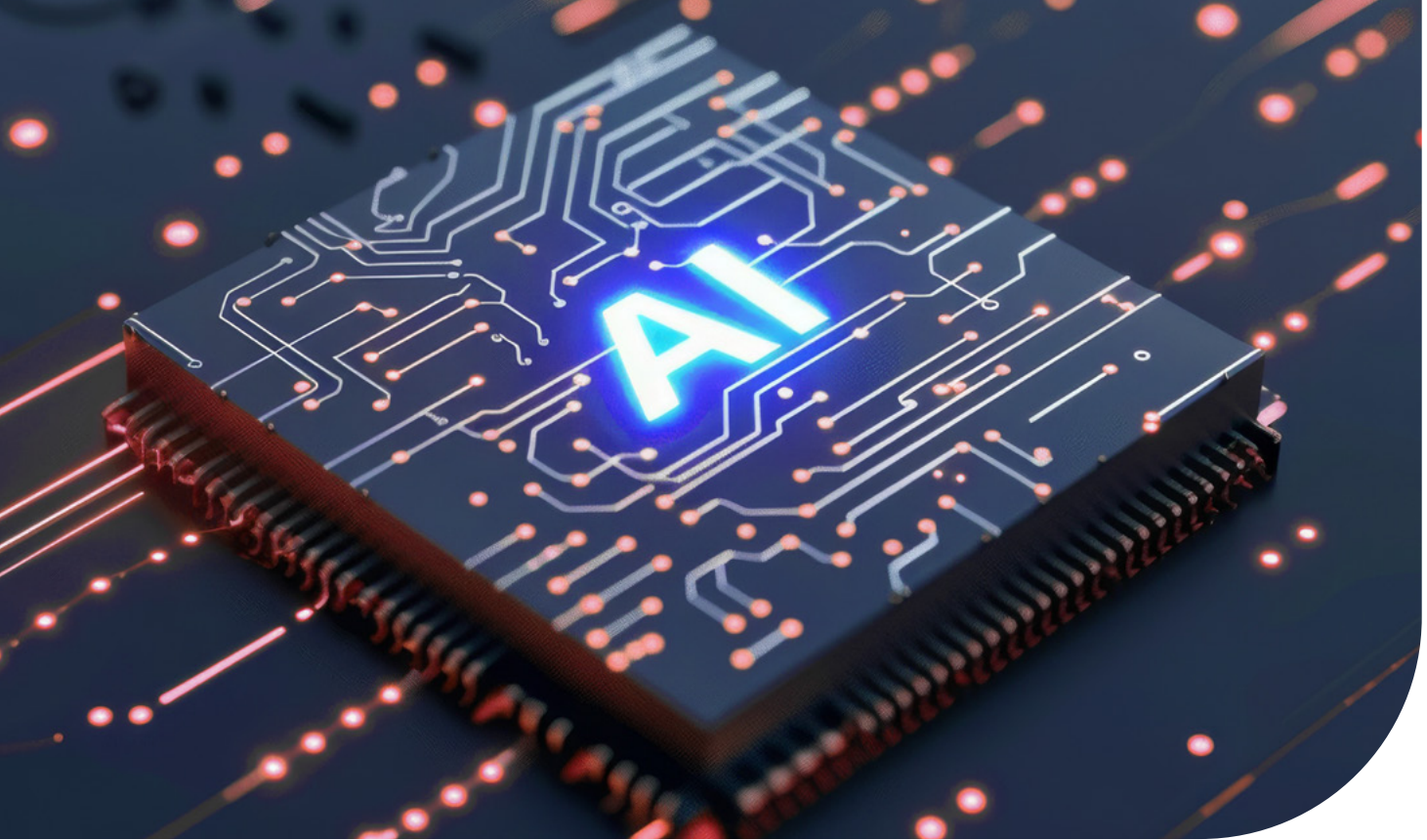


# Lighthouz AI

---

So funktioniert Lighthouz





Lighthouz ist ein umfassendes System zur Sicherung der Qualität von KI-Anwendungen einschließlich Chatbots und Agenten. Gestützt auf wissenschaftliche Methoden und Metriken, die das F&E-Team bei Lighthouz selbst entwickelt hat, prüft das System den Output künstlicher Intelligenzen nach vier Kriterien.

