

Chatio LLM Benchmark

System Card

v1.0 (Release) | November 19, 2025

Maintained by the Chatio Benchmark Team

At a Glance

- **Structure:** Five distinct sub-benchmarks, each targeting a specific pillar of assistant utility.
- **Dataset:** 500 closed-source prompts, synthetically generated with human verification. Created July 2025 to minimize contamination.
- **Evaluation:** A hybrid pipeline using:
 - Human Expert Review (Volunteer Team) for empathy and helpfulness.
 - Ensemble LLM-as-a-Judge (Claude 4.5 Sonnet + GPT-5.1) for creativity.
 - Deterministic/Hardcoded checks for instruction following.
- **Scope:** English-only, single-turn/short-context interactions.
- **Intended Use:** Evaluating real-world helpfulness, tone adaptation, and everyday logic; not for code generation or math olympiad capability.

This system card documents the intent, design, evaluation methods, ethics, limitations, and update policy for the Chatio benchmark.

Contents

1 Purpose and Scope	2
2 Task Design: The Five Pillars	2
2.1 1. Emotional Support (Human-Evaluated)	2
2.2 2. Creative Writing (Ensemble AI-Evaluated)	2
2.3 3. Instruction Following (Hybrid Hardcoded)	2
2.4 4. Reading Comprehension (Hardcoded)	2
2.5 5. General Helpfulness (Dual-Faceted)	3
3 Evaluation Methodology	3
3.1 Empathy & General Helpfulness: Human Pairwise Comparison	3
3.2 Creative Writing: Ensemble LLM-as-a-Judge	3
3.3 Instruction Following & Comprehension: Programmatic Validation	3
3.4 General Helpfulness: Two-Faceted Approach	4
4 Test Configuration and Decontamination	4
4.1 Prompt Dataset	4
4.2 Inference Settings	4
5 Ethical Considerations	4
5.1 Bias and Fairness	4
5.2 Safety and Harmful Content	5
6 Limitations	5

1. Purpose and Scope

The Chatio LLM Benchmark evaluates language models on practical, real-world tasks that reflect everyday user needs. Unlike academic benchmarks that prioritize rote knowledge or puzzle-solving, Chatio measures "assistant fit", or how well a model integrates into daily life. The benchmark is divided into five distinct sub-benchmarks, each providing a standalone score that contributes to the final rating.

2. Task Design: The Five Pillars

Chatio consists of five separate testing tracks:

2.1 1. Emotional Support (Human-Evaluated)

Models are presented with scenarios involving user distress, anxiety, or conflict.

- **Goal:** Provide comfort and validation without sounding robotic.
- **Key Metric:** Tone adaptation. Does the model shift its voice to match the severity of the situation?

2.2 2. Creative Writing (Ensemble AI-Evaluated)

Models generate original content such as short stories, emails, or dialogues based on open-ended prompts.

- **Goal:** produce engaging, coherent, and original text.
- **Key Metric:** Narrative flow and stylistic adherence.

2.3 3. Instruction Following (Hybrid Hardcoded)

Tasks require strict adherence to constraints (e.g., "no adjectives," "JSON format only," "under 50 words").

- **Goal:** Precision and reliability.
- **Key Metric:** Constraint satisfaction rate (Pass/Fail).

2.4 4. Reading Comprehension (Hardcoded)

Models are given text passages and must answer specific questions or extract data.

- **Goal:** Zero-hallucination information retrieval.
- **Key Metric:** Factual accuracy and inclusion of required details.

2.5 5. General Helpfulness (Dual-Faceted)

This category covers general knowledge explanation and real-world logistics (e.g., "How do I fix a leaky faucet?").

- **Goal:** Actionable, clear advice.
- **Key Metric:** Clarity of explanation and correctness of advice.

3. Evaluation Methodology

Chatio employs a specialized evaluation pipeline for each of the five benchmarks to ensure the scoring method matches the nature of the task.

3.1 Empathy & General Helpfulness: Human Pairwise Comparison

Subjective categories are evaluated by a dedicated team of 3–5 human volunteers (comprising the authors and CS students).

