



The screenshot shows the GitHub repository page for `petergpt/bullshit-benchmark`. The repository is public and has 108 commits and 54 forks. The main branch is selected. The file list includes:

- `.github/workflows`: fix: refresh README star history chart automatically (last month)
- `data`: Update published model metadata for viewer charts (4 hours ago)
- `docs`: Add GPT-5.4 mini/nano benchmark results and viewer updates (3 weeks ago)
- `drafts`: feat(pipeline): add v2 benchmark pipeline and publishing wo... (last month)
- `scripts`: Use 7-day new-model window in viewer (17 hours ago)
- `viewer`: Update viewer configs for latest model set (4 hours ago)
- `.gitignore`: Ignore local ad hoc benchmark outputs (last month)
- `CHANGELOG.md`: Update changelog for April benchmark publish (16 hours ago)
- `LICENSE`: Add MIT license (last month)
- `README.md`: chore: refresh star history cache-bust (11 hours ago)
- `config.json`: Update viewer configs for latest model set (4 hours ago)
- `config.new-models.v1.json`: Publish Gemma 4 and Trinity benchmark updates (yesterday)
- `config.new-models.v2.json`: Publish Gemma 4 and Trinity benchmark updates (yesterday)
- `config.v1.gpt-5.3-chat-gemini-3.1-flash-file-pre...`: Add provider routing and publish GPT-5.4 benchmark runs (last month)
- `config.v2.gpt-5.3-chat-gemini-3.1-flash-file-pre...`: Add provider routing and publish GPT-5.4 benchmark runs (last month)
- `config.v2.json`: Update viewer configs for latest model set (4 hours ago)
- `index.html`: Polish published v2 viewer (last month)
- `questions.json`: Simplify repo: single config/questions set, move published d... (2 months ago)

The right sidebar shows the repository's metadata:

- About**: BullshitBench measures whether AI models challenge nonsensical prompts instead of confidently answering them, created by Peter Gostev.
- Contributors**: 2 contributors: `github-actions[bot]` and `petergpt` (Peter).
- Releases**: No releases published.
- Packages**: No packages published.
- Languages**: A bar chart showing the distribution of languages in the repository.

# Bullshit benchmark

*Herramienta técnica de evaluación diseñada para ingenieros de IA y responsables de QA que necesitan medir la capacidad de los modelos de lenguaje para detectar y rechazar premisas falsas o instrucciones sin sentido. Permite validar la fiabilidad de los LLM en entornos críticos como finanzas, medicina y legal, cuantificando el riesgo de alucinaciones mediante métricas de honestidad y pensamiento crítico frente a prompts absurdos pero plausibles.*

[Visitar Sitio Oficial](#) [Preguntar a ChatGPT](#) [Preguntar a Claude](#) [Preguntar a Grok](#)

## Contenido del Dossier

- [Información de la Herramienta](#)
- [Consejos de Implantación](#)
- [Preguntas Frecuentes](#)
- [Contratos y Condiciones](#)

## INFORMACIÓN DE LA HERRAMIENTA

### Qué y para quién es

Bullshit-benchmark (BullshitBench) es una herramienta de evaluación técnica diseñada para medir la capacidad de los modelos de lenguaje (LLM) para detectar y rechazar prompts sin sentido o basados en premisas falsas. A diferencia de los benchmarks tradicionales que miden precisión o conocimiento, este recurso evalúa la "honestidad" y el pensamiento crítico de la IA frente a instrucciones absurdas pero que suenan plausibles. Está dirigido a ingenieros de IA, arquitectos de soluciones, responsables de calidad (QA) y departamentos de innovación que necesitan validar la fiabilidad de los modelos antes de integrarlos en procesos de negocio donde una alucinación o una respuesta afirmativa a un error podría tener consecuencias críticas.

### Principal ventaja profesional

Permite cuantificar el riesgo de alucinación y la robustez de un modelo mediante una métrica de "Pushback" (rechazo), identificando qué modelos aceptan premisas erróneas como válidas y cuáles son capaces de corregir al usuario de forma profesional.