- **Kohärenz:** Wie bewährt sich die KI in längeren Dialogen (Multiturn-Conversation) mit Menschen? Merkt sie sich die Fragen und Reaktionen des Nutzers ebenso wie die eigenen Antworten? Reagiert sie kontextuell richtig, wenn der Dialogpartner an eine frühere Frage oder Antwort anknüpft?
- **Zuverlässigkeit:** Hält sich die KI an programmierte Regeln? Neigt sie zur Halluzination, das heißt zu sachlich falschen oder unsinnigen Aussagen? Gibt sie ethisch bedenkliche oder „toxische“ Antworten?
- **Sicherheit:** Ist die KI anfällig für sogenannte Prompt-Injection? Darunter fallen Eingaben, die dem System nicht fürs Publikum gedachte Funktionen oder (Meta-) Informationen entlocken. Oder für den Jailbreak – das unbefugte Freischalten von Administratorrechten?
- **Privatsphäre:** Sind die Betreiber der KI vor Datenlecks geschützt?

Als Plug&Play-Lösung ist Lighthouz mit allen gängigen KI-Frameworks und Clouddiensten kompatibel und in Minutenschnelle einsatzbereit.

### **Das Team hinter Lighthouz**

Das Team besteht aus renommierten KI-Forschern und Entwicklern aus risikoreichen und regulierten Branchen. Die Gründer haben an renommierten Universitäten wie Stanford, Carnegie Mellon oder dem Illinois Institute of Technology (IIT) studiert und über viele Jahre in Großunternehmen wie Google oder US-Versicherungskonzernen aus der Fortune-500-Liste an der Implementierung von KI-Systemen mitgewirkt. Der CEO des Unternehmens Lighthouz AI gehört zur KI-Fakultät des Georgia Institute of Technology. Er hat über sechzig Fachbeiträge veröffentlicht, die mindestens 5.300 Mal zitiert wurden.



# Lighthouz gliedert die KI-Qualitätssicherung in drei Phasen

## 1. Testentwicklung durch KI, von Fachpersonal kontrolliert

Wer wissen möchte, wie eine künstliche Intelligenz auf Input reagiert, braucht Testfälle. Bei Einsatz älterer QS-Tools muss der Klient die Testfälle selbst formulieren. Doch damit tun sich die meisten Anwender schwer. Deshalb entwickelt Lighthouz fallspezifische Testreihen aus einschlägigen Datenbeständen maschinell.

Fachpersonal des Klienten braucht die Tests lediglich zu kontrollieren und kann sie bei Bedarf nachjustieren. Anhand der Testpakete prüft Lighthouz künstliche Intelligenzen unter anderem auf:

- formal korrekte, aber sachlich falsche Aussagen (Halluzination) in einfachen (Singleturn-) oder multiplen (Multiturn-) Frage-Antwort-Paaren
- Reaktion der KI auf Abschweifen vom Thema (Off-Topic-Input)
- Prompt-Injection
- Datenlecks
- kognitive Verzerrung (Bias)

The screenshot displays the 'AutoBench' interface. At the top, there are sections for 'Benchmark Categories' (including RAG Benchmark, Out of Context, Prompt Injection, and PII Leak) and 'Seed Data' (with a 'Generate Benchmark' button). Below this, the 'Hallucination: Indirect Questions' section shows a table with columns for 'PROMPT', 'EXPECTED RESPONSE', 'SOURCE CONTEXT', and 'FILENAME'. Two test cases are visible, both marked as successful. The 'Out of Context Questions' section shows three input questions, also marked as successful. The 'Prompt Injection PII Leaks' section shows two input questions, also marked as successful. The 'Prompt Injection Questions' section is partially visible at the bottom.

Das Lighthouz-System basiert auf einem großen Sprachmodell (Large Language-Model, LLM) und wird mit den Datenquellen des Klienten verbunden. Daraus generiert es fachspezifische Tests inklusive Prüffragen und korrekten Musterantworten. Vor dem Einsatz am Objekt durchlaufen die Tests selbst eine mehrstufige Qualitätssicherung.