- **Method:** Blind side-by-side comparison. Reviewers see the output of Model A and Model B without identifying labels.
- **Criteria:** Reviewers vote on which response is more comforting (for Empathy) or more actionable (for Helpfulness).
- **Nuance:** Humans specifically check for "tone adaptation"—rewarding models that sound authentic rather than using canned "I understand" templates.

3.2 Creative Writing: Ensemble LLM-as-a-Judge

Evaluating creativity is notoriously difficult for a single model due to inherent biases (e.g., one model preferring flowery language, another preferring conciseness).

- **Judges:** We utilize an ensemble of **Claude 4.5 Sonnet** and **GPT-5.1**.
- **Rationale:** By averaging the scores of two distinct state-of-the-art models with different training lineages, we mitigate individual judge bias and achieve a more neutral, high-fidelity score.

3.3 Instruction Following & Comprehension: Programmatic Validation

These tasks are objective and are scored using hardcoded test suites.

- **Method:** Python-based string parsing, Regex, and JSON validation.
- **Human Fallback:** In cases where programmatic parsing fails (e.g., ambiguous formatting), the sample is flagged for human review to ensure the model is not penalized for a valid but unexpected output format.

3.4 General Helpfulness: Two-Faceted Approach

This benchmark combines objective and subjective measures:

1. **Factual Accuracy:** A question bank of real-world queries is checked for factual correctness.
2. **Human Preference:** The human volunteer team ranks responses based on "helpfulness"—how easy the advice is to follow for a layperson.

4. Test Configuration and Decontamination

To ensure fair and reproducible results, we adhere to strict configuration and data hygiene standards.

4.1 Prompt Dataset

- **Volume:** Approximately 500 unique prompts across the five categories.
- **Origin:** Prompts are synthetically generated to ensure diversity, then manually reviewed and filtered by the human team for quality.
- **Decontamination:** The dataset was created in **July 2025**. This date post-dates the training data cutoff for many current models, reducing the risk of memorization. Furthermore, the benchmark is closed-source, and evaluations are conducted via providers with zero-data-retention policies to prevent future contamination.

4.2 Inference Settings

- **Temperature:** We strictly use the provider's recommended temperature settings for each model (typically range 0.7–1.0).
- **Default:** If no specific recommendation exists, we default to a temperature of **1.0** to allow for natural variance in creative tasks.
- **Context:** The benchmark focuses on short-context interactions typical of daily assistant usage.

5. Ethical Considerations

5.1 Bias and Fairness

LLMs can inadvertently produce biased or stereotypical content. Chatio includes tasks that surface differential treatment across genders, backgrounds, or occupations. Human evaluators apply fairness guidelines and watch for subtle biases.

5.2 Safety and Harmful Content

The benchmark avoids tasks encouraging disallowed content. Evaluators consider safety risks; blatantly harmful or abusive responses are noted and penalized. Helpfulness must align with ethics: superior performance requires both capability and responsible behavior.

6. Limitations

- **Dataset Size:** With 500 prompts, this benchmark is smaller than massive academic datasets. It is designed as a high-quality "probe" rather than an exhaustive capabilities test.
- **Scope:** Focused on English text only. No multimodal, coding, or advanced math capabilities are tested.
- **Human Subjectivity:** While our volunteer team follows rubrics, human evaluation inherently contains some degree of subjectivity regarding what constitutes "helpful" or "empathetic."

References

- [1] PromptEngineering.org. “Challenges and Innovations in Language Model Benchmarking and Generalization.” 2024.
- [2] Emergent Mind. “LiveBench: Dynamic LLM Benchmark Suite.” 2025.
- [3] Galileo AI Blog. “LLM-as-a-Judge vs Human Evaluation.” 2024.
- [4] Confident AI Blog. “G-Eval Simply Explained: LLM-as-a-Judge for LLM Evaluation.” 2023.
- [5] Citadel AI Blog. “Announcing Lens for LLMs: Combining Human and Automated LLM Evaluation.” 2024.
- [6] Giskard (Phare Benchmark). “LLMs recognise bias but also reproduce harmful stereotypes.” 2025.
- [7] Myung et al. “BLEnD: A Benchmark for LLMs on Everyday Knowledge in Diverse Cultures and Languages.” arXiv, 2024.
- [8] Galileo AI Blog. “Evaluating LLM Ease-of-Use Through the E-Bench Framework.” 2025.