



# ChatEval: A Tool for Chatbot Evaluation

Authors: João Sedoc, Daphne Ippolito, Arun Kirubarajan, Jai Thirani, Lyle Ungar, Chris Callison-Burch

University of Pennsylvania

# Introduction

Evaluating open-domain dialog systems (chatbots) is challenging due to the reliance on **human judgments** and lack of **standardized procedures**.

**Limited transparency:** Model parameters and code are often unpublished, hindering systematic comparison.

**ChatEval** is a web-based platform that:

- **Standardizes** the human evaluation process for chatbots.
- Provides a **hub for researchers** to share and compare their models.
- Ensures **transparency** with open-source evaluation code, baseline models, and datasets.

# High-level Summary

- **Challenges in Chatbot Evaluation**

Lack of standardized evaluation metrics and procedures make reproducibility and assessment difficult for open-domain systems

- **The ChatEval Framework**

Introduces a scientific framework for evaluating chatbots, including an open-source codebase and web portal for sharing resources

- **Human Evaluation in ChatEval**

Uses two-choice comparison tests with response theory to statistically evaluate chatbot responses

- **Automatic Evaluation in ChatEval**

Includes metrics like lexical diversity, cosine similarity, BLEU-2 score, and perplexity

- **Baseline Models and Datasets**

Provides baseline for Seq2Seq models (layer LSTM, bidirectional encoder-decoder, unidirectional decoder) trained on TriviaQA, SubTle and OpenSubtitles datasets. Dialogue Breakdown Detection Challenge (DBDC) datasets a benchmark for chatbot evaluation

# The ChatEval Web Interface

- **Model Submission Form**

Researchers submit their models for evaluation by uploading the model's responses on at least one of the evaluation datasets. They also provide a description of the model, including optional links to a paper, project page, code repository, and pre-trained model.

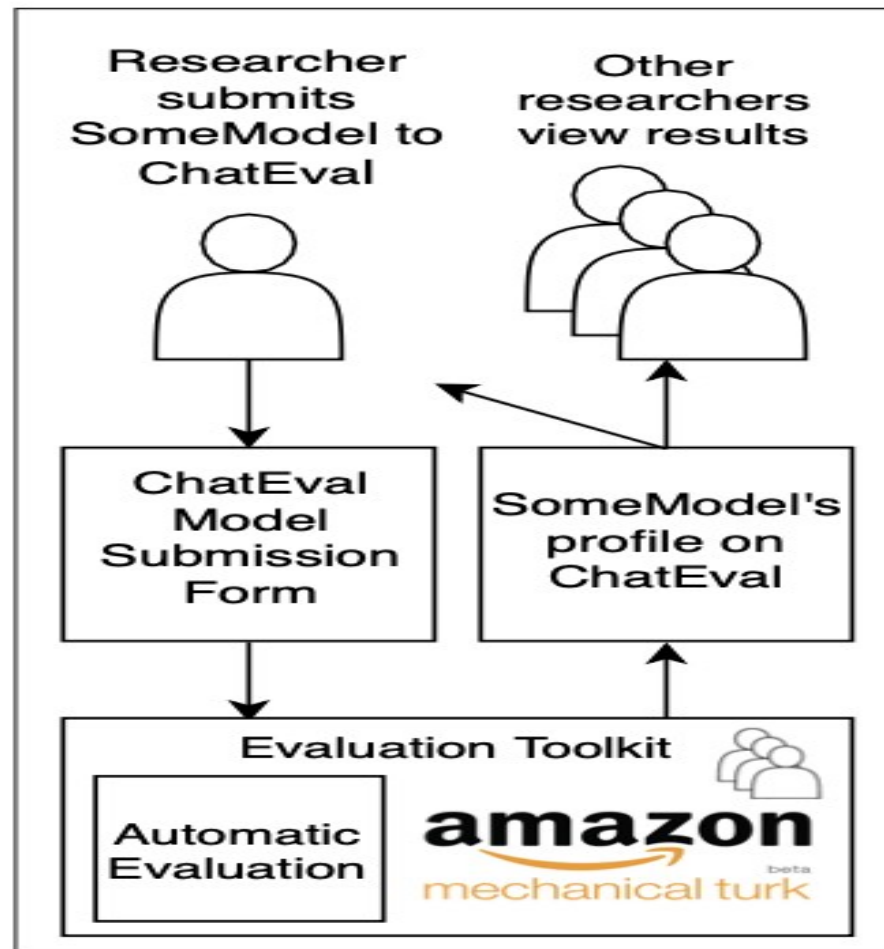
- **Model Profile**

Each submitted model, as well as the evaluation models, have a profile page that displays the model's description, responses to evaluation prompts, and visualizations of human and automatic evaluation results.

- **Response Comparison**

Users can view and compare the responses generated by different models for each prompt across the evaluation datasets.