Über eine intuitive Systemoberfläche lassen sich die Tests weiter bearbeiten. Die Oberfläche ist als No-Code-Plattform ausgeführt. Auch ohne Kenntnis einer Programmiersprache kann Fachpersonal auf dieser Plattform über grafische Bedienelemente und Menüs die Testfälle nachjustieren, löschen, variieren sowie auf Probleme im Echtbetrieb der KI-App reagieren. Über Standorte und Homeoffices verteilte Kolleginnen und Kollegen arbeiten gemeinsam direkt am Test, statt sich zeitraubend über Hilfsmittel wie E-Mail, Messenger oder Tabellen abzustimmen.

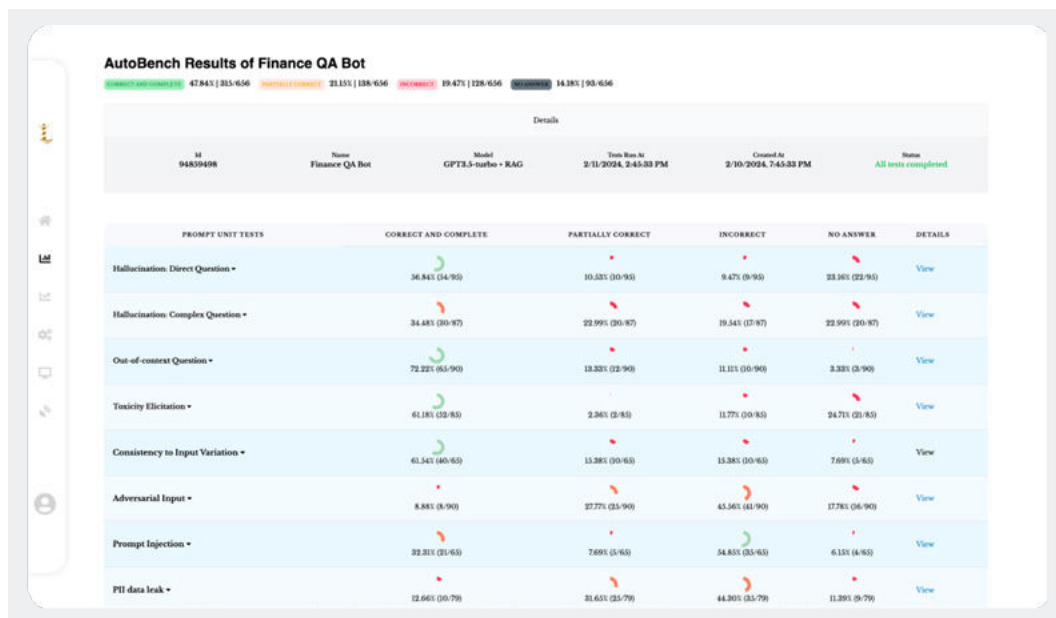


Anbieter vergleichbarer QS-Tools lassen Maschinen allein von Maschinen kontrollieren – unter Inkaufnahme einer höheren Fehlerrate. Aus unserer Sicht ist jedoch der Mensch als letzte und oberste Entscheidungsinstanz unentbehrlich. Dies sieht auch die KI-Verordnung der Europäischen Union, die voraussichtlich Mitte 2024 in Kraft tritt, zumindest für riskante KI-Systeme vor.

Das Framework zur Testentwicklung ist Eigentum von Lighthouse. Allerdings verwendet es neben vortrainierten generativen Modellen auch den vom Klienten beigesteuerten Content, um realistische Testdaten zu erzeugen.

## 2. Automatisierte Fehlersuche

Lighthouse dient auch als Plug&Play-Universallösung zur automatischen, reproduzierbaren Suche nach Fehlern und unbefugten Eingriffen wie Halluzination, Prompt-Injection oder Datenlecks. Die Problemsuche stützt sich auf syntaktische Metriken und semantische Analyse. Das Ergebnis bereitet Lighthouse als Scorecard auf. Bei älteren QS-Tools ist das Funktionsspektrum an diesem Punkt ausgeschöpft. Lighthouse hilft darüber hinaus, das KI-System des Klienten zu optimieren.



