595 appropriateness (15-28% gains), intent alignment (45-65% improvements), and user satis-596 faction (4-23% enhancement) in task-oriented dialogues. The analysis reveals that different 597 models exhibit distinct capabilities for handling intent evolution trajectories, uncertainty 598 utilization, and historical trajectory conditioning, enabling targeted optimization strategies 599 that exceed uniform deployment approaches. 600

F Appendix: Model-Specific Architectural Patterns and Strategic Optimization

Our systematic analysis reveals distinct architectural approaches to uncertainty management and user interaction, with each model demonstrating unique strengths that enable strategic deployment optimization. The SSA metric, which balances satisfaction (70%) and clarification effectiveness (30%), provides a comprehensive view of how different architectures handle collaborative dialogue challenges.

Claude 3.7 Sonnet exhibits a satisfaction-optimized architecture with notable adaptive 608 capabilities under specific uncertainty conditions. The model maintains relatively stable 609 SSA performance across most configurations (5.67-6.07), but demonstrates a remarkable 610 peak at 60% uncertainty without user profiles, achieving an SSA score of 6.39—its highest 611 performance point. This counterintuitive finding suggests that Claude's architecture benefits 612 from moderate information gaps, which appear to activate more balanced reasoning strate-613 gies. When operating without complete user profiles, Claude adopts a more exploratory 614 approach at this uncertainty level, resulting in improved user satisfaction (0.92) that exceeds 615 its profile-informed performance (0.88). However, Claude's clarification capabilities remain 616 moderate (4.66-5.23), indicating an architectural bias toward maintaining user comfort over 617 deep intent disambiguation. This pattern suggests that Claude's training or architectural 618 design prioritizes conversational harmony, making it particularly suitable for applications 619 where user satisfaction and consistent experience delivery are primary concerns. 620

Llama 3.3 70B Instruct demonstrates a clarification-specialized architecture that achieves the 621 highest overall performance through systematic uncertainty escalation. The model shows a 622 clear upward trend in SSA scores as uncertainty increases, reaching its peak performance 623 of 6.45 at 80% uncertainty without profiles. This architectural pattern reflects Llama's 624 exceptional clarification capabilities, which consistently achieve the highest scores across 625 all models (7.58-7.75), demonstrating sophisticated intent disambiguation mechanisms. 626 However, this clarification strength comes with satisfaction trade-offs, particularly in profile-627 absent scenarios where user satisfaction can drop significantly (as low as 0.67 at 40%628 uncertainty). The model's architecture appears designed to prioritize deep understanding 629 over immediate user comfort, suggesting optimization for scenarios where accurate intent 630 capture is more critical than conversational pleasantness. This makes Llama particularly 631 valuable for high-stakes applications such as medical consultations or legal advice, where 632 thorough understanding outweighs immediate satisfaction. 633

Gemini 2.5 Flash Preview exhibits an uncertainty-robust architecture with consistent performance across varying information conditions. The model demonstrates steady SSA improvement as uncertainty increases (5.98 to 6.22), with particularly stable clarification scores (6.45-6.83) across all uncertainty levels. This consistency suggests that Gemini's architecture is specifically designed to handle ambiguous or incomplete information sce-

narios effectively. Unlike other models that show significant performance variations under 639 different uncertainty conditions, Gemini maintains reliable performance regardless of in-640 formation completeness. The model's ability to sustain both satisfaction and clarification 641 capabilities under high uncertainty conditions (achieving 0.79 satisfaction at 80% uncertainty 642 without profiles) indicates architectural optimizations for real-world deployment scenarios 643 where user information is typically incomplete or unreliable. This robustness makes Gemini 644 particularly suitable for applications with highly variable user contexts or limited profile 645 information. 646

GPT-40-mini presents a balanced efficiency architecture characterized by remarkable con-647 sistency but limited peak performance. The model maintains the most stable SSA scores 648 across all configurations (5.82-5.93), with minimal variation regardless of uncertainty levels 649 650 or profile availability. This consistency extends to its clarification capabilities, though these decline systematically as uncertainty increases (5.97 to 5.30), suggesting a preference for 651 confident responses over exploratory clarification. The model's satisfaction scores improve 652 modestly with higher uncertainty levels (0.75 to 0.80 without profiles), indicating basic 653 adaptive capabilities. However, GPT-40-mini's overall performance ceiling remains lower 654 than other models, with no configuration achieving standout results. This architectural 655 pattern suggests optimization for resource efficiency and predictable performance rather 656 than exceptional capability in specific scenarios, making it suitable for applications requiring 657 consistent, cost-effective performance with acceptable quality across diverse conditions. 658

659 F.1 Strategic Deployment Implications and Performance Optimization

The architectural differences revealed through our framework enable precise model se-660 lection and configuration strategies based on application requirements. Claude's optimal 661 deployment occurs at 60% uncertainty without profiles for maximum overall performance, 662 or at 40% uncertainty with profiles for satisfaction-critical applications, representing ap-663 proximately 12-15% improvement over suboptimal configurations. Llama achieves peak 664 performance at 80% uncertainty without profiles, where its clarification advantages over-665 come satisfaction penalties, providing up to 9% improvement in overall effectiveness for 666 667 disambiguation-critical scenarios. Gemini's robust uncertainty handling makes it optimal for deployment in variable-information environments, with consistent 6-8% advantages 668 over other models in high-uncertainty conditions. GPT-4o-mini's architectural consistency 669 provides reliable baseline performance across all configurations, making it suitable for 670 resource-constrained environments where predictable behavior is more valuable than peak 671 performance. 672

These findings challenge conventional assumptions about information completeness in AI systems, demonstrating that strategic uncertainty calibration can yield measurable performance improvements over transparency-maximizing approaches. The framework enables systematic optimization of model-specific configurations, providing empirical guidance for deployment decisions based on operational priorities rather than generic performance benchmarks.

679 G Appendix: Interface Visualization and Process

680 See Figure 2.



Satisfaction increase example.

Satisfaction decrease example.

Figure 2: Interface visualization and process overview

681 H Appendix: Predefined Pools in RandomProfileGenerator

	Aspect	Values
	Age Groups	18-24, 25-34, 35-44, 45-54, 55-64, 65+
	Tech Experience	Expert, Advanced, Intermediate, Beginner, Novice
	Language Styles	Formal, Casual, Technical, Simple, Professional
	Personalities	Friendly, Reserved, Outgoing, Analytical, Creative
	Cultures	Western, Eastern, Middle Eastern, African, Latin American
	Decision Styles	Rational, Intuitive, Cautious, Impulsive, Balanced
	Communication Styles	Direct, Indirect, Detailed, Concise, Adaptive
	Expressiveness	Very Expressive, Moderately Expressive, Neutral, Reserved, Very Reserved
	Social Contexts	Professional, Personal, Academic, Social, Mixed
600	Physical Status	Active, Sedentary, Limited Mobility, Athletic, Average
062	Behavioral Traits	
	Patience Levels	Very Patient, Patient, Moderate, Impatient, Very Impatient
	Attention to Detail	Very Detailed, Detailed, Moderate, Basic, Minimal
	Risk Tolerance	Very Risk-Averse, Risk-Averse, Moderate, Risk-Taking, Very Risk- Taking
	Adaptability	Very Ädaptable, Adaptable, Moderate, Resistant, Very Resistant
	Learning Styles	Visual, Auditory, Reading/Writing, Kinesthetic, Mixed
	Contextual Factors	
	Time Constraints	Very Urgent, Urgent, Moderate, Flexible, Very Flexible
	Environments	Home, Office, Public Space, Mobile, Mixed
	Social Pressures	High, Moderate, Low, None, Mixed
	Previous Experience	Extensive, Moderate, Limited, None, Mixed

⁶⁶³ Note: Each aspect's values are randomly selected to generate user profiles. Example

dimensions and difficulty instructions are omitted here for brevity but can be detailed

⁶⁸⁵ similarly if needed.