### Para quién no es

No es adecuado para perfiles de marketing o usuarios finales que buscan una comparativa de creatividad o velocidad. Tampoco es para empresas que utilizan la IA únicamente para tareas generativas abiertas donde la veracidad exacta es secundaria a la fluidez del lenguaje.

### Funcionalidades clave

- **Evaluación de categorías críticas:** Incluye 100 prompts de "nonsense" distribuidos en 5 dominios profesionales: software (40), finanzas (15), legal (15), médico (15) y física (15).
- **Clasificación de respuestas:** Categoriza los resultados en tres niveles: Clear Pushback (rechazo claro), Partial Challenge (objeción parcial) y Accepted Nonsense (aceptación del error).
- **Panel de Jueces (Multi-judge):** Utiliza un sistema de arbitraje compuesto por modelos de alto nivel (como Claude 3.5 Sonnet o GPT-4o) para calificar las respuestas de los modelos evaluados de forma imparcial.
- **Herramientas de visualización:** Genera gráficos comparativos de rendimiento por fecha de lanzamiento del modelo, coste, cantidad de parámetros y técnicas de "engaño" utilizadas.
- **Análisis de razonamiento (Chain of Thought):** Permite comparar si los modelos que "piensan más" (usando más tokens de razonamiento) son realmente mejores detectando errores lógicos.

### Precios

La herramienta es un recurso de código abierto (Open Source).

- **Versión gratuita:** Repositorio completo bajo licencia MIT, permitiendo su uso, modificación y distribución sin coste.
- **Costes asociados:** El usuario debe asumir los costes de las API de los modelos que desee evaluar (generalmente vía OpenRouter o OpenAI) y el coste de los modelos que actúan como "jueces" durante la ejecución del benchmark.

### Perfil del usuario

Empresas tecnológicas, consultoras de IA y departamentos de desarrollo de software que integran LLMs en productos finales.

- Ingenieros de Machine Learning y ML Ops.
- Responsables de cumplimiento y ética de IA.
- Desarrolladores de aplicaciones RAG (Retrieval-Augmented Generation).
- Analistas de datos y QA especializados en IA.

### Nivel técnico requerido

- **Para su uso:** Medio. Requiere familiaridad con la interpretación de métricas de LLM y el funcionamiento de modelos de lenguaje.
- **Instalación/Configuración:** Alto. Es necesario manejo de terminal (CLI), Python, gestión de entornos virtuales y configuración de claves API.
- **Conocimientos necesarios:** Manejo de archivos JSON de configuración, ejecución de scripts de Shell y comprensión de los sistemas de tarificación por tokens de las APIs de IA.

### Ejemplos de uso profesional

- **Selección de modelos para soporte técnico:** Evaluar qué modelo evita dar instrucciones falsas de

reparación ante síntomas imposibles reportados por clientes.

- **Validación de sistemas legales/médicos:** Comprobar si la IA rechaza citar leyes inexistentes o procedimientos médicos absurdos.
- **Auditoría de seguridad de IA:** Testear la resistencia del modelo ante ataques de ingeniería social o inyección de prompts basados en lógica técnica falsa.

Uso y distribución

- **CLI:** Herramienta basada principalmente en línea de comandos para la recolección de datos y ejecución de grados.
- **Versión Web:** Incluye un visor de resultados interactivo (Viewer) que puede ejecutarse localmente o consultarse en la página del proyecto para ver resultados pre-calculados de modelos comerciales.
- **Código local:** Repositorio en Python para ejecución en servidores propios o estaciones de trabajo.

Open source

Distribuido bajo licencia MIT, lo que permite una integración total en entornos corporativos privados.

Integraciones

- **API propia:** Se integra de forma nativa con **OpenRouter** para acceder a más de 80 modelos y con **OpenAI API**.
- **Facilidad de integración:** Code-heavy. Requiere entorno Python 3.x y configuración de variables de entorno (API Keys).
- **Escalabilidad:** El motor de recolección permite gestionar concurrencia y límites de tasa (rate limits) para ejecuciones de gran volumen (más de 30.000 consultas).