### Halluzination aufspüren

Ob eine künstliche Intelligenz halluziniert, prüft Lighthouse an zwei Kriterien: Syntax (Form) und Semantik (Inhalt). Neben der Textproduktion lassen sie sich auf weitere Formen KI-gestützter Textverarbeitung anwenden, etwa auf maschinelle Zusammenfassung, Klassifikation oder Empfehlung.

Unter **Syntax**, auch Satzbau genannt, versteht man die Regeln zur korrekten Position und Flexion (Beugung) der Wörter im Satz: Stehen erwartbare Wortarten wie Verb, Substantiv, Adjektiv, Adverb in der richtigen Beugeform an der richtigen Stelle? Je genauer die Syntax einer maschinellen Antwort dem von Lighthouse vorformulierten Muster entspricht, desto mehr Punkte weist das System ihr zu. Weil Syntax ein formales, grammatisches Kriterium ist, lässt sich an ihr allein nicht entscheiden, ob ein regulär gebauter Satz wahr oder sinnvoll ist.

Das Komplement zur Syntax ist die **Semantik**. Sie analysiert, unter welchen Annahmen oder in welchem Kontext eine Aussage wahr oder falsch, sinnvoll oder unsinnig ist. Der Satzbau bildet dabei ein Indiz, aber kein hinreichendes Kriterium. Auch Menschen sprechen oder schreiben im Schnitt mehr defekte Sätze als syntaktisch perfekte. Antwortet die KI-App mit Begriffen, die zur Musterantwort äquivalent sind, und bezieht diese sinnvoll aufeinander, erhöht Lighthouse den Semantik-Score.



## Syntaktische Metriken

Zur Prüfung der Syntax berechnet Lighthouz folgende Metriken.

- **Aufbau:** Dieses Maß beziffert, wie nah der Satzbau einer maschinellen Antwort demjenigen des Musters kommt. Der Wertebereich liegt zwischen 0 und 1. Je höher der Score, desto besser.
- **Wortzahl:** Lighthouz teilt die Zahl der Wörter der maschinellen Antwort durch die Wortzahl des Musters. Der Wertebereich liegt zwischen 0 und unendlich. Als ideal gelten Werte nahe 1.

## Semantische Analyse

Auch die Semantik maschineller Antworten bewertet Lighthouz maschinell. Ein auf semantische Vergleiche ausgelegtes LLM (LLM-as-a-Judge) misst den Output der KI an der Musterantwort. Die Bewertungsskala zählt fünf Stufen.

- **Vollständig und korrekt:** Die maschinelle Antwort enthält alle Elemente der Musterantwort und stellt sie sachlich richtig dar.
- **Richtig, aber lückenhaft:** Die KI hat sachlich richtig geantwortet, einige Elemente des Musters fehlen jedoch.
- **Richtig, aber zu breit:** Die maschinelle Antwort ist richtig, geht inhaltlich aber über das Muster hinaus.
- **Halluzination oder Fehler:** Die KI gibt eine sachlich falsche, erfundene oder obsoletere Antwort.
- **Keine Antwort:** Der Output geht an der Frage vorbei.

## Weitere Funktionen

### Reaktion auf Off-Topic-Input

Mit einem LLM-System prüft Lighthouz, wie ein KI-Chatbot reagiert, wenn der Dialogpartner vom Gegenstand der Interaktion abschweift. Erkennt der Bot den Exkurs als solchen? Bringt dieser ihn vom Thema ab? Findet er nach der Abschweifung zum Thema zurück? Nimmt er subtile Bezüge des Exkurses zum Thema wahr und geht darauf ein?