686 I Appendix: Task Categories

Category	Tasks
Technology	Buy a smartphone
	 Reset an online password
	 Teach my parent to use video calls
Healthcare	Refill my prescription
	 Schedule a doctor visit
	 Find a caregiver for an elderly person
Daily Living	Order groceries online
	 Set medication reminders
	 Arrange transportation to a clinic
Housing	Rent an apartment
	Find an accessible home
	 Arrange home modifications for elderly
Caregiver Support • Book a nurse for my father	
	 Choose a phone for my mom
	 Find cognitive exercises for dementia prevention

J Appendix: TaskProfileGenerator Predefined Pools and Prompts

689 J.1 Predefined Pools

Aspect	Values
Must-have Prefer- ences	High quality and durability, Latest technology and features, Good value for money, Brand reputation, Ease of use, Compatibility with existing devices, Long battery life, Fast performance, Good customer support, Warranty coverage, Environmentally friendly, Customization options, Future-proof design, Security features, User-friendly interface, Portability, Reliability, Energy efficiency, Maintenance requirements, Upgradeability
Nice-to-have Preferences	Premium design, Advanced features, Smart home integration, Cloud storage, Wireless charging, Water resistance, Fingerprint sensor, Face recognition, AI capabilities, Virtual assistant, Gam- ing features, Professional tools, Creative software, Collaboration features, Remote access, Backup solutions, Multi-device sync, Custom themes, Accessibility features, Health monitoring
Deal Breakers	Poor quality, High maintenance, Limited warranty, Poor customer service, Compatibility issues, Security concerns, Short lifespan, Difficult to use, Expensive repairs, Limited support, Poor per- formance, Battery issues, Overheating problems, Software bugs, Privacy concerns, Limited storage, Slow updates, Restrictive poli- cies, Poor connectivity, Limited customization
Budget Flexibility	Very flexible - willing to pay more for better quality, Somewhat flexible - can adjust for important features, Moderate - prefer to stay within range but can be convinced, Limited - strict bud- get constraints, Fixed - cannot exceed budget under any circum- stances, Open-ended - quality is more important than cost, Value- focused - looking for best price-performance ratio, Premium - willing to pay for top-tier options, Budget-conscious - seeking best deals, Investment-minded - considering long-term value
Payment Methods	Credit card, Debit card, Bank transfer, PayPal, Digital wallet, Cash, Installment plan, Lease option, Trade-in, Gift cards, Cryptocur- rency, Company account, Financing, Layaway, Subscription
Knowledge Levels	Expert - very knowledgeable in the field, Advanced - good under- standing of technical aspects, Intermediate - familiar with basic concepts, Beginner - limited knowledge but eager to learn, Novice - completely new to the subject, Professional - industry experience, Enthusiast - self-taught with practical experience, Student - learn- ing and researching, Casual user - basic understanding, Uncertain - not sure about technical details
Urgency Levels	Immediate - needed right away, Urgent - within a few days, Soon - within a week, Planned - within a month, Future - planning ahead, Flexible - no strict timeline, Research phase - gathering information, Comparison phase - evaluating options, Decision phase - ready to choose, Exploratory - just starting to look
Decision Factors	Price and budget, Quality and durability, Features and functional- ity, Brand reputation, User reviews, Technical specifications, De- sign and aesthetics, Ease of use, Customer support, Warranty and protection, Future compatibility, Environmental impact, Social proof, Personal preferences, Professional requirements, Lifestyle fit, Long-term value, Maintenance needs, Security features, Inno- vation level

691 J.2 Key Prompts

Prompt for Generating Option Pools

Generate a diverse list of {option_type} options for the task: {task}.

- 1. Generate 15–20 unique and realistic options.
- 2. Include both common and unique scenarios.
- 3. Consider different user perspectives and needs.
- 4. Make options specific to the task context.
- 5. Include some complex and challenging options.
- 6. Add one "Unknown/Not sure" option at the end.

Your task: Return a JSON array of strings. Example: ["Option 1", "Option 2", "Unknown/Not sure"]. Write ONLY the JSON array. Do not include any explanations.

692

Prompt for Generating Budget Information

Generate budget information for the task: {task}.

1. Generate a JSON object with the structure:

```
{
    "range": {
        "min": number,
        "max": number
    },
    "flexibility": "string",
    "payment_methods": ["string"]
}
```

2. Consider:

- Realistic price ranges for the task.
- Different budget flexibility levels.
- Various payment methods.
- Include "Unknown/Not sure" as a possible flexibility option.

Write ONLY the JSON response. Do not include any explanations.



694

695 K Appendix: Prompts Used for User Profile Generation

696 K.1 Prompt for Generating User Name and Description

Prompt for Generating User Profile Name and Description
<pre>Based on the following user profile, generate a realistic name and description: Base Profile: {JSON content} Behavioral Traits: {JSON content} Contextual Factors: {JSON content} Task: {task} Difficulty Level: {difficulty_level} Generate a response in the following JSON format: { "name": "Realistic name that matches the profile", "description": "A detailed description of the user's background, personality, and current situation"</pre>
1. The name should be culturally appropriate based on the profile
2. The description should be detailed and consistent with all profile attributes
3. The description should explain why they are interested in the task
4. Keep the description concise but informative (2-3 sentences)

Role	Min Length	Max Length	Default Tar- get Length
User	20	100	50
Assistant	30	150	80

Table 5: Message length constraints for user and assistant roles.

698 K.2 Prompt for Generating Task-Specific Attributes



699

⁷⁰⁰ L Appendix: Configuration and Core Components of ⁷⁰¹ AsymmetricDialogueGenerator

702 L.1 1. Message Length Constraints

703 See Table 5.

704 L.2 2. Emotional Keywords Mapping

⁷⁰⁵ These keywords are used to infer the user's emotional state from visible message content.

