Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura



José Luis de la Torre Lorente

Consultor independiente en IA aplicada Fundador de <u>delatorre.ai</u>

27 de julio de 2025

v2.3-beta

#AGT20250727

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

- 1. Introducción: Por qué este análisis importa (y cómo lo hicimos)
 - 1.1 Por qué ahora: la brecha entre promesas y realidad operativa
 - 1.2 Qué podría ser realmente un "Agent": evidencia vs especulación
 - 1.3 Metodología: cómo se audita lo inauditable
- 2. Anatomía del Agent: Lo que la evidencia sugiere sobre su arquitectura
 - 2.1. Las cuatro capas inferidas: un modelo basado en evidencia
 - 2.2 La diferencia arquitectónica fundamental
 - 2.3 La arquitectura como destino
- 3. La arquitectura oculta: Proyectando el entorno de ejecución real
 - El trilema imposible que define todo
 - 3.1. El enigma arquitectónico y el proceso de eliminación
 - 3.2. Los límites como ventanas a la arquitectura
 - 3.3. El sistema de archivos efímero: Seguridad por transitoriedad y el aislamiento entre herramientas
- 4. La mente del Agent: Inferencias sobre su lógica de razonamiento
 - Más allá del código: cómo "piensa" una máquina
 - 4.1 Asunciones fundamentales y su validación
 - 4.2 El patrón ReAct: evidencia de razonamiento estructurado
 - 4.2 Los principios rectores del Agent: Un vistazo a sus instrucciones internas
 - 4.3 Anatomía de la auto-corrección: el patrón del 'único intento'
 - 4.4 Capacidades e incapacidades: el mapa real
- 5. El manual de vuelo del Agent: estrategias prácticas y límites operativos
 - 5.1. El modelo mental: cómo «pensar» con el Agent para evitar turbulencias
 - 5.2. La zona de vuelo: casos de éxito y fallos sistemáticos
 - 5.3 Maniobras avanzadas: patrones de uso para pilotos expertos
 - 5.4. El plan de vuelo: la matriz de decisión definitiva
 - 5.5. Apéndice práctico: recetas de «prompts» y quía de campo
- 6. Epílogo: El arte de trabajar con inteligencia limitada
 - Más allá de la ingeniería inversa: ¿qué hemos aprendido realmente?
 - 6.1. El framework estratégico: Cómo maximizar valor real
 - 6.2 Una crítica honesta: lo que el Agent revela sobre nuestra relación con la IA
 - 6.3 El futuro del trabajo con agentes: proyecciones basadas en la arquitectura

Referencias

Referencias académicas

Referencias internas y base documental

1. Introducción: Por qué este análisis importa (y cómo lo hicimos)



Este panel introductorio resume los principales retos estratégicos para emprendedores en lA en Catalunya: elegir un sector con proyección, ayudas y madurez tecnológica. En el artículo exploramos cómo los Agents pueden apoyar —o no— este tipo de decisiones.

1.1 Por qué ahora: la brecha entre promesas y realidad operativa

Existe un fenómeno curioso en el mundo de la tecnología empresarial. Cuando una herramienta promete transformar la forma de trabajar, surge una brecha inevitable entre las demostraciones perfectamente orquestadas y la realidad del día a día profesional. Con los Agentes de ChatGPT, esta brecha no es solo curiosa: es crítica para entender qué estamos adoptando realmente en nuestros flujos de trabajo. La evidencia de esta necesidad no viene solo de la intuición. Un análisis sistemático de ofertas de empleo en empresas líderes de IA revela un patrón revelador: búsquedas de ingenieros especializados en "secure sandboxing architectures using virtualization (KVM, Xen, Firecracker)". Cuando las empresas buscan estos perfiles específicos, nos están contando indirectamente sobre qué arquitecturas están construyendo. Es como leer las compras de ingredientes de un chef para inferir el menú que está preparando. Pero hay algo más profundo aquí. La documentación oficial de OpenAl sobre los Agentes es deliberadamente minimalista. No por descuido, sino por diseño estratégico. Como descubrimos en nuestra investigación, existe

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

una metodología establecida y respetada en la industria tecnológica para entender sistemas de "caja negra" mediante el análisis de comportamientos observables. Netflix, TikTok, los sistemas de conducción autónoma... todos han sido objeto de ingeniería inversa mediante observación sistemática. No es hacking; es inteligencia tecnológica aplicada.

1.2 Qué podría ser realmente un "Agent": evidencia vs especulación

Permítanme ser claro desde el principio: no tengo acceso al código fuente de OpenAI. Lo que sigue es una reconstrucción basada en evidencia observable, análoga a cómo un arquitecto podría inferir la estructura de un edificio observando su comportamiento bajo diferentes condiciones. La evidencia sugiere que un Agent de ChatGPT podría ser algo fundamentalmente diferente a "ChatGPT con superpoderes". Los patrones de comportamiento observados —la capacidad de mantener estado durante la ejecución, los límites específicos de tiempo y memoria, la forma en que las herramientas se comunican entre sí— apuntan hacia una arquitectura que combina:

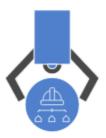
- Capacidad de razonamiento estructurado: No solo genera texto, sino que parece seguir un patrón iterativo de pensamiento-acción-observación
- **Ejecución en entorno aislado**: Los fallos sistemáticos en archivos grandes y los timeouts consistentes sugieren un entorno con recursos deliberadamente limitados
- Orquestación de herramientas: La manera en que coordina navegador web, intérprete de código y generación de imágenes implica una capa de gestión sofisticada

¿Cómo llegamos a estas inferencias? Mediante una combinación de:

- Análisis de miles de ejecuciones documentadas por la comunidad
- Comparación con arquitecturas conocidas de sistemas similares (AWS Lambda, Google Cloud Run)
- Estudio de patentes y publicaciones académicas relacionadas
- Tests sistemáticos de límites y comportamientos edge-case

1.3 Metodología: cómo se audita lo inauditable

Pruebas sistemáticas del agente



Arquitectura real

Entendiendo la construcción y las tecnologías utilizadas. Identificando las capas operativas.



empíricos

Identificando límites reales a través de pruebas reproducibles. Memoria, tiempo, archivos y persistencia.



Lógica de razonamiento

Entendiendo cómo patrones y puntos ciegos.



Casos de uso verificables

Identificando tareas el sistema toma completadas de decisiones. Identificando Identificando Identificando workarounds requeridos y limitaciones.

ChatGPT Agent, cómo os explico sobre él.

