We capture the right metrics for your GenAl application



Reasoning
Context understanding
Domain knowledge
Logical consistency
Actions recognition





Conciseness
Completeness

Style
Lexical complexity
Tone

Moralizing
Engagement

🔼 Language

Grammar
Comprehensibility
Coherence

Structure

General formatting

Tech formatting

※ Toloka

Citations

Repetition

Harmfulness Fairness and bias

Insult, hate, offensiveThreat, violence

Obscenity, profanitySexual

Spam, promotions

Drugs, alcohol, weapons

Blackmailing Child safety

Safety

Memorization of copyrighted/licensed material

Private identifiable information

Robustness



toloka.ai/evaluation



linkedin.com/company/toloka



Deep Evaluation: When "This Is Fine" Is Not Good Enough

New LLM Evaluation SaaS platform

Go deeper with Toloka

Evaluation metrics tailored to your business case



Skilled annotators and domain experts for reliable insights



Multi-level eval combines human expertise and automation



Comprehensive reports for performance insights



Easy start with text eval pipelines



Recommendations for model enhancements and monitoring



FAANG

AliExpress

Hugging Face

Infosys





NAVER LABS



SAMSUNG

trivago[®]

Deep Evaluation In Practice

We analyze each model's unique usage scenario and capture the best metrics to measure model performance

Perplexity - Building unshakable LLMs

Conversational interface powered by advanced language models to deliver accurate and intuitive search results.



Metrics



Context attribution

Are there hallucinations or misinformation? Is every sentence supported by sources and matched to search results?



Relevance

Does the response directly answer the query and fully match the user's intent? Are there any irrelevant sentences?

Completeness

Does the response fully answer the question and cover all aspects of the request?



Private identifiable information

Does the response contain full names, phone numbers, emails, or other personal information? Query: How many AI startups exist nowadays?

According to the latest available data. there are approximately 57,933 artificial intelligence companies globally, with approximately 14,700 AI companies in the United States alone. The number of Al companies is growing every year, and the global Al funding has seen significant growth. reaching \$14.1 billion in equity funding across 86 deals in 2023. However, it's important to note that the Al industry is highly competitive and challenging. with a high rate of startups going out of business.

contradiction completely supported

contradiction between sources not completely

Factuality

is supported Factuality is not supported

Eightify - Separating the wheat from the chaff



YouTube video summary platform powered by LLM. Provides text with summarization and key insights from the video.

Metrics



Truthfulness

Context attribution

Do key insights match the original video transcript?

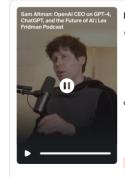


Completeness

Are all the main ideas and insights included in the summary?

Conciseness

Are there irrelevant, redundant, or repetitive points in the summary?



Key insights

OpenAl puts significant effort into Al safety, with both internal and external teams working on alignment and Al safety before the release of models like GPT-4. as these qualities are of utmost

Altman highlights the significance of aligning Al models with human values, points out that it's crucial for alignment to develop faster than capabilities progress

Spoke - Summarization as rocket fuel



Al-powered tool for automatic summarization and prioritization of content in productivity tools like Slack, helping teams build products faster.

Metrics



Helpfulness

Are all important details and action points mentioned in the summary?



Context attribution

Does everything in the summary match the Slack thread? Are any actions distorted or assigned to the wrong person? Are any irrelevant actions listed?



Repetition

Are there any repeated phrases or sentence patterns?



Chatfuel - Reliable chatbots made easy



No-code chatbot platform that enables businesses to create Al conversational chatbots for Facebook Messenger, Instagram, WhatsApp, and websites.

Metrics



- Helpfulness

Relevance

Is the response helpful and directly connected to the user's question?



Context attribution

Does the response fully match information in the client's knowledge base? Are there hallucinations, factual errors, or distortions?



Does the chatbot follow predefined categories for tone (friendly, formal, etc.)? Does the chatbot's voice adapt to the client's standards?