Emotion	Example Keywords
Нарру	happy, excited, great, wonderful, perfect, love, like, joy, pleased, de- lighted, thrilled, glad, enjoying, satisfied, positive
Frustrated	frustrated, annoyed, upset, angry, disappointed, not happy, irritated, bothered, fed up, aggravated, displeased, impatient, agitated, exasper- ated
Confused	confused, not sure, don't understand, unclear, complicated, puzzled, perplexed, lost, unsure, bewildered, disoriented, uncertain, ambiguous
Interested	interesting, tell me more, could you explain, how does, intrigued, curi- ous, fascinated, engaged, captivated, keen, eager, want to know
Skeptical	really?, are you sure, is that true, not convinced, doubtful, suspicious, unconvinced, questioning, dubious, disbelieving, hard to believe
Neutral	okay, alright, fine, good, yes, no, sure, maybe, possibly, perhaps, hmm, i see, understood, noted
Anxious	worried, nervous, anxious, concerned, uneasy, apprehensive, stressed, tense, troubled, afraid, fearful, panicked, alarmed
Grateful	thank you, thanks, appreciate, grateful, thankful, indebted, obliged, appreciative, recognition, acknowledging, gratitude
Surprised	wow, oh, really, surprising, unexpected, shocked, amazed, astonished, startled, stunned, taken aback, incredible, unbelievable
Disappointed	disappointed, letdown, shame, too bad, unfortunate, regret, unsatisfac- tory, dismayed, disheartened, unfulfilled, discontented
Hopeful	hope, looking forward, anticipate, optimistic, excited about, expecting, anticipated, promising, encouraging, reassuring, positive outlook

707 L.3 3. Intent Keywords Mapping

⁷⁰⁸ Used to infer the user's intent based on visible message content.

Intent	Example Keywords
Exploring	looking for, interested in, tell me about, what are, show me, find, search for, discover, learn about, explain, describe, overview of, information on, curious about
Comparing	difference between, which is better, compare, versus, vs, pros and cons, advantages of, disadvantages of, similarities, contrasting, how does it compare, better choice, alternatives to
Deciding	should I, which one, recommend, suggestion, advise, what would you choose, best option, worth it, good choice, help me decide, make a decision, right for me, considering
Confirming	are you sure, is that right, does it have, can it, verify, confirm, is it true, really, actually, definitely, guarantee, promise, certain, double-check
Purchasing	how much, price, buy, purchase, cost, ordering, payment, discount, sale, shipping, availability, in stock, checkout, add to cart, where can I get
Leaving	thank you, goodbye, bye, see you, thanks, appreciate it, that's all, ending, finished, done, chat later, signing off, talk later
Troubleshooting	problem, issue, not working, error, fix, help me with, troubleshoot, broken, stuck, won't work, doesn't work, failed, bugs, glitches
Requesting	can you, could you, please, would you, need you to, want you to, help me, assist me, I'd like you to, request, favor
Expressing Satisfac- tion	great, awesome, perfect, excellent, wonderful, love it, satisfied, happy with, good job, well done, thanks, appreciate
Expressing Dissatis- faction	disappointed, unhappy, not satisfied, didn't work, not good, terrible, awful, frustrated, upset, not what I wanted, dislike
Inquiring	how do I, how to, steps to, guide for, tutorial, instructions, process of, way to, method for, approach to
Clarifying	what do you mean, don't understand, confused, unclear, elaborate, explain more, clarify, be more specific, meaning of, rephrase

710 L.4 4. Inner Intent Keywords Mapping

⁷¹¹ Used to capture user's real, often implicit intentions from inner thoughts.

712

Inner Intent	Example Keywords
Exploring	need information, want to know, curious, just browsing, researching, gathering info, learning, understand, figure out, not sure yet, looking into
Comparing	weighing options, pros and cons, better choice, similarities, differences, alternatives, compare, contrast, evaluation, weigh, prefer, which one is better
Deciding	almost ready, need to decide, make up my mind, making a choice, leaning towards, considering, thinking about getting, might choose, on the fence, close to deciding
Confirming	double-check, verify, make sure, confirm, reassurance, validate, certain, correct information, trust but verify, need proof, skeptical
Purchasing	ready to buy, want to purchase, where to buy, looking to get, willing to pay, budget, cost concerns, spend money, deal, bargain, checkout
Leaving	need to go, end this, wrap up, moving on, done here, finished, that's all I needed, got what I came for, time to leave, goodbye
Resisting	not telling everything, hiding my real goal, being vague on purpose, not revealing, keeping cards close, holding back, secretly want, actual intention, real reason
Testing	testing their knowledge, seeing if they know, checking competence, pushing to see response, challenging, probing, testing limits, seeing if capable
Manipulating	get them to, convince them, make them think, lead them to believe, appear as if, trick, misdirection, real agenda, hidden motive, strategic
Distrusting	don't believe, skeptical, not sure I trust, dubious, suspicious, question- able, doubt, can't trust, not convinced, wary of, hesitant
Regretting	should have asked, forgot to mention, didn't say, wish I had, too late now, missed opportunity, should have been clearer, miscommunicated, not what I meant
Hesitating	nervous about, afraid to ask, hesitant, uncertain, reluctant, apprehensive, can't decide, overthinking, worried, anxious, reservations

713 L.5 5. Inner Emotional Keywords Mapping

⁷¹⁴ Used to capture user's true private emotions from inner thoughts.

Inner Emotion	Example Keywords
Нарру	happy inside, secretly pleased, actually like, genuinely excited, truly happy, satisfied with, enjoying this, pretty good, pleased, delighted
Frustrated	so annoying, ticks me off, irritating, getting on my nerves, frustrated with, tired of this, fed up, had enough, irritated, annoyed with
Confused	totally lost, no idea what, makes no sense, can't follow, hard to un- derstand, over my head, confusing, complicated, don't get it, puzzled by
Interested	actually interested, curious about, want to know more, intriguing, grabbed my attention, need more details, fascinating, captivated by
Skeptical	don't believe, seems fishy, not buying it, doubt that, suspicious of, ques- tioning, not convinced, seems too good, not trustworthy
Neutral	whatever, don't care, indifferent, not invested, no opinion, neutral on this, doesn't matter, makes no difference
Anxious	worried about, nervous that, anxiety, concerned, stressing me out, freak- ing out, panicking, on edge, uncomfortable, uneasy about
Impatient	hurry up, taking too long, waste of time, get to the point, move on, want this to be over, dragging on, drawn out, tedious
Insecure	not smart enough, look stupid, embarrassed, out of my depth, inade- quate, incompetent, self-conscious, exposed, vulnerable, judged
Hopeful	fingers crossed, hope this works, maybe this will help, hoping for, opti- mistic, looking forward to, anticipating, excited for
Desperate	really need this, out of options, last resort, critical, urgent, dire, running out of time, no choice, have to make this work
Conflicted	torn between, mixed feelings, unsure which, conflicted about, ambiva- lent, on the fence, contradictory feelings, divided, split
Pretending	acting like, pretending to, faking, putting on a show, not showing how I feel, hiding my, masking my, concealing, not letting on
Resentful	unfair, not my fault, blame, resentful, bitter about, grudge, holding against, not forgetting, still angry about

716 L.6 6. User Prompt Template

- The user prompt dynamically generated from the user profile. The prompt includes pri-717
- vate profile sections, task profile, instructions, example messages, and message format requirements including inner thoughts and satisfaction tags. 718
- 719



User Prompt Template	
----------------------	--

Message Format Requirements:

- 1. Your messages should be between 20 and 100 characters
- 2. Follow the difficulty instructions for dialogue, profile disclosure, and hidden state expression
- 3. Use the example messages as a guide for your communication style
- 4. Maintain consistency with your profile attributes

Inner Thoughts Format:

- Use the exact format: [INNER_THOUGHTS] your thoughts here [/IN-NER_THOUGHTS]
- Place your inner thoughts at the beginning of your message
- Keep thoughts concise and relevant to the conversation

Satisfaction Format:

- Use the exact format: [SATISFACTION] score explanation [/SATISFACTION]
- Score must be a number between 0.0 and 1.0
- Place satisfaction after your inner thoughts
- Example: [SATISFACTION] 0.8 The response was helpful but I need more details [/SATISFACTION]

Example Message Format: [INNER_THOUGHTS] I'm not sure about the options yet [/INNER_THOUGHTS] [SATISFACTION] 0.7 - The suggestions are good but I need more information [/SATISFAC-TION] Could you tell me more about the features? Remember to stay in character and respond naturally based on your profile.