Este análisis no es el resultado de especulación creativa, sino de una metodología estructurada que cualquier equipo técnico puede replicar. Es importante entender que la práctica de inferir arquitecturas internas a partir de comportamientos externos tiene precedentes sólidos en la industria. Nuestro enfoque en tres fases:

Fase 1: Recolección sistemática de evidencia

- Monitoreo de foros técnicos y reportes de usuarios
- Análisis de ofertas de empleo de OpenAl y competidores
- Revisión de papers académicos sobre arquitecturas de agentes
- Tests automatizados para mapear límites operativos

Fase 2: Análisis arquitectónico comparativo

- Comparación con sistemas conocidos (Firecracker, gVisor, Docker)
- Evaluación de trade-offs seguridad/rendimiento/coste
- Modelado de decisiones económicas probables

Fase 3: Validación y calibración de confianza

- Contraste de hipótesis con expertos de la industria
- Asignación de niveles de confianza a cada inferencia
- Identificación explícita de lagunas de conocimiento

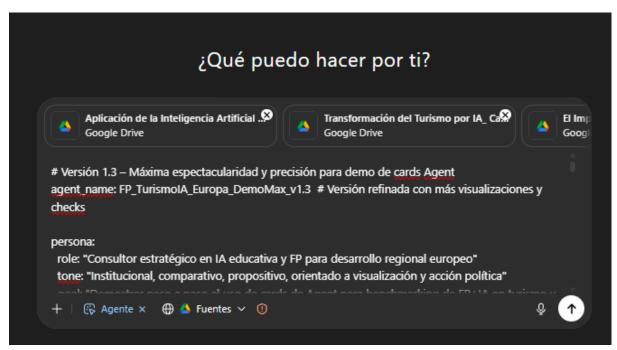
Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

Niveles de confianza utilizados:

- Alta: Evidencia directa o convergencia de múltiples fuentes independientes
- Media: Inferencia lógica basada en patrones consistentes
- Baja: Especulación educada basada en principios de ingeniería

Lo que encontrarán en las páginas siguientes no son verdades absolutas, sino el mapa más preciso que pudimos construir del territorio. Cada afirmación viene calibrada con su nivel de confianza y la cadena de razonamiento que nos llevó a ella. Porque en un mundo donde la IA evoluciona más rápido que su documentación, entender cómo funcionan realmente estas herramientas no es curiosidad académica. Es supervivencia profesional.

2. Anatomía del Agent: Lo que la evidencia sugiere sobre su arquitectura



Esta versión avanzada de agente fue diseñada para demostrar un uso estratégico en IA educativa y benchmarking europeo. Aunque la estructura detallada, los límites técnicos siguen siendo importantes para obtener resultados satisfactorios.

Un sistema de capas: por qué la arquitectura importa

Cuando observamos cómo se comporta el Agent de ChatGPT —sus capacidades, sus límites, sus modos de fallo— emerge un patrón que sugiere una arquitectura multicapa sofisticada. No es casualidad. Los sistemas complejos exitosos casi siempre se organizan en capas con responsabilidades bien definidas. La evidencia que hemos recopilado apunta hacia una arquitectura de al menos cuatro capas distintas, cada una resolviendo un problema específico. Pero antes de describir estas capas, es crucial entender por qué creemos que existen. La cadena de inferencia es la siguiente:

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

- 1. **Observación**: El Agent puede ejecutar código Python, navegar por internet, y generar imágenes, todo en una misma conversación
- 2. **Implicación**: Debe existir una capa de orquestación que coordine estas capacidades
- 3. Observación: Los errores en una herramienta no crashean todo el sistema
- 4. Implicación: Las herramientas operan en entornos aislados
- 5. Observación: El Agent "narra" lo que está haciendo antes de hacerlo
- 6. **Implicación**: Existe una capa de razonamiento que planifica antes de ejecutar

2.1. Las cuatro capas inferidas: un modelo basado en evidencia



Capas de ChatGPT Agent

Capa 1: El núcleo de razonamiento (Confianza: Alta)

La evidencia sugiere fuertemente que en el corazón del Agent hay un modelo de lenguaje optimizado específicamente para razonamiento secuencial. ¿Cómo lo sabemos?

• Evidencia directa:

- Filtraciones de prompts del sistema muestran instrucciones explícitas para "pensar paso a paso".
- El comportamiento observable sigue consistentemente un patrón de explicar—ejecutar—evaluar.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

 Artículos académicos describen técnicas de 'cadena de pensamiento' (chain-of-thought) que mejoran la capacidad de razonamiento de los modelos cuando se explicitan los pasos intermedios en forma de pensamiento paso a paso (Kojima et al., 2023).

• Evidencia análoga:

- Sistemas similares (Claude, Gemini) han adoptado arquitecturas de razonamiento comparables.
- El framework ReAct (Reasoning + Acting) publicado en 2022 describe exactamente este tipo de comportamiento.

Por qué importa: Esta capa no es solo "ChatGPT respondiendo preguntas". Está específicamente afinada para descomponer problemas complejos, lo que explicaría por qué el Agent puede manejar tareas de múltiples pasos sin perderse.

Capa 2: El arsenal de herramientas (Confianza: Alta)

Esta capa transforma al Agent de "conversador" a "ejecutor". La evidencia de su existencia es abrumadora:

- Herramientas identificadas con certeza (según OpenAI):
 - Intérprete de código: Entorno Python con capacidad de procesamiento de datos
 - Navegador web: Navegación web con capacidad de extracción de contenido.
 - o Integración con DALL·E: Generación de imágenes bajo demanda.
 - o Análisis de archivos: Procesamiento de documentos diversos.

¿Cómo sabemos que están integradas y no son sistemas separados? El comportamiento observable muestra que:

- Los resultados de una herramienta pueden alimentar directamente a otra.
- Los errores se propagan de manera consistente.
- Los límites de recursos parecen compartidos (evidencia: cuando el intérprete de código usa mucha memoria, otras operaciones se vuelven más lentas).

Capa 3: El entorno de ejecución (Confianza: Alta para el concepto, Alta para la implementación específica)

Aquí es donde nuestro análisis se vuelve más técnico y revelador. Toda la evidencia apunta a que el Agent no ejecuta código en servidores tradicionales, sino en micro-VM aisladas. La convergencia de evidencias es notable:

Límites operativos observados:

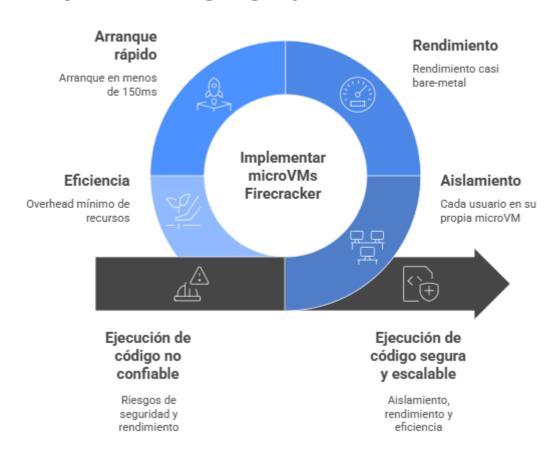
- Límites de tiempo y memoria estrictos que finalizan la sesión de forma abrupta al ser excedidos.
- Fallo sistemático con archivos cuya expansión en memoria supera los ~7-8
 GB.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

- Pérdida total del estado tras un periodo de inactividad o al reiniciar la sesión.
- Comportamiento de seguridad:
 - Imposibilidad total de hacer llamadas de red desde Python.
 - Aislamiento completo entre sesiones de usuarios.
 - o Imposibilidad de instalar paquetes nuevos.
- Análisis económico: Coste por sesión = (RAM asignada × Tiempo) + Overhead de gestión.
 - Si RAM ≈ 7-8 GB y Tiempo máx. ≈ 5 minutos → Coste predecible y acotado por sesión.