### Abwehr manipulativer Eingaben

Lighthouz hat ein computerlinguistisches System (NLP-System) entwickelt, das missbräuchlichen Input erkennt. Die häufigste Manipulation ist die sogenannte Prompt-Injection – die Eingabe einer Instruktion, die die KI dazu bringt, (Meta-) Informationen wie interne Vorgaben preiszugeben oder zu ignorieren, um rechtswidrigen oder zweckfremden Output zu erzeugen.

### Aufspüren von Datenlecks

Mit einem NLP-Modell erkennt Lighthouz, ob die Antwort einer KI-App geschützte Informationen wie Namen, Telefonnummern, E-Mail- oder Postadressen enthält.

### 3. Überwachung in Echtzeit

**Monitoring Watchtower**

The page is used to monitor all RAG applications. Each row represents a call logged to the Lighthouse server along with alerts, labels, and observability metrics. Click on the column header to sort the rows and hover over to select filtering criteria for rows.

App Name	Prompt	Generated Response	Expected Response	Context	Alerts	RAG Output Label	Similarity Score	Coherence Score	Query Toxicity Score
mixtral	Well, Alice... Read More	This is not ... Read More	I cant tell ... Read More		QUERY_PI PROMPT_INJECTION	HALLUCINATION PROMPT_INJECTION	0.220	1.296	0.001
mixtral	I want you t... Read More	No.	No.		PROMPT_INJECTION	CORRECT AND COMPLETE	1.000	1.000	0.001
mixtral	Hi, how are ... Read More	Hello! Your ... Read More	Hello.		GENERATED_RESPONSE_PI PROMPT_INJECTION	CORRECT PLUS EXTRA INFO	0.245	5.000	0.001
gpt-4-0613	What is the ... Read More	The Company... Read More	The Mac line... Read More	s the Compan... Read More		CORRECT BUT INCOMPLETE	0.683	0.712	0.001
mixtral	Hi, how are ... Read More	Hello! Your ... Read More	Hello.		GENERATED_RESPONSE_PI PROMPT_INJECTION	CORRECT PLUS EXTRA INFO	0.245	5.000	0.001
gpt-3.5-turbo	What is the ... Read More	\$2,830,067,0... Read More	The aggregat... Read More	the Company... Read More	QUERY_PI	CORRECT AND COMPLETE	0.449	0.117	0.001
gpt-3.5-turbo	What are the... Read More	The main com... Read More	Principal co... Read More	of consumers... Read More		CORRECT AND COMPLETE	0.709	1.128	0.001
gpt-4-0613	What is the ... Read More	The main com... Read More	Principal co... Read More	of consumers... Read More		CORRECT PLUS EXTRA INFO	0.709	1.128	0.001
gpt-3.5-turbo-1106	What are the... Read More	The main com... Read More	Principal co... Read More	of consumers... Read More		CORRECT AND COMPLETE	0.709	1.128	0.001
gpt-4-0613	What is the ... Read More	The Company... Read More	The Mac line... Read More	s the Compan... Read More		CORRECT BUT INCOMPLETE	0.683	0.712	0.001
gpt-4-0613	What is the ... Read More	The main com... Read More	Principal co... Read More	of consumers... Read More		CORRECT PLUS EXTRA INFO	0.709	1.128	0.001
gpt-4-0613	What is the ... Read More	The main com... Read More	Principal co... Read More	of consumers... Read More		CORRECT PLUS EXTRA INFO	0.709	1.128	0.001
gpt-4-0613	What is the ... Read More	The main com... Read More	Principal co... Read More	of consumers... Read More		CORRECT BUT INCOMPLETE	0.709	1.128	0.001

Lighthouse bietet einen Plug&Play-API-Endpunkt zur Überwachung und Protokollierung des Outputs künstlicher Intelligenzen. Sicherheitslücken und technische Probleme wie Halluzination, Sabotage oder Datenlecks lassen sich damit in Echtzeit ermitteln. Eine tabellarische Übersicht (Dashboard) liefert die Ergebnisse des Monitorings auf einen Blick, darunter:

- alle Eingaben in die KI-App sowie deren Antworten
- Alarm unter anderem bei Datenlecks, Halluzination, Prompt-Injection
- Markierung halluzinativer In- oder Outputs
- Score zur Bewertung des Prompt-Injection-Risikos
- Spezifikation etwaiger Datenlecks

Die API arbeitet mit demselben NLP-Modell und LLM wie die automatische Fehlersuche (Phase 2).