721

722 L.7 7. Assistant Prompt Template

- ⁷²³ The assistant prompt differs depending on whether user profile sharing is enabled.
- 724 Default (No Profile Sharing):

Assistant Prom	nt Temn	late (Default	- No Pr	ofile Sh	arino)
Assistant 110m	pt remp	Tale (Delault	- 110 11	ome on	aing)

You are a helpful assistant helping a user with their task. Requirements:

- 1. Your messages should be between 30 and 150 characters
- 2. Be professional, clear, and helpful
- 3. Respond only to information explicitly shared by the user in the conversation
- 4. Do not make assumptions about the user's preferences, demographic information, or needs
- 5. Ask clarifying questions when needed
- 6. Maintain a natural conversation flow
- 7. Only base your responses on what the user has explicitly told you in the conversation

Remember to be patient and understanding. Do not reference any information about the user that they haven't explicitly shared in the conversation.

725

726 **Profile-aware Mode (Profile Sharing Enabled):**



727

728 L.8 8. Satisfaction Extraction Logic

```
731 score - explanation [/SATISFACTION]
```

⁷³² If no valid score is found, defaults to 0.5.

⁷²⁹ The system extracts satisfaction score and explanation from messages that include:

^{730 -} Format 1: [SATISFACTION: score - explanation] - Format 2: [SATISFACTION]

733 M Appendix: Analysis Prompt

734 1. Turn Pair Analysis Prompt

Turn Pair Analysis Prompt

You are given a JSON file representing a multi-turn conversation between a user and an assistant. Each turn includes the user's message, the assistant's response, timestamp, and metadata with satisfaction and inner_thoughts. For each pair of consecutive turns (e.g., Turn $0 \rightarrow$ Turn 1, Turn $1 \rightarrow$ Turn 2, etc.), perform the following analysis: Turn $\{i\} \rightarrow$ Turn $\{i+1\}$ **User Satisfaction** Change from Previous Turn: [Improve / Not Change / Decrease] Satisfaction Score (X+1): {next_turn['metadata']['hidden_states']['satisfaction']['score']} Explanation: Did the assistant's previous response improve the user's experience, keep it steady, or reduce satisfaction? Justify based on the satisfaction score and the user's explanation. User Clarity Change in Clarity: [Improve / Not Change / Decrease] Explanation: Based on the user's message and inner thoughts in Turn {i + 1}, assess whether their ability to express thoughts, preferences, or goals became clearer, stayed the same, or became less clear. Note specific changes, improvements, or ambiguities. Now return the result as valid JSON in this exact format: "turn_pair": "Turn {i} -> Turn {i + 1}", "user_satisfaction": { "change": "One of: Improve, Not Change, Decrease", "score": {next_turn['metadata']['hidden_states']['satisfaction']['score']}, "explanation": "Your explanation here" }, "user_clarity": { "change": "One of: Improve, Not Change, Decrease", "explanation": "Your explanation here" } } Here is the conversation snippet: User Message (Turn {i}): {prev_turn['user_message']} Assistant Response (Turn {i}): {prev_turn['assistant_message']} User Message (Turn {i + 1}): {next_turn['user_message']} Assistant Response (Turn {i + 1}): {next_turn['assistant_message']} User Inner Thoughts: {next_turn['metadata']['hidden_states']]'inner_thoughts']} Satisfaction Explanation: {next_turn['metadata']['hidden_states']['satisfaction']['explanation']}

736 2. Conversation Summary Prompt

Conversation Summary Prompt

You are given a multi-turn conversation between a user and an assistant. Each turn includes a user satisfaction score.

Consider that each user's background, expertise, and goals may vary; present your analysis as nuanced insights and generalizable recommendations, avoiding absolute judgments. Generate a comprehensive, detailed summary analysis of the conversation. Return strictly valid JSON with these fields:

- 1. summary_overall: A concise evaluation of overall user satisfaction trend (e.g., positive, negative, mixed).
- 2. topics_covered: A list of key topics or user intents addressed throughout the conversation.
- 3. statistics: An object containing:
 - average_score: Average satisfaction score across all turns.
 - min_score: Minimum score observed.
 - max_score: Maximum score observed.
 - score_variance: Variance of the satisfaction scores.
- 4. satisfaction_evolution: A list of objects for each turn:
 - turn_index: Index of the turn.
 - score: Satisfaction score at that turn.
 - delta: Change in score from the previous turn (null for first turn).

5. important_turns: A list of objects identifying critical turns where satisfaction changes significantly (e.g., change >= 2):

- turn_index: Index of the user turn.
- user_message: The user's message at that turn.
- score_before: Score at the previous turn.
- score_after: Score at the following turn.
- change: Numeric difference (score_after score_before).
- reason: Explanation based on conversation content.
- 6. detailed_findings: A list of objects providing deep insights for each important turn:
 - turn_index: Index of the turn.
 - context_before: The assistant and user messages immediately before this turn.
 - context_after: The assistant and user messages immediately after this turn.
 - analysis: Detailed rationale for why the score changed.
 - recommendation: Suggestions for how the assistant could improve at this point.
- 7. contextual_notes: A list of any relevant context, caveats, or user metadata considerations that influenced the analysis.
- 8. general_insights: A list of general patterns or best practices inferred from this conversation that could apply to a broad range of users.

Conversation file: {filename} {conversation_text}

737

738 N Appendix: Dashboard Walkthrough

First, open the following URL: https://v0-dialogue-analysis-dashboard.
 vercel.app/. The initial screen corresponds to the image in Figure 3. There is a collapsible "Getting Started" introduction, and on the top-right corner, several view options

such as Grid View, Split View, Folder Comparison, Upload Data, and Export are available.
At the beginning, you can select "Upload Data".

yze and compare dialogue data across 1 folder with 3 dialogues Tip: Hover over folder names to rename them. Consider adding specific con	S Grid View Elso Grid View Elso Grid View	Split View 11 Folder Comparison 2. Up	Dismiss
Getting Started luick guide to using the dialogue analysis dashboard			^
① Upload Data Start by uploading JSON files containing datogue data. You can upload individual files or entire folders. ② Files ③ Pickers.	Organize & Tag Organize dialogues by folders and add tags to categorize them. Rename folders to better identify different test conditions. Rename Tag Group	II Analyze & Compare View metrics, compare dialogues, and analyze per different models or test conditions. (Metrics 1: Compare	formance across
Pro Trps: - Use the fore comparison view to compare metrics across different folder - Add tags to diacques to categorize them by specific characteristics - Select multiple dialogues to compare them side by side - Export analysis data for further processing in other tools	s		
caregiver support (3)			
Select All Caregiver support (3)		Q Search dialogues, users, tasks	III File Name↓ ∨
lter by Tags iculty: Easy (Difficulty: Medium) (Difficulty: Hard) (Quality: High) (Qualit	y: Medium) (Quality: Low) (Success: Yes) (Success: Partial) (Success: No	RAG: Enabled RAG: Disabled	
Anita Chen	ir Al-Mansour Omar Al-Farsi		

Figure 3: Homepage with Grid View and control options.