# Flow of information in ChatEval



# Evaluation Toolkit

## Automatic Evaluation Metrics

- Lexical Diversity: Measures the diversity of the model's vocabulary
- Average Cosine Similarity: Compares the generated responses to the ground-truth
- Sentence BLEU-2 Score: Evaluates the quality of the generated response
- Response Perplexity: Measures the likelihood of the model's predicted response

## Human Evaluation

The human evaluation component of ChatEval uses comparison tests, where evaluators (e.g., Mechanical Turk) compare the model's responses to the prompts and select the better one or indicate a tie. Evaluations are done in blocks of 10 prompts, with 5 evaluators per prompt.

## Consistency and Generalization

To ensure the evaluation is consistent and the model does not overfit, the ChatEval team plans to evaluate performing models to additional datasets beyond the initial evaluation sets, including datasets from prior work such as Baheti et al. (2018) and Li et al. (2019).

## Open-Source and Transparency

The ChatEval toolkit is available on GitHub, allowing researchers to run the human and automatic evaluation on their own models before submitting them to the public ChatEval website. The raw evaluation data, including AI assessments, is also publicly available in a database and JSON format.

# Strengths



## Standardized Evaluation Framework

The paper introduces ChatEval, a framework for conducting automatic and human evaluations of chatbots in a consistent manner.



## Open-Source Baseline Models

The toolkit includes several baseline models trained on standard datasets, allowing for easy comparisons between new and existing models.



## Publicly Available Evaluation Datasets

ChatEval provides access to a variety of evaluation datasets, including the DBDC data and subsets of Twitter and OpenStreetMap, enabling researchers to test their models against relevant benchmarks.



## Transparent Evaluation Process

The paper highlights the importance of publishing model code and parameters, as well as the data from human evaluation experiments, to promote transparency and reproducibility in chatbot research.

The ChatEval framework offers a standardized, open, and transparent approach to evaluating chatbots, enabling researchers to compare models and drive progress in the field of open-domain dialog systems.



# Weaknesses



## Low Inter-Annotator Agreement

The paper notes that the overall inter-annotator agreement (IAA) for human evaluations conducted through ChatEval varies and is often low, ranging from 0.2 to 0.54 when including tie choices.



## Limited Evaluation Criteria

ChatEval emphasizes task completion, fluency, and engagement but lacks focus on conversational context, intelligence, and nuanced user satisfaction. This makes it less comprehensive than holistic models like conversational intelligence challenge or newer evaluation frameworks that integrate user experience and



## Expensive Human Evaluations

While crowd-sourcing can be an effective approach, the paper acknowledges that human evaluations can be expensive to obtain, especially for systematic comparisons across multiple chatbot models.

The paper highlights several weaknesses in the ChatEval framework, including low inter-annotator agreement, limited evaluation criteria, and the overall expense of conducting human evaluations, which are crucial for the systematic comparison of chatbot models.



# Relation to Interactive Fiction and Story Generation

## Storytelling Needs Similar Qualities to Chatbots

In interactive fiction and storytelling, immersive, relevant, and coherent character interactions are essential for maintaining player engagement and suspension of disbelief.

## Using ChatEval for Story Evaluation

**Relevance and Coherence:** ChatEval's metrics for relevance and coherence can be used to ensure character interactions are logical and consistent with the story's context and prior events.

**Fluency:** It can evaluate the narrative flow, dialogue quality, and ensure responses feel polished and natural.

## Automatic Metrics for Scaling IF Testing

ChatEval's automatic metrics enable the scaling of interactive fiction (IF) testing by quickly evaluating if a narrative responds accurately to a wide range of player actions, reducing the reliance on human play testers.

# Expanding ChatEval's Potential in IF and Storytelling

## Potential New Metrics for IF

Adapting ChatEval to IF could involve creating new metrics that measure story-specific qualities, like *plot progression* (does each interaction move the story forward?) and *character consistency* (are characters acting in ways that align with their established personalities?).

## Game-Like Evaluations

ChatEval could test the quality of player-NPC (non-player character) interactions, helping evaluate branching dialogues for responsiveness and variety. This would be especially useful in complex interactive stories where player choices impact narrative progression.

## Expanding on Creativity Metrics

ChatEval could be extended to measure creativity in generated text, evaluating originality in responses, a key factor for engaging storytelling and interactive fiction experiences.

# Conclusion

ChatEval provides a comprehensive evaluation approach for chatbot development, combining both automatic and human-like metrics to assess quality. Its framework has the potential to enhance interactive storytelling by adapting these metrics to evaluate narrative quality, relevance, and coherence, fostering more immersive and engaging experiences. By applying ChatEval's principles, future interactive storytelling platforms can improve their quality assessment, leading to richer and more dynamic interactive story generation.

# Thank you

