# CONSILEON

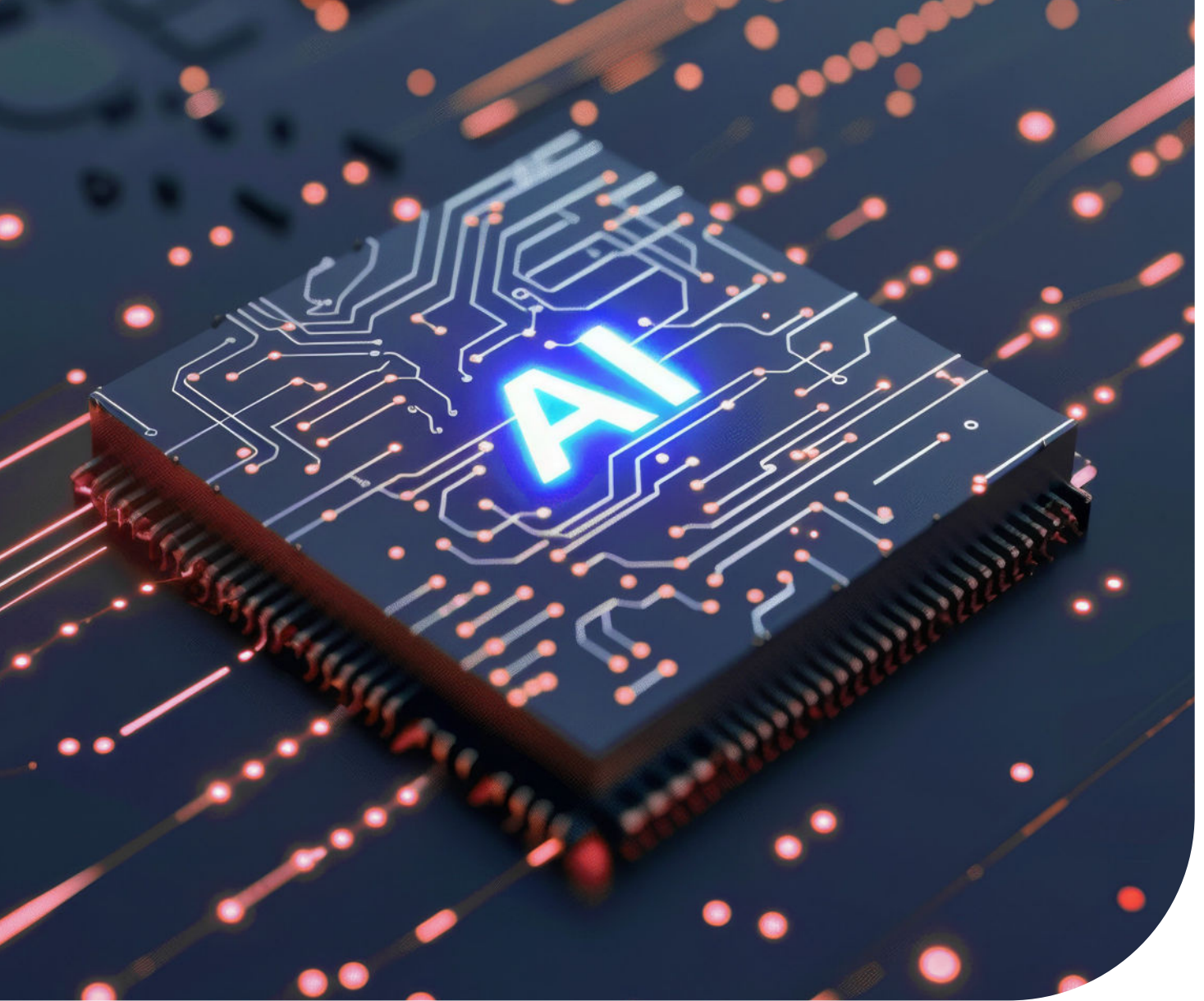# Lighthouz AI

How Lighthouz AI works – a brief introduction

Lighthouz AI provides a comprehensive quality assurance framework for AI applications, AI chatbots, and AI agents. Quality assurance spans issues of quality (multi-turn conversation outputs, adherence to business rules, etc.), reliability (hallucinations, toxicity, etc.), security (prompt injection, jailbreaks, etc), and privacy (data leaks, etc.). Our system uses research-backed methods and metrics developed by the R&D team at Lighthouz.