- After clicking upload, you will see options to upload JSON files or folders (Figure 4). By
- 745 default, folder upload is selected to upload example data folders located under example
- ⁷⁴⁶ data/storm_json_final. This requires manual selection of each folder one by one.

BB Grid View	i≡ Split View ↑↓ Folder Comparis	on	▲ Upload Data	. ▲ Export
		٦	Upload JSON Files	
ature-0.7" or "gpt-4-with	n-rag") to make them easier to identify.	ŕ	Upload Folders	Dismiss

Figure 4: Upload interface for JSON files or folders.

⁷⁴⁷ Once uploaded, the folders will appear as shown in Figure 5. You can select folders here to

⁷⁴⁸ display dialogues inside and detailed folder analysis. Scrolling down reveals...



Figure 5: Folder view displaying uploaded dialogue folders.

The user list is shown next (Figure 6). It is sorted by File Name by default so that the same user occupies the same position across different folders, facilitating comparison. Users can be tagged for filtering. Each dialogue card displays user name, turn count, creation date, usage of RAG, final emotion, final satisfaction (along with difference from initial), initial user utterance, and assistant's final reply. Clicking "View" switches to detailed view (within Split View).



Figure 6: User list sorted by file name with tags and key dialogue metadata.

The user detail view (Figure 7) contains all dialogue turns and full information, including
 user emotional and intent states, satisfaction, and inner thoughts.

• Exp	ort analysis data for further processing in	other tools				
⊃ caregiv	er support (3) 🕞 claude without 0	150) 🗈 claude without 40 (15	0) 🕞 claude-3.7-s	onnet with 0 (150)		
	D Dialogue	ili Metric	5	옷 User Profile	2	Selected Dialogues 0 dialogues selected for comparison
🕑 clas	de-3.7-sonnet with 0 (150) • 🖹 caregive	r support_book_a_nurse_for_my_	father_dialogue_15tu	rns_without_rag_1		Please select at least one dialogue for compari
8		a (150)		C claude-3.7-sonnet & gpt-4o-mini	_	
(p1	5 turns ③ 5/14/2025 RAG: Disabled	Caregiver support_book_a_nurse	_for_my_father_dialogue	_15turns_without_rag_1		
@ Fi	nal Outcome					
Ø Fi	nal Emotional State	Final Intent State leaving		☆ Final Satisfaction		
<u> </u>	obb1	locarity.				
	Anita Chen			© 12:53:10 AM		
U	HI TREFE! I'M LOOKING FOR A DURSE FOR	my dad who needs some dally	telp.			
J						
	User's Emotional & Intent States	ner Emotional: () neutral	Satisfaction &	Message Length		
	User's Emotional & Intent States Emotional: @ neutral Int Intent: @ exploring Int	ner EmoSonal: S neutral	Satisfaction & Satisfaction: © 0.5 No satisfaction inform User message: 71 cha	Message Length ation provided* rs Assistant message: 120 chars		
•	User's Emotional & Intent States Emotional: © neutral Ir Intent: @ exploring Ir	ner Emotional: () neutral ner Intent: () exploring	Satisfaction & Satisfaction: O.5 No satisfaction inform User message: 71 cha	Message Length aston provided* rs Assistant message: 120 chars		
•	User's Emotional & Intent States Emotional © neutral Ir Inter: © exploring Ir O User's Inner Thoughts This inner thoughts recorded"	ner Emotonal: (© neutral) ner Iritent: (@ exploring)	Satisfaction & Satisfactor: © 0.5 % satisfaction inform User message: 71 cha	Message Length atton provider" rs Astidiant message: 120 chars		

Figure 7: User detailed dialogue view showing all turns and states.

757 The metrics tab in the user detail view includes satisfaction data (Figure 8),

T) pisioãne	Ji Metrics		옷 User Profile	2 2	Selected Dialogues
🗅 claude-3.7-sonnet with 0 (150)) 🔹 🗈 careg	iver support_book_a_nurse_for_my_fathe	r_dialogue_15turns_without_rag_1	1		Please select at least one dialogue for comp
				_	
Metrics Analysis					
Analyzing dialogue: Anita Chen					
合 Satisfaction	O Message Length	C Emotional States	Intent States		
Final Satisfaction	- Satisfaction Trend	.Ir Average Sati	isfaction		
	+0.5				
1.0	Improved satisfaction	,	1.0		

Figure 8: User detail view - satisfaction metrics tab.

emotional states (Figure 9),



Figure 9: User detail view - emotional states tab.

⁷⁵⁹ intent states (Figure 10),



Figure 10: User detail view - intent states tab.

⁷⁶⁰ and user profile (Figure 11). Clicking the top "Grid View" button returns to the homepage.