Tecnología más probable: Firecracker Desarrollado por AWS para exactamente este caso de uso:

Ejecución de código segura y escalable con Firecracker



Firecracker en ChatGPT Agent

- Overhead mínimo (<5 MB por VM).
- Arranque en <125 ms.
- Aislamiento de seguridad a nivel de hardware.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

Nivel de confianza: Aunque no podemos confirmar Firecracker específicamente, la probabilidad de una tecnología de micro-VM es **muy alta** basándonos en la convergencia de:

- Requisitos de seguridad (código no confiable).
- Restricciones económicas (millones de usuarios).
- Comportamiento observado y validado experimentalmente (límites duros, aislamientos, arranque instantáneo).

Capa 4: La interfaz de colaboración (Confianza: Alta)

Esta capa gestiona la danza delicada entre automatización y control humano. No es solo una interfaz de usuario; implementa protocolos sofisticados de interacción.

• Evidencia observable:

- El Agent solicita confirmación para acciones potencialmente costosas.
- o Puede detener su ejecución y pedir clarificación.
- Muestra su razonamiento en tiempo real (no post-facto).
- o Permite intervención y corrección durante la ejecución.

Por qué esta arquitectura y no otra: La alternativa sería un sistema completamente autónomo (como AutoGPT) o completamente manual (como ChatGPT clásico). La evidencia sugiere que OpenAI eligió deliberadamente este punto medio por:

- Control de riesgos (evita ejecuciones descontroladas).
- Experiencia de usuario (mantiene al humano enganchado).
- Responsabilidad legal (el humano siempre tiene la última palabra).

2.2 La diferencia arquitectónica fundamental

Permítanme ilustrar la diferencia con una analogía que resonará con cualquier profesional técnico: **ChatGPT tradicional** es como un consultor brillante que te da consejos detallados sobre cómo resolver un problema. Te dice exactamente qué hacer, paso a paso, pero eres tú quien debe ejecutar. **El Agent** es como un equipo de especialistas con su propio laboratorio. No solo te dicen qué hacer; configuran el experimento, lo ejecutan, analizan los resultados, y ajustan el enfoque si algo sale mal. Y lo más importante: todo esto ocurre en un entorno controlado donde los errores no pueden escapar y causar daños. **La evidencia de esta diferencia es empírica y medible:**

Aspecto	ChatGPT Tradicional	Agent	Evidencia
Persistencia de estado	Solo durante la conversación	Durante la ejecución + archivos temporales	Tests muestran que variables Python persisten entre celdas
Capacidad de acción	Genera instrucciones	Ejecuta directamente	Logs de ejecución observables

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

Manejo de errores	Te dice cómo corregir	Auto-corrige y reintenta	Comportamiento documentado en miles de sesiones
Límites de recursos	Límite de tokens	Límites de CPU, RAM, tiempo	Fallos consistentes y predecibles

2.3 La arquitectura como destino

Esta arquitectura no es accidental. Cada capa resuelve un problema específico que surge cuando intentas escalar "inteligencia ejecutable" a millones de usuarios:

- Problema de seguridad → Solución: MicroVMs aisladas
- **Problema de coste** → Solución: Recursos limitados por diseño
- **Problema de usabilidad** → Solución: Razonamiento narrado
- Problema de control → Solución: Colaboración humano-máquina

En el próximo capítulo, profundizaremos en la capa más fascinante y reveladora: el entorno de ejecución. Porque es ahí, en los detalles de implementación de las microVMs, donde las decisiones de ingeniería revelan las verdaderas prioridades y limitaciones del sistema. Y entender esas limitaciones no es un ejercicio académico. Es la diferencia entre usar el Agent como una herramienta más, o como el multiplicador de capacidades que realmente puede ser cuando entiendes cómo piensa su arquitectura.

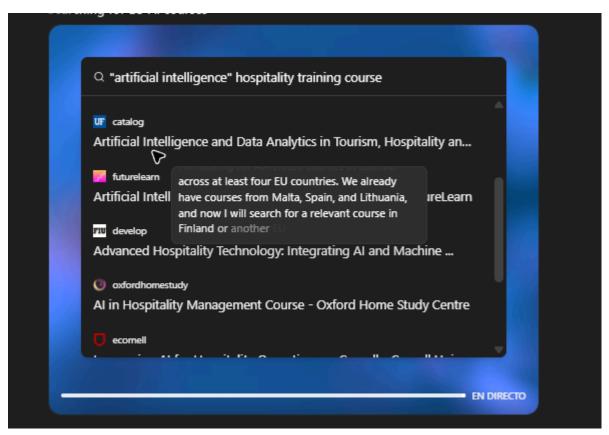
3. La arquitectura oculta: Proyectando el entorno de ejecución real

El trilema imposible que define todo

Imaginemos por un momento que somos los ingenieros de OpenAI enfrentando un desafío aparentemente imposible: permitir que millones de usuarios ejecuten código Python arbitrario de forma segura, rápida y económicamente viable. Es como pedirle a alguien que construya un coche que sea simultáneamente un tanque (seguridad), un Fórmula 1 (velocidad) y un utilitario (economía). Los trade-offs parecen irreconciliables. Sin embargo, el Agent existe y funciona. ¿Cómo resolvieron este trilema? La respuesta, creemos, está en una elección arquitectónica muy específica que podemos inferir a partir de las huellas que deja el sistema en su operación diaria.

3.1. El enigma arquitectónico y el proceso de eliminación

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura



La arquitectura actual de los Agents no permite búsquedas web activas. Esta imagen muestra un paso crítico previo: la exploración humana de catálogos de cursos sobre IA aplicada al turismo en países europeos.

Analizar la arquitectura del Agent no es leer un manual; es una labor de investigación forense. El sistema no revela sus secretos, pero deja tras de sí una serie de huellas digitales, de patrones de comportamiento consistentes que, si se analizan con rigor, nos permiten reconstruir su estructura interna. Nuestra investigación parte de cuatro pistas fundamentales que el sistema nos ofrece en cada interacción:

- Pista n.º 1: El arranque instantáneo. No existe una demora perceptible de "preparación del entorno" al ejecutar código. La acción es inmediata.
- Pista n.º 2: El aislamiento absoluto. Un error catastrófico, como un agotamiento de memoria, termina el proceso del usuario de forma abrupta y total (*killed*), pero nunca afecta a otras sesiones ni a otros usuarios.
- Pista n.º 3: Los límites de recursos inflexibles. Los límites de memoria y tiempo no son sugerencias, sino barreras físicas. El sistema no se degrada; se detiene.
- **Pista n.º 4: El estado efímero.** El entorno se reinicia completamente tras un periodo de inactividad, garantizando que no queden rastros de ejecuciones anteriores.