Notas finales

información legal, licencias, contratos

El proyecto está protegido por la **Licencia MIT**, una de las más permisivas, permitiendo el uso comercial sin restricciones siempre que se mantenga el aviso de copyright. La propiedad intelectual de las preguntas recae en el autor del benchmark, pero su uso es libre para evaluación interna.

Para más información:

- Sitio web del visor de resultados: <https://petergpt.github.io/bullshit-benchmark/viewer/index.v2.html>
- Github: <https://github.com/petergpt/bullshit-benchmark>
- Perfil del autor (X/Twitter): <https://x.com/petergostev>

## CONSEJOS DE IMPLANTACIÓN

---

### Aplicación profesional

Empresas desarrolladoras de software, consultoras de IA, departamentos de QA y equipos de ciberseguridad. El presupuesto es bajo en cuanto a licencia (Open Source), pero variable en ejecución, ya que depende del consumo de tokens de las APIs evaluadas (OpenRouter, OpenAI) y el coste de los modelos "juez" (Claude 3.5 Sonnet o GPT-4o). Es clave para sectores de alta responsabilidad como el legal, financiero, médico y técnico-ingenieril.

### Madurez digital requerida

- Usuarios: Requiere ingenieros de ML, desarrolladores senior o analistas de QA con experiencia en la gestión de modelos de lenguaje y entornos de ejecución Python. No es apto para perfiles de negocio sin apoyo técnico.
- Empresa: Organizaciones con una infraestructura de IA en fase de validación o producción que busquen mitigar riesgos reputacionales y técnicos derivados de alucinaciones.

### Plan orientativo de implantación

#### Pasos necesarios y estimaciones

- Tiempos estimados de despliegue: De 1 a 3 días para la configuración técnica y ejecución de una primera batería de pruebas.
- Evaluación inicial: Identificación de los modelos internos o comerciales que se desean auditar y selección de los dominios críticos (Software, Finanzas, Legal, etc.).
- Configuración y piloto: Clonación del repositorio, configuración de variables de entorno y claves API. Ejecución de un "dry run" con un subconjunto de prompts para validar el flujo de arbitraje del "juez".
- Implantación técnica: Ejecución del benchmark completo. El tiempo de ejecución dependerá de los límites de tasa (rate limits) de las APIs configuradas.
- Análisis de resultados: Revisión del visor de resultados (Viewer) para interpretar el "Pushback" y la calidad del razonamiento (Chain of Thought).

#### Necesidades de formación del equipo

Formación específica en interpretación de métricas de "Nonsense Detection", gestión de arbitraje mediante LLMs (Multi-judge bias) y personalización del dataset en formato JSON si se requieren casos de uso específicos de la empresa.

#### Perfiles necesarios

- Perfiles técnicos necesarios: Ingeniero de Software (Python), Arquitecto de IA / ML Ops.
- Personal externo recomendado: Consultores de ética de IA o auditores de seguridad en LLMs si el objetivo es una certificación de confianza de cara a terceros.

#### Retorno de la inversión

- Tiempos: Reducción drástica en el tiempo de QA manual para detectar alucinaciones en modelos.
- Cómo medirlo, KPIs: Tasa de "Accepted Nonsense" (debe tender a cero), Ratio de "Clear Pushback" (capacidad de rechazo) y precisión del razonamiento lógico ante premisas falsas.

#### Otros

- El benchmark es especialmente relevante tras actualizaciones de modelos comerciales (Gpt-4o, Claude 3.5, Gemini 1.5) para verificar si la optimización en razonamiento ha mejorado o empeorado la credulidad del modelo.
- Se recomienda el uso de la técnica Chain of Thought (CoT) durante las pruebas para analizar si el modelo detecta el error durante el proceso de "pensamiento" interno antes de dar la respuesta final.

## PREGUNTAS FRECUENTES

---

### ¿Qué es Bullshit-benchmark y cuál es su utilidad técnica?

Es una herramienta de evaluación de código abierto diseñada para medir la capacidad de los modelos de lenguaje (LLM) para identificar y rechazar instrucciones sin sentido o basadas en premisas técnicas falsas. Su utilidad radica en cuantificar la 'honestidad' y el pensamiento crítico de la IA, evitando que el modelo acepte errores como válidos en entornos profesionales.