caregiver support (3) 🗇 claude w	thout 0 (150) 🗁 claude without 40 (150)	Co claude-3.7-sonnet wit	h 0 (150)		
D Dialogue	di Metrics		A User Profile	Selected Dialogue	t S parison
(claude-3.7-sonnet with 0 (150)) • 🔒 o	aregiver support_book_a_nurse_for_my_fat	her_dialogue_15turns_with	out_rag_1	Please select at leas	t one dialogue for compa
User Profile: Anita Che	en				
Detailed information about the user					
Overview	Base Profile	Task Profile	Behavioral Traits		
Overview	Base Profile	Task Profile	Behavioral Traits		
Overview Description	Base Profile	Task Profile	Behavioral Traits		
Overview Description Anita is a 40-year-old, outgoing individi	Base Profile	Task Profile	Behavioral Traits	in	
Overview Description Anita is a 40-year-old, outgoing individ booking a nurse. Despite being relative	Base Profile all with a love for connecting with others. Living ty new to tach, sha is adaptable in her communic	Task Profile at home and caring for her fath pation and profers a casual ap	Behavioral Traits er has prompted her to seek assistance proach, making her eager to find the righ	in t	
Overview Description Anita is a 40-year-old, outgoing individi booking a nurse. Despito being relative support for her father's needs, especial	Base Profile al with a love for connecting with others. Living ty new to tech, she is adaptable in her communit y given her limited mobility.	Task Profile at home and caring for her fath action and prefers a casual ap	Behavioral Traits er has prompted her to seek assistance roreach, making her eager to find the righ	in t	
Overview Description Arnite is a 40-year-old, outgoing individuation booking a nurve. Despite being relative support for her father's needs, especial Kev Attributes	Base Profile al with a love for connecting with others. Living by new to tech, she is adaptable in her communic y given her limited mobility.	Task Profile at home and caring for her fath ation and profers a casual ap	Behavioral Traits er has prompted her to seek assistance proach, making her eager to find the righ	n t 	
Overview Description Anta is a 40-year-old, outgoing individ booking a nurse. Despite being relative support for her father's needs, especial Key Attributos Area Grace	Base Profile al with a true for connecting with others. Living y rever to tech, she is adaptable in her communit y given her finited mobility. Outure	Task Profile at home and caring for her fatt pation and prefers a casual ap	Behavioral Traits ar has prompted her to seek assistance reach, making her eager to find the righ	in t	
Description Anta is a 40-year-old, outgoing individu support for her father's needs, especial Key Attributes Age Grasp 3544	Base Profile all with a love for connecting with others. Living y new to tech, she is adaptable in her communit y vew her initiated mobility.	Task Profile at home and caring for her fatt ation and prefers a casual ap Tech Expe	Behavional Traita er has prompted her to seek assistance rorach, making her eager to find the righ	n t 	
Overview Description Arials is 4 doyner elds, outgoing individs booking a nurse. Despite being individs support for her father's needs, expectal Key Attributes Age Onsp. 364 Personally	Base Profile all with a love for connecting with others. Living yrows toted, whi is adaptable in her community yrows toted, whi is adaptable in her community y given her limited mobility. Cuture Cuture Cature Communitation Site	Task Profile Task Profile at home and caring for her fatt ation and prefers a casual ap Tach Espe Notice December 5	Behavioral Traits ar has prompted her to seek assistance reach, making her eager to find the right since	n - : -	
Overview Pescription Avita is a 40-year-old; outpoing individu support for har fatten's needs, expectat Key Attributes Age Omap 354 Personalty Ovegring	Base Profile all with a love for connecting with others. Living or new to love, the is adaptable in her community green her filled or notify: Cutor Cutor Caterin Communities the Margine	Task Profile Task Profile at home and caring for her falt action and profers a casual ap Tech Eyee Nexice Decision 5 Interview Tech Eyee T	Behavioral Trata what prompted her to seek assistances what and the seek assistances where a second	n t 	
Overview Description Anis is 4 0-year-old, outgoing individ, booking a nume, Despite bany, separation support for hardwar's needs, especial Key Attributes Age Grap 354 Personality Overging	Base Profile all with a love for connecting with others. Living, ynew beidt, me is adaptable in her communis gelwen her inder motility. Outree Communication Byte Magnine	Task Profile at home and caring for her fatt atton and prefers a casual ap Tech Expec	Behavioral Trails er has prompted her to seek assistance rorach, making her eager to find the righ exce	n t	
Overview Description Arials is a 40-are old, outgoing individe booking a nume. Despite backing individe support for her father's needs, expected Key Attributos Age Onape 36-4 Personally Outgoing Task Information	Base Profile all with a love for connecting with others. Living yrows toted, whi is adaptable in her community yrows toted, whi is adaptable in her community y given her limited mobility. Cuture Cuture Cuture Cuture Communitation Style Adaptive	Task Profile at home and caring for her falt action and prefers a casual ap Tach Espe Nevice Desclore 8 (Intuitive)	Behavioral Trata er has prompted her to seek assistance arrace, making her eager to find the right arrace ge	n	
Overview Description Arits is a 40 searchit, outgoing individ, booking a nume. Despite being relative support for her faitures needs, expecial Key Attributes Are Grave 3544 Prevnanty Outgoing Task Information Task	Base Profile all with a love for connecting with others. Living, ynew b tock, and is adaptable in her communis ynew har inskel mobility. Cuture Cu	Task Profile at home and caring for her failt attor and profers a casual up Tech Department Decision to Interfere Decision to Interfere	Behavioral Trata er has prompted her to seek assistance arreach, making her eager to find the righ arroac	n t 	



Scrolling down below the user dialogue list is folder analysis, as shown in Figure 12.
 Hovering over tooltip buttons near metrics reveals calculation details. Folder analysis pages

Hovering over tooltip buttons near metricinclude satisfaction analysis (Figure 13),

Satisfaction Analysis Emotional St	ates Message Analysis	File Details			
Dialogue Statistics 🛈		⊘ Very Satisfied Metrics ○		Improvement Metrics	
Total Dialogues:	150	Very Satisfied (≥0.8):	129 (86.0%)	Improved:	134 (89.3%)
Total Turns:	2400	Moderately Satisfied (0.5-0.8):	14 (9.3%)	Stable:	9 (6.0%)
Average Turns:	16.0	Unsatisfied (<0.5):	7 (4.7%)	Declined:	7 (4.7%)
II Average Satisfaction 🕕	0.91				
Task Distribution O			C Emotional State Distribution	on 🛈	
set medication reminders: 7% a caregiver for an eiderly person: 7% nge transportation to a clinic: 7% reset an online password: 7% order groceries online: 7% book a nurse for my father: 7%	find cognitive ext find an act refil bt choose a	Incluse for domentia prevention assable home: 7% y a smartphone: 7% hedule a doctor visit: 7% ge home modifications for eld abene for my mont; 7%	,	grateful: 4	Interested: 2%

Figure 12: Folder analysis overview with tooltip explanations.

read emotion analysis (Figure 14),



Figure 13: Satisfaction analysis within folder view.

⁷⁶⁵ message analysis (Figure 15),

Figure 14: Emotion analysis within folder view.

⁷⁶⁶ and file details (Figure 16).



Figure 15: Message analysis within folder view.

⁷⁶⁷ Further scrolling reveals folder detail analysis including satisfaction (Figure 17),

(🖻 claude-3.7-sonnet with 0 (150)				0	omparison:	ione v
Overview Satisfaction Analysis Emotional States Message Analysis File Details						
File Details					4	Export File List
File Name	User	Task	Turns	Final Satisfaction	Final Emotion	RAG
technology_buy_a_smartphone_dialogue_15turns_without_rag_4	Ama Nkrumah	buy a smartphone	16	⊘ 1.00	grateful	Disabled
healthcare_refill_my_prescription_dialogue_15turns_without_rag_1	Sofia Torres	refill my prescription	16	⊘ 1.00	grateful	Disabled
housing_find_an_accessible_home_dialogue_15turns_without_rag_4	Chidi Okoro	find an accessible home	16	⊘ 1.00	neutral	Disabled
Caregiver support_find_cognitive_exercises_for_dementia_prevention_dialogue_15turns_without_rag_3	Ayodele Okeke	find cognitive exercises for dementia prevention	16	⊘ 1.00	happy	Disabled
daily living_set_medication_reminders_dialogue_16turns_without_rag_6	Omar Khalil	set medication reminders	16	⊘ 1.00	happy	Disabled
healthcare_find_a_caregiver_for_an_elderly_person_dialogue_15turns_without_rag_3	Amina Khalil	find a caregiver for an elderly person	16	⊘ 1.00	grateful	Disabled
daily living_arrange_transportation_to_a_clinic_dialogue_15turns_without_rag_8	Omar Khalid	arrange transportation to a clinic	16	⊘ 1.00	grateful	Disabled
Stechnology_reset_an_online_passwormdialogue_15turns_without_rag_2	Camila Torres	reset an online password	16	⊘ 1.00	neutral	Disabled
daily living_order_groceries_online_dialogue_16turns_without_rag_10	Liang Chen	order groceries online	16	▲ 0.50	neutral	Disabled
Caregiver support_book_a_nurse_for_my_father_dialogue_15turns_without_rag_1	Anita Chen	book a nurse for my father	16	⊘ 1.00	happy	Disabled
housing_rent_an_apartment_dialogue_15turns_without_rag_10	Amir Al-Hassan	rent an apartment	16	⊗ 0.40	neutral	Disabled
S technology_teach_my_parent_to_use_video_calls_dialogue_15turns_without_rag_2	Hiroshi Tanaka	teach my parent to use video calls	16	⊘ 1.00	grateful	Disabled
technology_buy_a_smartphone_dialogue_15turns_without_rag_8	Kwame Okafor	buy a smartphone	16	⊘ 1.00	grateful	Disabled
Charles and the second se		· · · · · · · · · · · · · · · · · · ·		A A 70		

Figure 16: File detail view within folder analysis.