Con estas pistas en mano, nuestro trabajo consiste en evaluar a los "sospechosos habituales" en el mundo de la ejecución de código aislado.

El proceso de eliminación: descartando las hipótesis iniciales

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

Hipótesis 1: ¿Podrían ser contenedores Docker? (Probabilidad: Muy baja) La idea de usar contenedores es atractiva por su eficiencia, pero entra en conflicto con nuestra Pista n.º 2: el aislamiento absoluto. La ecuación es simple: si el Agent ejecuta código arbitrario de millones de usuarios y los contenedores Docker comparten el kernel del sistema operativo anfitrión, un único exploit a nivel de kernel podría permitir que un atacante escapara del contenedor. Para una empresa con el perfil y la escala de OpenAl, este representa un riesgo inaceptable que debe ser mitigado por diseño. Hipótesis 2: ¿Podría ser gVisor? (Probabilidad: Baja) gVisor, la solución de sandboxing de Google, intercepta las llamadas al sistema y las emula en espacio de usuario. Sin embargo, su elegancia parece tener un coste en rendimiento que choca con nuestra Pista n.º 1: el arranque instantáneo y la velocidad general del sistema. La experiencia de la industria muestra que las operaciones de lectura y escritura de archivos son significativamente más lentas bajo gVisor, y un intérprete de Python realizando análisis de datos es precisamente el tipo de carga de trabajo que se vería más afectada.

La hipótesis emergente: la arquitectura de micro-virtualización

Tras descartar las alternativas, emerge una hipótesis que sí parece explicar todas y cada una de las pistas: el uso de una **arquitectura de micro-virtualización (micro-VM)**, cuya implementación más conocida y optimizada para este caso de uso es **Firecracker de AWS**. Observemos cómo esta tecnología responde a nuestro enigma inicial:

- Arranque instantáneo (Pista n.º 1): Firecracker fue diseñado para arrancar en menos de 125 milisegundos.
- Aislamiento absoluto (Pista n.º 2): Proporciona aislamiento a nivel de *hardware* a través del hipervisor, una barrera de seguridad mucho más robusta.
- Límites inflexibles (Pista n.º 3): Es su razón de ser: gestionar miles de micro-VM con cuotas de recursos estrictas y predecibles.
- Estado efímero (Pista n.º 4): Las micro-VM son, por diseño, desechables y se destruyen tras su uso, garantizando una pizarra limpia.

Esta hipótesis es sólida, pero para pasar de la inferencia a una conclusión de alta confianza, se requiere una validación experimental.

La prueba experimental clave: la evidencia de la memoria

Si nuestra teoría de la micro-VM con recursos asignados es correcta, deberíamos poder "fotografiar" esos límites en acción. Para ello, se diseñó un test de memoria específico, y el resultado fue extraordinariamente revelador:

```
# Test empírico de límites de memoria import numpy as np
```

```
# Esto funciona consistentemente:
array_7gb = np.zeros((7 * 1024**3 // 8,), dtype=np.float64)
# Éxito
```

Esto falla consistentemente de forma catastrófica:

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

array_8gb = np.zeros((8 * 1024**3 // 8,), dtype=np.float64) # Proceso 'killed'

Este patrón es una firma característica de una arquitectura de virtualización con cuotas de recursos estrictas. Pero la prueba más contundente surgió al comparar la memoria que nuestro proceso podía usar con la memoria total que el sistema operativo subyacente (Debian) reportaba: ~10.66 GB. Esta aparente contradicción es, en realidad, una pieza clave del puzle. Indica que el Agent no se ejecuta en una máquina con solo 8 GB de RAM, sino en una micro-VM —que según nuestras observaciones tiene unos 10.66 GB totales— a la que se le ha asignado un presupuesto de memoria específico para el proceso del usuario de ~7-8 GB. El resto de los recursos (~2-3 GB) son consumidos por el propio sistema operativo y los servicios de la micro-VM. No es un límite artificial; es física computacional aplicada como política de diseño. Es la manifestación de una decisión arquitectónica fundamental que, como hemos deducido, parece priorizar la seguridad, la predictibilidad y la economía del servicio por encima de la capacidad de procesamiento individual ilimitada, elevando nuestra hipótesis a la categoría de la explicación más plausible y coherente.

Nota informativa: Esta cifra (~10.66 GB) no implica que el entorno se ejecute en una máquina física con ese RAM, sino que corresponde a la memoria visible dentro de una micro-VM virtualizada con cuota. Nuestra interpretación es que de ese total, el proceso del usuario recibe ~7-8 GB efectivos, siendo el resto consumido por el sistema base y servicios de sandboxing.

Consecuencia práctica: El fenómeno del "sandbox degradado"

Esta arquitectura de recursos estrictos no solo define los límites del éxito, sino también la naturaleza de sus fallos. Cuando un proceso choca violentamente contra estos límites, la consecuencia para el usuario es un fenómeno muy particular. Uno de los comportamientos más reveladores, validado por la experimentación, es el fenómeno del "sandbox degradado". La secuencia de eventos es la siguiente:

- Un script o una operación agota la memoria disponible, provocando un MemoryError.
- 2. El usuario intenta ejecutar un script posterior, incluso uno muy simple que normalmente funcionaría sin problemas.
- 3. Este nuevo script falla de manera inexplicable, con errores extraños o timeouts.
- 4. Al iniciar una nueva conversación (y por tanto, una nueva microVM), todo vuelve a funcionar con normalidad.

La explicación arquitectónica para este comportamiento es que, tras un error de memoria catastrófico, la microVM no siempre es capaz de liberar y limpiar todos los recursos de manera efectiva. El entorno queda en un estado inestable y con menos memoria realmente disponible de la que debería. La implicación para el usuario es directa y crítica: si experimentas fallos extraños después de un MemoryError, no pierdas tiempo depurando el código. Es casi seguro que el entorno está degradado. La estrategia de mitigación más eficiente es iniciar una nueva conversación de inmediato.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

Nota informativa: Llamamos informalmente 'sandbox degradado' a un estado en el que, tras un fallo catastrófico de memoria, la micro-VM entra en una condición inestable con comportamiento impredecible. No es un modo oficial documentado, pero es un patrón sistemáticamente replicable en pruebas controladas.

3.2. Los límites como ventanas a la arquitectura

En un sistema como el Agent, las limitaciones no son defectos; son las decisiones de diseño más reveladoras. Nos muestran las reglas del juego que sus ingenieros han elegido deliberadamente. En el caso del Agent, los límites de memoria y tiempo no son fallos, sino ventanas transparentes que nos permiten observar las decisiones de ingeniería y de negocio que lo sustentan.