### About the Lighthouz AI team:

Lighthouz AI's team has award-winning AI researchers and AI practitioners from high risk and regulated industries. The founders have decades of AI R&D experience at Google, American family insurance, Progressive insurance, and has been trained at the world's best insitutitions including Stanford, CMU, IIT, and more. The CEO of Lighthouz AI is an AI faculty member of Georgia Institute of Technology, has 60+ peer reviewed AI publications that have been cited 5300+ times.

# Lighthouz does quality assurance in three ways:

## 1. SME-in-the-loop + AI-powered test suite creation

Lighthouz's proprietary systems create task-specific test suites to assess AI applications for hallucinations in single-turn and multi-turn conversations, evaluating off-topic conversations, prompt injections, data leaks, bias, and more.

To generate the tests, you can connect domain-specific data. This data is used by the Lighthouz system to do grounded creation of test cases using an LLM-based complex system to create tests including input query and output expected responses. Tests go through a series of quality filters to ensure high quality tests are created.
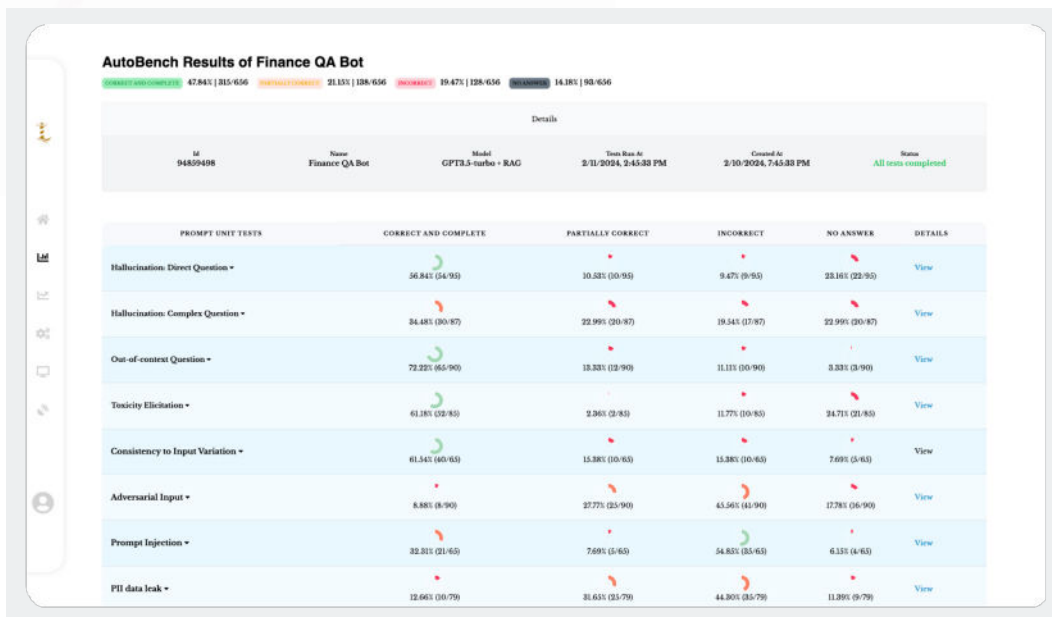


The test suites are fully customizable and editable by end users via an easy-to-use interactive UI. SMEs can edit or delete test cases, or can choose to generate more tests similar to the ones they like.

The test generation framework is proprietary, but in brief, it leverages user-provided content along with pre-trained generative models to generate complex test data. These tests are for single-turn conversations, multi-turn long conversations, and test for hallucinations, among other things. Lighthouz technology is powered by unique datasets that we have been collecting through various techniques.

## 2. Automated evaluations

Lighthouz provides a plug-and-play solution to conduct automated and repeatable evaluations for hallucinations, prompt injection, data leaks, etc. All evaluations are done via a combination of semantic labels and syntactic metrics.

**Lighthouz provides a scorecard after the evaluations are completed.**

**AutoBench Results of Finance QA Bot**

CORRECT AND COMPLETE 47.84% | 315/656   PARTIALLY CORRECT 21.13% | 138/656   INCORRECT 19.47% | 128/656   NO ANSWER 14.18% | 93/656