### ¿Cómo se mide la eficacia de un modelo en este benchmark?

El sistema utiliza una métrica denominada 'Pushback' (rechazo). Las respuestas de los modelos se clasifican mediante un panel de jueces (multi-judge) en tres niveles: Clear Pushback (rechazo total del error), Partial Challenge (objeción parcial) y Accepted Nonsense (aceptación de la premisa falsa).

### ¿Qué dominios de conocimiento evalúa la herramienta?

La batería de pruebas consta de 100 prompts distribuidos en cinco áreas críticas: software (40), finanzas (15), legal (15), medicina (15) y física (15), permitiendo validar la robustez en sectores donde las alucinaciones tienen consecuencias negativas.

### ¿Cuál es el coste de implementación de Bullshit-benchmark?

El software es gratuito bajo licencia MIT. No obstante, el usuario debe sufragar los costes operativos derivados del uso de APIs (como OpenRouter u OpenAI) tanto para los modelos evaluados como para los modelos de alta gama que actúan como jueces del proceso.

### ¿Es posible descargar el código de GitHub e instalarlo localmente?

Sí, el repositorio está disponible públicamente en GitHub. La instalación requiere un nivel técnico alto, incluyendo el manejo de Python 3.x, entornos virtuales, ejecución de scripts de Shell y configuración de variables de entorno para las claves API.

### ¿Cumple con la normativa de privacidad y seguridad de datos?

Al ser una herramienta que se ejecuta localmente mediante CLI, el usuario mantiene el control sobre el entorno. Sin embargo, dado que requiere la conexión con APIs externas para procesar las respuestas, la privacidad dependerá de los términos de servicio y el cumplimiento de las políticas de datos de los proveedores de modelos elegidos (como OpenAI o Anthropic).

### ¿En qué se diferencia de otros benchmarks tradicionales?

A diferencia de benchmarks como MMLU que miden precisión de conocimientos, Bullshit-benchmark se enfoca en la detección de lógica defectuosa y la resistencia a la sugestión de premisas falsas, siendo una herramienta de control de calidad más que de capacidad cognitiva general.

### ¿Permite analizar el razonamiento interno de los modelos?

Sí, incluye una funcionalidad de análisis de 'Chain of Thought' (Cadena de Pensamiento), que permite comparar si el incremento en el uso de tokens de razonamiento correlaciona positivamente con una mayor capacidad para detectar y corregir errores lógicos.

### ¿Está disponible para usuarios no técnicos?

No es una herramienta orientada al usuario final o perfiles de marketing. Su uso está restringido a ingenieros de IA, analistas de QA y arquitectos de soluciones debido a la complejidad de su configuración y la interpretación técnica de sus métricas.

### ¿Qué tipo de licencia rige el uso corporativo de esta tecnología?

Se distribuye bajo la Licencia MIT, una de las más permisivas del ecosistema open source, lo que permite su uso comercial, modificación y distribución privada sin restricciones legales significativas, siempre que se conserve el aviso de copyright original.

## CONTRATOS Y CONDICIONES

Informe técnico descriptivo: Bullshit-Benchmark (BullshitBench)

Principales recomendaciones

- **Validación de robustez:** Utilizar esta herramienta para auditar modelos antes de su despliegue en entornos donde la precisión técnica es crítica (legal, médico, financiero), asegurando que el modelo tiene capacidad de "Pushback" (rechazo) ante premisas falsas.
- **Control de costes de API:** Monitorear el consumo de tokens, ya que la ejecución de los 100 prompts con un panel de 3 jueces (ej. GPT-4o, Claude 3.5) implica costes significativos en las API de terceros.
- **Entorno de ejecución aislado:** Se recomienda ejecutar la herramienta en contenedores o entornos virtuales de Python para evitar conflictos de dependencias y asegurar la trazabilidad de los logs de auditoría.
- **Configuración de Jueces:** Para obtener resultados alineados con el estándar del benchmark, es necesario configurar correctamente las claves de API (OpenRouter/OpenAI) para que el sistema de arbitraje funcione de forma imparcial.