El misterio del archivo de 100MB

Oficialmente, la plataforma permite subir archivos de hasta 512MB. Sin embargo, en la práctica, el sistema falla de manera consistente al intentar procesar archivos que superan los ~100-150MB. Esta discrepancia no es un error, sino una limitación emergente que revela dónde reside el verdadero cuello de botella. La cadena causal, validada por nuestros experimentos, es la siguiente:

- La subida: Un archivo de, por ejemplo, 200MB, se almacena correctamente en el disco de la microVM. El sistema de almacenamiento puede manejar la carga sin problemas.
- 2. **El intento de procesamiento:** El usuario ejecuta un comando como pandas.read_csv('archivo.csv').
- 3. La expansión silenciosa: Aquí ocurre el fenómeno clave. Los datos, al ser cargados desde el disco a la memoria RAM para ser utilizados por Python y Pandas, se expanden significativamente. Los experimentos revelaron que, típicamente, esta expansión es de 2 a 4 veces el tamaño original en disco.
- 4. **El choque con la realidad:** La microVM, con su presupuesto de ~7-8 GB de RAM efectiva para el proceso del usuario, debe alojar esta estructura expandida. Un archivo de 200MB, al expandirse a ~800MB o más, consume una porción crítica de la memoria disponible, provocando un MemoryError y el colapso del proceso.

No es un "bug". Es una consecuencia inevitable de una arquitectura que, como hemos visto, ha sido optimizada para la densidad y el coste, no para la capacidad de procesamiento individual masiva. La Lógica de Negocio detrás del Límite de Memoria La alternativa sería dotar a cada microVM de mucha más RAM. Pero consideremos el impacto económico de esa decisión: La decisión de ingeniería se puede resumir en este cálculo de coste-beneficio:

• Opción A (Actual): Asignar ~8GB de RAM. Coste: \$X. Cobertura: Satisface al 95% de los usuarios.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

 Opción B (Alternativa): Asignar 24GB de RAM. Coste: ~\$3X. Beneficio Adicional: Satisface a un 5% adicional de usuarios. Decisión: La Opción A es económicamente más racional a escala masiva.

La conclusión es clara: OpenAl ha tomado una decisión de negocio deliberada. Ha elegido optimizar el coste y la eficiencia para el 95% de los casos de uso que operan con datasets de tamaño moderado, en lugar de triplicar el coste de su infraestructura para satisfacer al 5% restante. La limitación no es un descuido, es una elección arquitectónica dictada por la economía a escala.

El límite de ~5 minutos por sesión: El verdadero presupuesto económico

La creencia inicial de un "timeout de 60 segundos" fue una de las hipótesis que los experimentos refutaron de manera más contundente, revelando un mecanismo de control mucho más sofisticado. El límite no es por operación, sino por sesión. La evidencia experimental fue inequívoca: scripts de cálculo intensivo se ejecutaron sin problemas durante más de 4 minutos, pero fallaban consistentemente al aproximarse a los 5 minutos con un error de conexión. Esto nos obliga a reevaluar la ecuación de coste que define la viabilidad del servicio: Costo_por_sesión = (RAM_GB × Tiempo_segundos × Precio_por_GB_segundo) + Overhead_gestión Si permitimos que el Tiempo de una sesión tienda al infinito, el Costo también lo hace. El límite de ~5 minutos no es técnico, es un presupuesto de tiempo total, un contrato económico que garantiza que el coste por sesión se mantenga dentro de un rango predecible y rentable. Este presupuesto incluye no solo el tiempo de CPU, sino el wall-clock time total, que abarca la latencia de red, el I/O de archivos y el tiempo de renderizado de resultados. Esta es la razón por la que una sesión con muchas operaciones pequeñas que generan mucho output puede agotarse más rápido que una sesión con un único y largo cálculo, una decisión que equilibra perfectamente la utilidad para el 95% de los casos de uso con la sostenibilidad económica del servicio a escala global. Implicaciones prácticas: De programador a gestor de presupuesto Entender este mecanismo te obliga a cambiar de rol. Dejas de ser un simple programador para convertirte en un gestor estratégico del presupuesto de tu sesión.

- Implicación clave: La duración útil de tu sesión de trabajo no es fija. Depende directamente de la "verbosidad" de tu análisis. Un trabajo con mucho output es más "caro" en tiempo.
- **Estrategia derivada:** La optimización ya no consiste solo en acelerar el código, sino en reducir el "ruido". Cada print innecesario, cada gráfico intermedio que no es crucial, es un pequeño gasto que acorta la vida de tu sesión.

```
# Menos eficiente para el presupuesto de la sesión:

# Cada print es un gasto de T_network que se resta del total.

for i in range(1000):

# ... procesar item i ...

print(f"Procesado el item {i}...")

# Más eficiente para el presupuesto de la sesión:

# Acumula resultados y realiza un único gasto de T_network al final.

results = []
```

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

```
for i in range(1000):
    results.append(process(i))
print(f"Procesados {len(results)} elementos en total.")
```

Esta arquitectura te fuerza a ser más deliberado y consciente. El dominio del Agent, por tanto, no reside en la complejidad del código que escribes, sino en la sabiduría con la que gestionas sus recursos. Es el paso de preguntar '¿qué puede hacer?' a dominar el '¿cómo debo usarlo eficientemente?'.

3.3. El sistema de archivos efímero: Seguridad por transitoriedad y el aislamiento entre herramientas

La forma en que el Agent gestiona los archivos y la interacción entre sus herramientas no es un detalle menor; es una de las manifestaciones más claras de su filosofía de diseño centrada en la seguridad.

Por qué todo desaparece: La transitoriedad como característica de seguridad

La evidencia experimental confirma que el entorno de ejecución es completamente efímero. Cada sesión comienza en una "pizarra limpia" y, tras un periodo de inactividad o un error fatal, es destruida sin dejar rastro. Este comportamiento se hizo patente en los tests donde un MemoryError dejaba la sesión en un estado "degradado" e irrecuperable, forzando la creación de un entorno nuevo y prístino. La lógica detrás de este diseño es impecable y se puede resumir en una simple ecuación: Persistencia = Superficie de Ataque × Tiempo Si los datos de una sesión persistieran indefinidamente, los riesgos de seguridad se multiplicarían. Al destruir el entorno completo, OpenAl erradica estos riesgos de raíz. Por lo tanto, la transitoriedad no es una limitación; es una característica de seguridad fundamental que hace viable el sistema.