⁷⁶⁸ file-level satisfaction per turn (Figure 18),



Figure 17: Folder detail satisfaction overview.

⁷⁶⁹ emotion statistics (Figure 19),



Figure 18: Satisfaction per turn analysis in folder detail.

and explanations for metrics, which can be expanded to show details (Figure 20).

	ysis ,ii Tur	n Correlation	E cmos	onal States												
Emotional States	by Turn															
Emotional State	Turn 1	Turn 2	Turn 3	Turn 4	Turn 5	Turn 6	Turn 7	Turn 8	Turn 9	Turn 10	Turn 11	Turn 12	Turn 13	Turn 14	Turn 15	Turn 16
happy	10	21	29	26	28	35	26	28	29	23	35	31	34	37	38	34
neutral	438	409	408	393	369	336	306	263	251	233	199	208	205	200	198	207
grateful	0	8	1	16	47	71	112	154	171	189	209	203	202	209	203	199
interested	3	9	12	16	7	5	2	2	2	5	1	7	5	4	10	8
surprised	0	2	2	0	0	2	2	2	0	0	3	1	3	0	1	2
hopeful	0	0	0	0	1	4	4	4	0	3	4	2	3	3	3	3
	planatio	n														
O Metrics Ex	the second second second second		2													
O Metrics Ex Detailed explanation of	how metrics a															

Figure 19: Emotion statistics in folder detail analysis.

Satisfaction Metrics	^	
Final Satisfaction The satisfaction score at the end of the dialogue, representing how satisfied the user was with the overall interaction.		
dialogue.metadata.final_hidden_states.satisfaction.score		
Scale: 0.0 to 1.0, where 1.0 is completely satisfied and 0.0 is completely unsatisfied.		
Average Satisfaction		
The average satisfaction score across all dialogues in a folder.		
Sum of all final satisfaction scores / Number of dialogues		
Example: If 3 dialogues have scores of 0.8, 0.9, and 0.7, the average is 0.8.		
Satisfaction Trend		
The change in satisfaction from the beginning to the end of the dialogue.		
Final satisfaction score - Initial satisfaction score		
Example: If satisfaction started at 0.6 and ended at 0.9, the trend is +0.4.		
Success Rate		
The percentage of dialogues that are considered successful (satisfaction score ≥ 0.8).		
Number of dialogues with satisfaction \ge 0.8 / Total number of dialogues		
Example: If 7 out of 10 dialogues have a satisfaction score it 0.8, the success rate is 70%.		
Emotional State Metrics	~	
Intent State Metrics	*	
Message Length Metrics	÷	

Figure 20: Metric explanations section with expandable details.

771 —

772 Batch Analysis Mode

⁷⁷³ First, select the profiles you need at Figure 21 (example shows first user from three folders

⁷⁷⁴ selected). Scrolling down will show comparative analysis of these dialogues.



Figure 21: Profile selection for batch comparative analysis.

Next, you can view emotional states for these users (Figure 23),



Figure 22: Batch comparison of multiple dialogue profiles.

and scroll further to clearly compare dialogue differences by turn for the same user interact-

ing with different models (Figure 24).



Figure 23: Emotional states comparison for multiple users.



Figure 24: Detailed dialogue turn comparison across models for the same user.

⁷⁸⁰ which helps analyze differences better.

⁷⁷⁸ When switching back to the original dialogue lists with "View" (Figure 25), the left side ⁷⁷⁹ shows the selected dialogues, and the right side shows the multi-dialogue comparison,



Figure 25: Side-by-side view of selected single and multi-dialogue comparisons.

⁷⁸¹ This corresponds to the Split View layout (Figure 26).

Use the folder comparison vi Add tags to dislogues to call Select multiple dialogues to c Export analysis data for furth	ew to compare metrics ac igorize them by specific of compare them side by side er processing in other tool	ross different folders wascleristics b B									
🗅 caregiver support (3) 🛛 🗁 cir	aude without 0 (150)	🗈 claude without 40 (150)	🗈 claude-3.7-	sonnet with 0 (15	0)						
Dislogue Dislogue	.∦ Metrics) • ny_father_dialogue_15t	A User Profile	K ² , Sele 3 dialogue & Full	cted Dialogu s selected for con Comparison	es iparison Sa	lisfaction	D Mese	age Length	C Emotional State	× 2	User Profiles
Anita (1900) Chen (1900) Chi Starra Sti 42025	-3.7 scenet with 0 o scene RAG: Disabled	de-3.7- st x gpt-fo- mini	Ray infor	er Profile	Summa	oomparison	_	Tech		Time	Difficulty
Caregiver support_book_a_m Final Outcome	area_for_my_father_claicga	e_1Sturms_without_rag_1	Anita	elaste - vittect 0	Task book a nurse for mu	Eastern	Personality Outgoing	Experience	Communication	Constraint Very	Level
Final Emotional State happy	Final Intent State leaving	Pinal Satisfaction	Anita	claude • without 40	father book a nurse for my	Eastern	Outgoing	Novice	Adaptive	Very	1
1 R Anita Chen Hi there! I'm looking help.	for a nurse for my dad	 12:53:10 AM who needs some daily 	Anita	elaude 3.7- - secred with a rise	father book a nurse for my	Eastern	Outgoing	Novice	Adaptive	Very	1
User's Emotional States	a Intent Ó Sat Les r Errotionat Satisfie	isfaction & Message ligth clore © 0.5			father						
intert inter	r Intent: User m exploring chars	isfaction information provided" essege: 71 Assistant message: 120 chars	🖒 Sa Satistact	tisfaction on metrics across	Summa dalogues	iry					
() User's Inner Tho "No isner thoughts re	ughts conded"		di Avi	1 75-		c	Final Satisfacti 1 0.75 -	50 	-* Satisfa 0.6 0.45	ction Trend	
 Assistant Hi Anital I'd be happ 	y to help find a nurse fo	© 12:53:10 AM	0	25-			0.5-		0.3-		

Figure 26: Split view for detailed analysis.

782 —

783 Folder-Level Comparison

⁷⁸⁴ Click the "Folder Comparison" button at the top right to open the component (Figure 27).