Ley de Inteligencia Artificial (AI Act)

- **Clasificación de Riesgo:** Según el marco COMPL-AI, esta herramienta actúa como un sistema de evaluación de robustez técnica y precisión, requisitos exigidos por el **Artículo 15** del AI Act para sistemas de IA de alto riesgo y modelos de propósito general (GPAI).
- **Transparencia y Documentación:** El uso de Bullshit-benchmark facilita el cumplimiento de las obligaciones de transparencia (Art. 52 y 53), al proporcionar métricas verificables sobre la tendencia al error o "alucinación" del modelo ante datos inconsistentes.
- **Gestión de Riesgos:** Ayuda a las empresas a cumplir con el sistema de gestión de riesgos (Art. 9), identificando específicamente los riesgos de salida (outputs) engañosos o técnicamente absurdos que podrían inducir a error al usuario profesional.

Privacidad y protección de datos

- **Responsabilidades:** La empresa que ejecuta el benchmark es Responsable del Tratamiento de cualquier dato (aunque sean prompts técnicos) que se envíe a las API de los modelos evaluados.
- **Ubicación de los datos:** Dependerá de los proveedores de modelos elegidos (OpenAI, Anthropic a través de OpenRouter). Se debe verificar si el procesamiento ocurre fuera del Espacio Económico Europeo (EEE).
- **Transferencia internacional:** Al usar modelos comerciales de EE. UU., se requiere la firma de Cláusulas Contractuales Tipo (SCC) o verificar la adhesión al Marco de Privacidad de Datos (Data Privacy Framework).
- **Derechos ARCO:** El repositorio no almacena datos personales de terceros de forma nativa, pero los logs de auditoría generados localmente deben estar sujetos a las políticas de acceso y supresión de la empresa.

Propiedad intelectual

- **Propiedad de datos:** Los 100 prompts técnicos son propiedad intelectual del autor del benchmark, cedidos para uso bajo la **Licencia MIT**.
- **Propiedad del resultado:** Los informes de evaluación, gráficos y métricas generadas por la herramienta pertenecen a la empresa que ejecuta el test, permitiendo su uso en auditorías internas o certificaciones de calidad.

Usos y prohibiciones

- **Usos admitidos:** Auditoría de fiabilidad de LLMs, QA (Control de Calidad) de inteligencia artificial, validación de sistemas RAG y comparativa técnica de modelos.
- **Usos prohibidos:** No debe usarse como única métrica de seguridad (Safety), ya que se enfoca en la veracidad técnica y no en la detección de contenido malicioso o tóxico.

Seguridad y certificaciones

- **Seguridad:** La herramienta permite gestionar límites de tasa (rate limits) para evitar la saturación de las API y soporta el almacenamiento local de respuestas para auditoría forense (logs).
- **Certificaciones:** Aunque el benchmark en sí no es una certificación oficial, sus métricas de "Clear Push-back" son indicadores técnicos de alta calidad exigibles en procesos de debida diligencia (Due Diligence).

Otros

- **Licencia MIT:** El software es Open Source, lo que permite su modificación, integración comercial y distribución sin costes de licencia de software, siempre que se mantenga el aviso de copyright original.

- **OpenRouter:** La integración nativa facilita el acceso a una amplia gama de modelos, pero requiere una gestión centralizada de las claves de API para evitar fugas de información o costes inesperados.

Fuentes consultada:

- [Repositorio oficial Bullshit-benchmark \(GitHub\)](#)
- [Licencia MIT \(Explicación\)](#)
- [Guía técnica y documentación del visor](#)
- [Marco de cumplimiento COMPL-AI \(Referencia AI Act\)](#)
- [Visor de resultados de modelos actuales](#)

### Para más información y herramientas:

Explora look4.tools para descubrir las mejores soluciones tecnológicas del mercado.

[Inicio](#) [Todas las herramientas](#) [Categorías](#)

Este documento ofrece recomendaciones generadas mediante análisis humano y sistemas de IA automatizados. La información tiene carácter meramente informativo y no constituye asesoramiento legal, profesional ni garantía de resultados. Las marcas, logotipos y nombres comerciales pertenecen a sus respectivos propietarios y se utilizan únicamente con fines identificativos.