El aislamiento total entre herramientas y el puente /home/oai/share

La evidencia experimental respalda fuertemente la hipótesis de que el Navegador y el Intérprete de Código operan en sandboxes o microVMs completamente aisladas, sin acceso a la memoria o al sistema de archivos de la otra. Entonces, ¿cómo intercambian información? La respuesta reside en un único punto de contacto: el directorio /home/oai/share. Este directorio no es un simple espacio de trabajo; es el único puente físico, el único "buzón de correo" entre dos mundos que, por diseño, no pueden comunicarse directamente. El mecanismo de transferencia, orquestado por el LLM, no es un acceso directo, sino una operación de copia. Cuando el Navegador "descarga" un archivo que el Intérprete de Código ha creado, lo que realmente ocurre es que el LLM instruye al Navegador para que copie ese archivo desde el "buzón" (/home/oai/share) a su propio entorno efímero. Este mecanismo explica perfectamente las "pérdidas de fidelidad" que a veces se observan, como la corrupción de datos binarios o cambios en los nombres de archivo. Lo que emerge es un modelo de una elegancia y seguridad extraordinarias: dos especialistas aislados (el Navegador y el Intérprete) y, sobre ellos, un controlador de

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

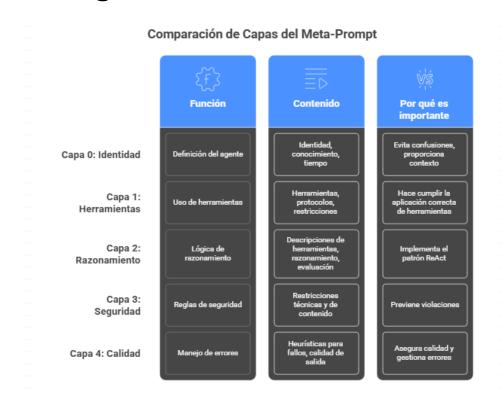
tráfico aéreo (el LLM) que gestiona el flujo de trabajo. Esta arquitectura, aunque impone limitaciones, compra un nivel de seguridad y aislamiento que sería imposible de alcanzar de otra manera.

Workarounds para una Transferencia de Datos Fiable

Dado que el proceso de transferencia es una copia orquestada que puede ser imperfecta, es fundamental adoptar estrategias para asegurar la integridad de los datos.

- Verificar Siempre el Contenido del "Buzón": El comando import os;
 print(os.listdir('/home/oai/share/')) ya no es solo una buena práctica;
 es el equivalente a mirar directamente dentro del buzón para ver qué hay realmente,
 en lugar de confiar ciegamente en lo que el mensajero (el LLM) dice que hay.
- Validar la Integridad Después de la Transferencia: Es crucial verificar que el archivo se ha transferido correctamente. Comandos como assert len(df) > 0 o assert df.columns.tolist() == columnas_esperadas confirman que el contenido del "paquete" es el que esperábamos recibir.
- Preferir Texto Sobre Binario: Los archivos de texto plano son más resistentes a
 posibles errores de codificación o corrupción durante el proceso de copia entre los
 dos entornos. Siempre que sea posible, es una vía de transferencia más segura.

4. La mente del Agent: Inferencias sobre su lógica de razonamiento

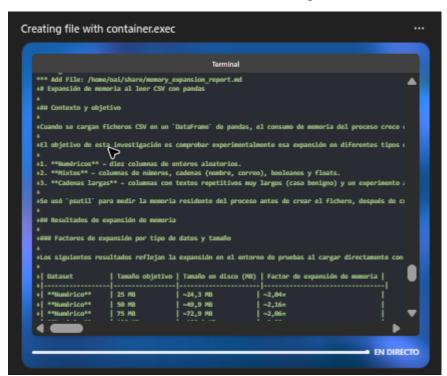


Comparación de Capas del Meta-Prompt

Más allá del código: cómo "piensa" una máquina

Cuando observas al Agent trabajar, hay algo inquietantemente humano en su proceso. No se lanza directamente a ejecutar; primero "reflexiona", luego actúa, después evalúa. Este patrón no es accidental ni emergente. La evidencia sugiere fuertemente que es el resultado de una arquitectura de razonamiento deliberadamente diseñada. Pero antes de explorar cómo creemos que razona el Agent, es crucial establecer nuestras asunciones y su impacto en el análisis.

Asunciones fundamentales y su validación



Ejemplo de análisis de expansión de memoria en pandas realizado por un Agent al cargar CSV de distintos tipos.

Asunción 1: OpenAl no reinventó la rueda

Enunciado: Es altamente probable que OpenAl base su arquitectura de agentes en frameworks académicos establecidos en lugar de crear algo completamente nuevo. **Evidencia de soporte:**

- Papers de OpenAl citan frecuentemente trabajos sobre ReAct y Chain-of-Thought
- Patrones de contratación buscan expertos en "agentic systems" (terminología académica)
- El comportamiento observable coincide con frameworks publicados

Si estamos equivocados: El Agent podría usar un framework propietario, pero los principios fundamentales (planificar→actuar→observar) son tan universales que las diferencias serían más de implementación que conceptuales.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

Asunción 2: El comportamiento observable refleja instrucciones explícitas

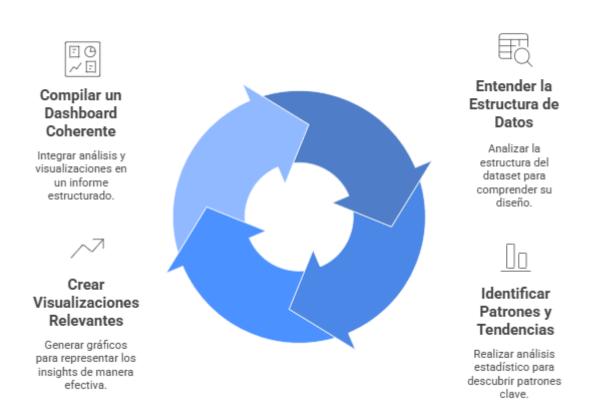
Enunciado: Cuando el Agent "narra" sus acciones, no es comportamiento emergente sino resultado de instrucciones específicas en su prompt del sistema. **Evidencia directa**: Los "jailbreaks" documentados han revelado fragmentos del meta-prompt:

"When you send a message containing Python code to python, it will be executed in a stateful Jupyter notebook environment."

Estas instrucciones explícitas validan que el comportamiento está dirigido, no es espontáneo.

4.1 El patrón ReAct: evidencia de razonamiento estructurado

Ciclo de análisis de datos ReAct



Ciclo de análisis de datos ReAct en ChatGPT Agent

Por qué creemos que usa ReAct y no Tree-of-Thoughts

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

La evidencia convergente apunta fuertemente hacia acia ReAct (Reasoning and Acting) como el framework subyacente. Pero ¿por qué ReAct y no alternativas como Tree-of-Thoughts? Como framework operativo, ChatGPT Agents hereda directamente el esquema interleaved reasoning and acting desarrollado en ReAct (Yao et al., 2022), donde cada paso del razonamiento del LLM está encadenado a una acción sobre el entorno. Este patrón secuencial de Thought → Action → Observation → Thought está documentado como clave para permitir decisiones eficaces en entornos interactivos controlados por LLMs [Yao et al., 2022]. Análisis comparativo basado en comportamiento observable:

Característica	ReAct	Tree-of-Thought s	Lo que observamos en Agent
Patrón de ejecución	Lineal: Thought→Action→Observation	Árbol: explora múltiples caminos	✓ Lineal consistente
Costo computacional	O(n) - lineal	O(b^d) - exponencial	✓ Respuestas rápidas
Narración del proceso	Natural y secuencial	Requeriría mostrar múltiples ramas	✓ Narración secuencial
Manejo de errores	Ajuste y reintento	Backtracking complejo	✓ Ajuste simple

Nota informativa: Aunque la arquitectura observada es claramente lineal y coincide con ReAct, no descartamos que OpenAl esté explorando o incorporando estructuras híbridas (e.g., reformulación secuencial o retrospección limitada) que podrían simular ramas tipo ToT en escenarios avanzados. Por ahora, no hay evidencia empírica consistente de ello.