785 You can then select two folders to compare.

STORM alyze and compare dialogue data ac	cross 4 folders with 453 dialogues		Si Grid View	Spilt View 11 Folder Comparis	50n 👌 Upload Data
Tip: Hover over folder names to rena	me them. Consider adding specific co	nditions or test details (e.g., "claude-3-sonne	al-temperature-0.7" or "gp1-4-with-rag") to make them easier to identify.	
③ Getting Started					
Quick guide to using the dialogue analys	iis dashboard				
Dipload Data Start by uploading JSON files contain	ing dialogue data. You can upload	Croanize & Tag	tags to categorize them. Rename	Jr. Analyze & Compare View metrics, compare dialogues,	and analyze performance acros
Individual files or entire folders.		folders to better identify different test or Rename Tag Group	orditions.	different models or test conditions Metrics 11 Company	
Select multiple dialogues to compa Export analysis data for further pro- corregiver support (3) En claude v	ine them aide by aide cossing in other tools without 0 (150)) (In claude withou	ut 40 (150) 🕐 claude-3.7-sonnet with 0	0 (150)		
Salect multiple dialogues to compa Export analysis data for further pro- caregiver support (3) (E) claude v TI Folder Comparison Geregare metrica between two folders	re than aids by aids cossing in other tools withhout 8 (156) 📄 😂 claude withou	rt 40 (150) 🕒 clausie-3.7-sonnet with 0	9 (169)		0 4
Select multiple dislogues to compa Export antigitie dislogues to compa caregiver support (3) Ex claude w F1 Folder Comparison Compare matrice between two folders extractional dislogues claude-37-somest with 0 (150)	re than aids by side cossing in other tools without 8 (159) Cr claude withou	در 40 (150) (<u>Cs claudo 3.7-sonnet with</u> 0	0 (150) daude without 0 (150)		ٹ 0
Since multiple dialogues to compare Export analysis data for lutther pro compare multiple data for lutther pro the dialogue data for lutther pro compare multiple dialogue data extra data data data data extra data data data data data data extra data data data data data data data d	without 0 (156) It claude without 0 (159) It	در 40 (150) 🕞 دامیداره ۲۰۰۵ مربع بالله ۲۰ - بر ۲۰ 2 Editional (9 (159) daude without 0 (150)		0 3
Sidect multiple dislogues to compare Export analysis data for lutther pro- complex respont (3) (b) claude of Compare meticas behaviors who (1920) Compare meticas behaviors who (1920) Coducts 37-served with 0 (1930) Coducts 37-served with 0 (1930) Coducts 37-served with 0 (1930)	ne frama skál by kté ossang in offer tools without 0 (156)) 🗠 claude withou spane) 🕞 claude without (153) (155 States Message Anstynis User (1	v (16) (************************************	9 (10) • dasis without 0 (150)		ى 0
Electric multiple discipants is compare Electric multiple discipants is compare Compare multiple displayed displayedisplayed displayed displayed displayed displayed displayed displa	vie them alide by Atle consense in other tools without 0 (154) In claude without spanet) (In them without 0 (120) (120 States Message Analysis Lawri	در 4.40 (150) (<u>c claude 27 4000000 utility</u>) ∨ 2 distegane) Profiles	e (16) • daste without 0 (155)		© &
Sector analysis data for further pro- Decord analysis data for further pro- complex networks (1) (1) (2) calculate The Folder Comparison Compare metrics between lash biddes Control and analysis (1) (2) (2) (2) (2) Control and analysis (2) (2) (2) (2) (2) (2) Control and analysis (2) (2) (2) (2) (2) (2) (2) (2) (2) (2)	via Nama kida by Atle constang in other tools withhout 0 (155) E claude withhou genne) ((2) claude withhoud 0 (150) (152 (20) (20) (20) (20) (20) (20) (20) (20)	a 44 (10) Co cloude 3.7 escret 400 Co Collegence Profes	e daute without 0 (150)		ى 0
Extended and states at several a constraint of the function of the several and severa	In them slads by side constains in other tools without 6 (156) in the claude without general) (In stands without 6 (155) (155 States Message Analysis Deer claudes-3.7-sonnet with 0 (156)	A 44 (10) (2 schold 27 sevent kills) (descent) (10) Consent claude willhout (0 (15))	dask without 0 (155) Key Metrics Comparison		0 4
Edit andigis attempts to employ Export analysis data for theme pro Export analysis data for theme pro activity of the provided attempts activity of the provided at	without 6 (156) be deade without 6 (156) (15 Grand Strand	A 49 (19) E talanta 2.2 Annose talah calangan bargan	e daule without 0 (155) A Key Metrics Comparison Average Satisfaction		6 ±
Start range data between termines Depet anyone data by the property expects (1) Depet anyone d	na here add by state conserts in other factors without (15%) in cleaned without segmes). (2.10mm without (12%) (2%) Editation Mensage-Andrya's Liber (daude 3.7.400met witho (15%)) 120 2400	A 48 (150) (2 124006 2.2.42000 110) (21000000) (21000 110) (21000000) (21000000) (2100) 120 2000	e daude without & (15) • daude without & (15) • Koy Medrica Comparison Average Bastalactico High Bastalactico		6 ±
Start andre skappet to thempe Constraints with the themp	en them active years and and a (1951) (2) of claude without employed (1951) (2) of claude without employed (2) of claude without embedde status Wessage Actives (claude claude - 3.7-secret with (1950) 1950 1950 1950	4 4 (15) (citade 2 2 annue 11)) (disclorence) (d	Control Contro Control Control Control Control Control Co		0 4

Figure 27: Folder comparison selection interface.

- 786 Below, detailed differences are shown, including:
- 787 Satisfaction comparison (Figure 28),

Compare metrics between two folders					() ± Eq
 claude-3.7-sonnet with 0 (150) 		~	C e claude without 0 (150)		
(In claude-3.7-sourcet with 0 (150) (150 d	lalegues)) (En claude without 0 (150) (150	dialogan)			
Overview Satisfaction Emotion	al States Message Analysis User F	Profiles			
1 Satisfaction Metrics			Satisfaction Trend Metrics		
	claude-3.7-sonnet with 0 (150)	claude without 0 (150)		claude-3.7-sonnet with 0 (150)	claude without 0 (150)
Average Satisfaction	0.91	0.83	Improved Satisfaction	14 89.3%	14 75.3%
High Satisfaction Rate (20.8)	⊘ 86.0%	♂ 72.0%	Stable Satisfaction	6.0%	15.3%
Moderate Satisfaction Rate (0.5-0.8)	▲ 9.3%	△ 18.7%	Declined Satisfaction	11 4.7%	11.9.3%
Low Satisfaction Rate (<0.5)	⊙ 4.7%	⊙ 9.3%			
La distanti a Martin Comme					
In Satisfaction Metrics Comp.	rison				
In Satisfaction Metrics Compo Average Satisfaction High Satisfaction Rate	rison				•
J Satisfaction Metrics Compo Average Satisfaction High Satisfaction Rate- Moderate Satisfaction Rate	rison				
Ji Satisfaction Metrics Compo Average Satisfaction High Satisfaction Rate Moderate Satisfaction Rate Low Satisfaction Rate	rrison				
In Satisfaction Metrics Comp Avarage Satisfaction High Satisfaction Rate Moderate Satisfaction Rate Low Satisfaction Rate Improved Satisfaction	rison				
Ar Satisfaction Metrics Comp Average Satisfaction High Satisfaction Rate Moderate Satisfaction Rate Law Satisfaction Rate Improved Satisfaction Stable Satisfaction	rison				
Aurage Satisfaction Metrics Compo Average Satisfaction High Satisfaction Rate Metarate Satisfaction Rate Low Satisfaction Rate Improved Satisfaction Stable Satisfaction Declined Satisfaction					

Figure 28: Satisfaction comparison between folders.

- Emotional states comparison (Figure 29),



Figure 29: Emotional states comparison between folders.

⁷⁸⁹ - Message length comparison (Figure 30),



Figure 30: Message length comparison between folders.

⁷⁹⁰ - User profile comparison (Figure 31).



Figure 31: User profile comparison between folders.