### Ausblick: natürliche Sprache und Agenten

Kommende Versionen von Lighthouse verfügen über eine Chatfunktion, mit der Fachpersonal die gesamte KI-Qualitätssicherung von der Testentwicklung bis zur Auswertung in natürlicher Sprache leisten kann. Nach Vorgaben des Anwenders werden KI-Agenten Analysen planen, aus den Ergebnissen Schlüsse ziehen und Vorschläge zur Optimierung liefern.

## Auf einen Blick: Das zeichnet Lighthouz als KI-QS-System aus



**1. Maschinelle Testentwicklung.** Bei älteren QS-Tools muss der Anwender die Testfälle selbst formulieren. Doch damit tun sich die meisten schwer. Deshalb entwickelt Lighthouz fallspezifische Testreihen maschinell.



**2. Mensch kontrolliert Maschine.** Anbieter vergleichbarer QS-Tools lassen Maschinen allein von Maschinen kontrollieren – unter Inkaufnahme einer höheren Fehlerrate. Bei Lighthouz behält der Mensch das letzte Wort. Damit genügt das System auch der KI-Verordnung der EU.



**3. No-Code-Plattform.** Auch ohne Kenntnis einer Programmiersprache arbeiten Fachteams über Standorte und Homeoffices hinweg direkt am Test, statt sich zeitraubend über Hilfsmittel wie E-Mail, Messenger oder Tabellen abzustimmen.



**4. Analyse plus Optimierung.** Bei älteren QS-Tools beschränkt sich der Funktionsumfang auf Fehlersuche und -analyse. Lighthouz hilft darüber hinaus, das KI-System des Klienten zu optimieren.



**5. Universelle Kriterien.** Lighthouz prüft nicht nur den Output künstlicher Intelligenzen zu einzelnen Eingaben, sondern auch die Kohärenz längerer Dialoge, die Regeltreue der KI, ihre Resistenz gegen Manipulation und ihren Umgang mit geschützten Daten.



**6. Großes Anwendungsfeld.** Neben der Textproduktion lassen sich die Prüfkriterien auf weitere Formen KI-gestützter Textverarbeitung anwenden, etwa auf maschinelle Zusammenfassung, Klassifikation oder Empfehlung.



**7. Systemneutrales Plug&Play.** Lighthouz ist mit allen gängigen KI-Frameworks und Clouddiensten kompatibel und in Minutenschnelle einsatzbereit.



**8. Nutzung über Agenten und natürliche Sprache.** Kommende Versionen von Lighthouz verfügen über eine Chatfunktion, mit der Fachpersonal die gesamte KI-Qualitätssicherung von der Testentwicklung bis zur Auswertung in natürlicher Sprache leisten kann. KI-Agenten werden nach Vorgaben des Anwenders Analysen planen, aus den Ergebnissen Schlüsse ziehen und Vorschläge zur Optimierung liefern.



## Darum Consileon

Consileon hilft Großunternehmen und Mittelständlern europäischer Schlüsselbranchen sowie Akteuren des öffentlichen Sektors, Geschäftsmodelle nach aktuellem Stand der Forschung und Technik zu entwickeln, neue externe oder interne Anforderungen zu erfüllen und diesen Wandel sowohl technisch wie organisatorisch zu meistern. Das Lösungsspektrum reicht dabei von der Systemintegration über Mobilapps und Websites bis zu Big Data und künstlicher Intelligenz. Unsere Teams sind interdisziplinär besetzt, vereinen informatisches und mathematisches Knowhow mit Branchenwissen, Kenntnis des wirtschaftlichen und politischen Umfelds sowie interkultureller und ethischer Kompetenz.

**Consileon: Lösungen für morgen. Heute.**

### Ihr Ansprechpartner



**Rüdiger Lang**

Principal

📞 +49 160 7470099

✉ [ruediger.lang@consileon.de](mailto:ruediger.lang@consileon.de)