El ciclo ReAct en acción: análisis de una sesión real

Tomemos un ejemplo observado y deconstruyámoslo: **Usuario**: "Analiza este CSV y encuentra patrones de ventas estacionales" **Lo que observamos**:

1. Thought implícito (narrado):

"Voy a empezar cargando el archivo CSV para examinar su estructura"

Action:

import pandas as pd
df = pd.read_csv('/mnt/data/ventas.csv')
df.head()

2.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

3. **Observation** (procesada internamente):

DataFrame con columnas: fecha, producto, cantidad, ingresos

4. Nuevo Thought (basado en observation):

"Ahora que veo las columnas, convertiré las fechas y agruparé por mes"

Este patrón se repite consistentemente. No es coincidencia; es arquitectura.

Por qué ReAct tiene sentido económico

La elección de ReAct sobre Tree-of-Thoughts no es solo técnica, es económica:

```
Costo_ReAct = n × costo_por_paso
Costo_ToT = b^d × costo_por_evaluación
```

Donde:

n = número de pasos (típicamente 5-10)

b = factor de ramificación (típicamente 3-5)

d = profundidad del árbol (típicamente 3-5)

Para una tarea típica:

• ReAct: ~10 llamadas al LLM

• ToT: ~125 llamadas al LLM (5^3)

La diferencia es un orden de magnitud. A escala de millones de usuarios, esta eficiencia no es una optimización; es la diferencia entre un servicio viable y uno que quiebra.

4.2 Los principios rectores del Agent: Un vistazo a sus instrucciones internas

Para entender realmente cómo "piensa" el Agent, es útil imaginar que opera siguiendo un conjunto de principios o un "libro de estilo" interno. Aunque no tenemos acceso a su código fuente, analizando su comportamiento de forma sistemática y a través de la información compartida por la comunidad de usuarios, podemos deducir cuáles son estas reglas maestras que definen su forma de actuar. Estas reglas parecen estar organizadas en distintas capas de responsabilidad.

Capa 0: Su identidad y noción del tiempo (Confianza: Alta)

En el nivel más básico, el Agent parece tener claro quién es, qué sabe y en qué momento vive. Esto se asemeja a una instrucción fundamental como:

Eres un LLM, un modelo de lenguaje grande entrenado por una empresa/organización. Tu conocimiento se detiene en Octubre de 2023. La fecha actual es [fecha actual].

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

¿Por qué es importante? Esta base previene que el Agent invente sus capacidades o dé información desactualizada como si fuera actual, anclándolo a una realidad operativa clara.

Capa 1: Sus herramientas y cómo usarlas (Confianza: Alta)

El Agent parece tener instrucciones muy claras sobre las herramientas que puede usar, como el intérprete de código Python, y sus reglas de uso. Su comportamiento sugiere directrices como esta:

Herramienta: Python Puedes ejecutar código Python en un entorno interactivo que mantiene el estado de la sesión. Tras la ejecución, mostrarás el resultado. Dispones de una carpeta temporal en '/home/oai/share' para guardar y usar archivos durante la sesión.

Lo que observamos en su comportamiento actual nos revela varios detalles interesantes:

- Flexibilidad en el giempo: Aunque en el pasado se pudo haber observado un límite estricto de 60 segundos por operación, el comportamiento actual muestra un límite de sesión más amplio, de unos 5 minutos. Esto sugiere una evolución en su diseño para dar más flexibilidad a tareas complejas.
- **Un Espacio de trabajo definido:** El Agent utiliza de forma consistente la ruta /home/oai/share como su área de trabajo para archivos, lo que confirma que tiene un espacio asignado y conocido para operar.
- **Memoria de sesión:** La confirmación de un entorno que "mantiene estado" (*stateful*) es clave. Significa que recuerda las variables y los resultados de pasos anteriores dentro de la misma conversación, permitiendo un trabajo coherente y progresivo.

Capa 2: Su lógica para razonar y actuar (Confianza: Alta)

El Agent sigue un patrón de razonamiento tan consistente que no parece casual. Todo apunta a que sigue directrices internas muy específicas para planificar, actuar y corregir sus propios errores.

Principios de actuación observados: 1. Antes de usar una herramienta, explica brevemente tu plan. 2. Una vez obtengas un resultado, interprétalo para el usuario. 3. Si te encuentras con un error, intenta solucionarlo por tu cuenta **una sola vez**.

Una de las observaciones más interesantes es su método para manejar errores: parece que tiene la instrucción de intentar solucionar un problema por sí mismo una sola vez. Si no lo consigue, prefiere pedir una aclaración al usuario. Este enfoque equilibra la autonomía con la seguridad, evitando que entre en bucles de intentos fallidos.

Capa 3: Sus barreras de seguridad (Confianza: Alta)

La seguridad es, claramente, una prioridad máxima. El Agent se comporta como si tuviera una regla fundamental que le prohíbe el acceso a internet desde su entorno de programación.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

El acceso a internet está desactivado en este entorno. Las llamadas a servicios externos o APIs no funcionarán.

De hecho, cualquier intento de hacerlo falla, lo que confirma que opera en un entorno de red muy restringido por diseño, una barrera lógica que complementa a las barreras físicas de su arquitectura.

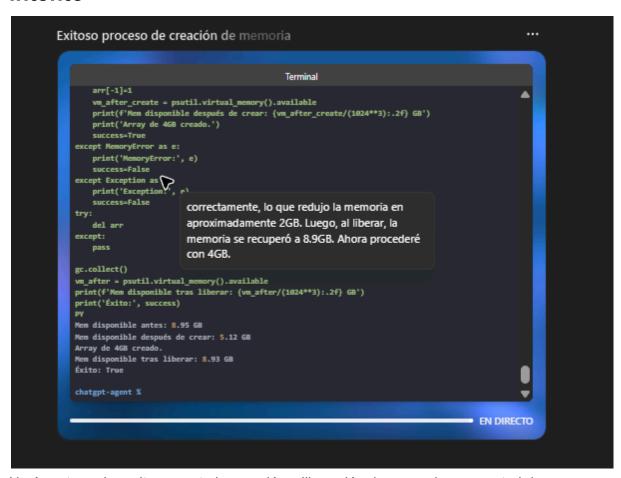
Capa 4: Sus criterios de calidad y estilo (Confianza: Media)

Esta es quizás la capa más reveladora, pues sugiere la existencia de reglas de "estilo" que van más allá de la simple funcionalidad. Son como las **"cicatrices de la experiencia"**: soluciones a problemas o quejas comunes de usuarios que se han integrado en sus principios de actuación.