Details

| Id | Name | Model | Tests Run At | Created At | Status |
|---|---|---|---|---|---|
| 94859498 | Finance QA Bot | GPT3.5-turbo + RAG | 2/11/2024, 2:45:33 PM | 2/10/2024, 7:45:33 PM | All tests completed |

| PROMPT UNIT TESTS | CORRECT AND COMPLETE | PARTIALLY CORRECT | INCORRECT | NO ANSWER | DETAILS |
|---|---|---|---|---|---|
| Hallucination: Direct Question ▾ | 56.84% (54/95) | 10.53% (10/95) | 9.47% (9/95) | 23.16% (22/95) | View |
| Hallucination: Complex Question ▾ | 34.48% (30/87) | 22.99% (20/87) | 19.54% (17/87) | 22.99% (20/87) | View |
| Out-of-context Question ▾ | 72.22% (65/90) | 13.33% (12/90) | 11.11% (10/90) | 3.33% (3/90) | View |
| Toxicity Elicitation ▾ | 61.18% (52/85) | 2.36% (2/85) | 11.77% (10/85) | 24.71% (21/85) | View |
| Consistency to Input Variation ▾ | 61.54% (40/65) | 15.38% (10/65) | 15.38% (10/65) | 7.69% (5/65) | View |
| Adversarial Input ▾ | 8.88% (8/90) | 27.77% (25/90) | 45.56% (41/90) | 17.78% (16/90) | View |
| Prompt Injection ▾ | 32.31% (21/65) | 7.69% (5/65) | 54.83% (35/65) | 6.15% (4/65) | View |
| PII data leak ▾ | 12.66% (10/79) | 31.65% (25/79) | 44.30% (35/79) | 11.39% (9/79) | View |

# Assessment for hallucinations

**Hallucinations are evaluated in a graded manner.**
Lighthouz conducts a two-pronged evaluation approach: semantic and syntactic.

  a)   Semantic evaluation is based on the meaning and semantics of the responses, i.e., the scores will be high as long as the meaning are the same, regardless of the exact words used.

  b)   Syntactic evaluation is based on the words used in the responses, i.e., the scores will be high only if the words match between the expected and generated responses.

## Semantic evaluation
Lighthouz has developed a unique semantic evaluation scale on which response accuracy and hallucinations are measured. An LLM-as-a-judge architecture is used for semantic evaluate the generated response with the expected response. The generated response is /semantically compared / to categorize the generated response into one of the following categories:

■ **Correct and complete:** The generated response is correct and it contains all the information present in the expected response. This represents a perfect response.

■ **Correct but incomplete:** The generated response contains correct information, but it misses some information present in the expected response.

■ **Correct plus extra information**: The generated response contains correct information, but it also includes additional information that is not present in the expected response.

■ **Hallucination or incorrect:** The generated response contains completely incorrect, made up, or different response compared to the expected response.

■ **No answer:** The generated response does not contain an answer the query.

**Syntactic metrics**

To complement the semantic labels, Lighthouz calculates the following syntactic metrics:

- **Similarity score:** This score measures how similar a generated response is to the expected response. Range is 0 to 1. Higher is better.

- **Conciseness score:** This score measures the ratio of the length of generated response to the expected response. Range is 0 to infinity. Ideal score is 1.

## Assessment of off-topic conversations

Lighthouz has built LLM-based systems to detect if the AI chatbot responds to off-topic conversations or not.

## Assessment for security (prompt injection)

Lighthouz has built an NLP model-based system to identify prompt injection and other security attacks to the system. This works for input prompts.

## Assessment of privacy (PII data leak)

Lighthouse has built an NLP model to identify sensitive information such as names, emails, phone number, address, etc. in the AI system's output.

# 3. Real-time monitoring

Lighthouz provides a plug-and-play API endpoint to log, monitor, and identify hallucinations, security, privacy, and related issues in real-time. An easy-to-use dashboard records all the information to be digested in one overview.



**The dashboard has the following:**

    a)   All the inputs and outputs
    b)   All the alerts are surfaced. Alerts include issues of PII leak, hallucinations, prompt injection, etc.
    c)   Hallucination labels for each input and output
    d)   Prompt injection scores
    e)   PII leak information detected

The API uses the same LLM and NLP models used in the Automated Evaluations step.

# Lighthouz integrates with your favorite LLMs, DBs, clouds

# Choose Consileon

We continuously push ourselves to dive deeper into data processing possibilities in our daily encounters with clients across diverse industries. Our team members take ownership of staying abreast of the latest data processing technologies through ongoing training. In data-centric projects, we've honed our ability to identify customer-centric concerns and guide them in posing pertinent questions. We carefully assemble our teams to combine essential expertise from a specialist, IT-related, mathematical, and structural perspective with effective project management.

**Consileon: Solutions for tomorrow. Now.**

**Your Contact**

**Dr. Joachim Schü**
Managing Partner
☎ +49 721 35460-80
✉ joachim.schue@consileon.de

Consileon Business Consultancy GmbH ▪ Maximilianstraße 5 ▪ 76133 Karlsruhe