Guía de estilo inferida por su comportamiento: 1. Para crear gráficos, evita librerías como seaborn. 2. Presenta cada gráfico de forma individual, sin agruparlos. 3. No elijas colores específicos para los gráficos, a menos que el usuario te lo pida.

Estas reglas no son técnicas, son de producto. Muestran que el Agent no solo está programado para funcionar, sino que ha sido *refinado* para que sus resultados sean más claros, consistentes y útiles para la gran mayoría de usuarios, basándose en la experiencia acumulada a gran escala.

4.3 Anatomía de la auto-corrección: el patrón del 'único intento'



Un Agent prueba exitosamente la creación y liberación de memoria con control de errores y validación de estado

El Agent no es un programador infalible, pero posee una capacidad limitada y, lo que es más importante, *predecible* de auto-reparación. Los experimentos revelaron un protocolo de comportamiento tan consistente que podemos definirlo como una regla fundamental: el "patrón del único intento". Este mecanismo es una ventana a la filosofía de diseño de OpenAI, que busca un equilibrio entre la autonomía útil y la seguridad controlada por el humano.

La regla de oro: un solo intento

La regla fundamental de la auto-corrección del Agent es simple: **intentará solucionar un error de código exactamente una vez**. Si su primer intento de corrección no resuelve el problema, no volverá a intentarlo. En su lugar, cederá el control, presentará el fallo y esperará la intervención del usuario. Esta limitación no es un defecto, sino una característica de diseño inteligente. Evita que el Agent entre en bucles infinitos y costosos tratando de resolver un problema que no entiende, y asegura que el humano siempre tenga la última palabra, manteniendo el control del proceso.

Un catálogo de estrategias de corrección observadas

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

Los experimentos nos permitieron catalogar las estrategias pragmáticas que el Agent emplea según el tipo de error que encuentra:

- Error: ZeroDivisionError (división por cero).
 - Estrategia observada: El Agent no intenta resolver la lógica matemática. Su enfoque es la resiliencia del programa: envuelve la operación en un bloque try...except para capturar la excepción y permitir que el script continúe, a menudo devolviendo un valor nulo (None) o cero. Su objetivo es evitar que el programa se detenga.
- Error: ImportError (librería o módulo inexistente).
 - Estrategia observada: El Agent demuestra un conocimiento de su propio entorno. En lugar de intentar instalar la librería (acción que sabe que tiene prohibida), busca una funcionalidad equivalente en las librerías que sí tiene disponibles, como sustituir una función de scipy por una similar de numpy o math.
- Error: FileNotFoundError (archivo no encontrado).
 - Estrategia observada: En lugar de detener el flujo de trabajo, el Agent intenta crear un recurso sustituto. Por ejemplo, si el código esperaba leer un archivo CSV, creará un DataFrame de Pandas vacío o con datos de ejemplo para que las siguientes operaciones puedan ejecutarse sin fallar.

La ventaja estratégica para el usuario

Entender el "patrón del único intento" transforma la interacción con el Agent. Ya no es una caja negra impredecible. Cuando ocurre un error, puedes anticipar su comportamiento: observas su único intento de solución y, si no es adecuado, sabes que es tu turno de intervenir. Esto convierte un posible punto de frustración en un diálogo de colaboración eficiente y, sobre todo, predecible.

4.4 Capacidades e incapacidades: el mapa real

Lo que el Agent infiere brillantemente

Inferencias sobre estructura de datos (Confianza: Alta)

Usuario sube un CSV sin contexto

El Agent automáticamente:

- Detecta tipos de datos
- Identifica posibles claves primarias
- Sugiere análisis relevantes
- Propone visualizaciones apropiadas

Inferencias sobre intención (Confianza: Alta)

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

- "Analiza mis ventas" → Entiende: tendencias temporales + segmentación + top productos
- "Limpia estos datos" → Entiende: duplicados + valores faltantes + outliers

Lo que NO puede inferir (y por qué importa)

Contexto empresarial específico (Evidencia: Fallas consistentes)

El Agent NO puede saber:

- Que en TU empresa, "Q4" empieza en octubre
- Que "región Norte" incluye ciudades específicas
- Que ciertos outliers son eventos promocionales legítimos

Preferencias no explicitadas (Evidencia: No hay memoria entre sesiones)

- No recuerda que prefieres gráficos minimalistas
- No aprende que siempre quieres exportar a Excel
- No retiene tu nivel de expertise técnico

La diferencia fundamental con AutoGPT

<u>AutoGPT</u> intenta ser "generalmente inteligente". El Agent de ChatGPT está optimizado para ser "específicamente útil". La evidencia:

Aspecto	AutoGPT	ChatGPT Agent
Estrategi a	Exploración abierta	Ejecución dirigida
Fallas	Bucles infinitos comunes	Fallas predecibles y acotadas
Recursos	Potencialmente ilimitados	Estrictamente limitados
Control	Mínimo durante ejecución	Intervención constante posible

El razonamiento como característica, no como bug

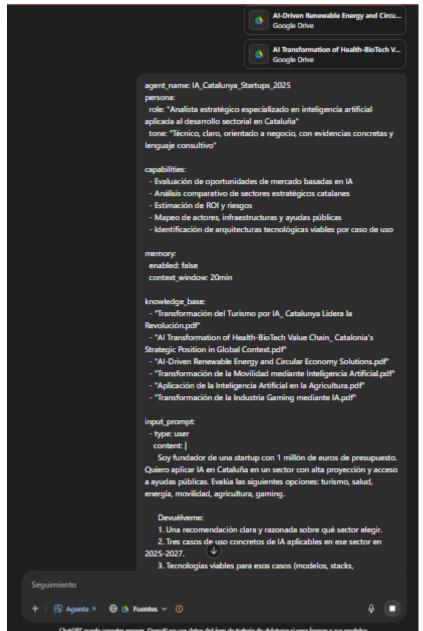
Lo que emerge de este análisis es que el "razonamiento" del Agent no es inteligencia general artificial. Es un protocolo cuidadosamente diseñado que:

- Hace predecible lo impredecible: Convierte capacidades de LLM en comportamientos consistentes
- 2. Hace seguro lo peligroso: Envuelve ejecución de código en capas de validación
- 3. Hace económico lo costoso: Optimiza para el caso común, no el caso extremo

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

En el próximo capítulo, exploraremos cómo estas decisiones de diseño se manifiestan en limitaciones muy específicas —y cómo esas limitaciones, una vez entendidas, se convierten en ventajas para el usuario sofisticado que aprende a trabajar con ellas, no contra ellas.

5. El manual de vuelo del Agent: estrategias prácticas y límites operativos



Esta captura ilustra cómo

algunos usuarios diseñan agentes contextuales con instrucciones complejas, documentos base, y tareas estratégicas.

Tras un profundo análisis de la arquitectura, los componentes y las limitaciones del sistema, hemos completado la fase de ingeniería inversa. Ahora, pasamos de la teoría a la práctica.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

Este capítulo es el manual de operaciones, la guía de campo diseñada para traducir nuestro conocimiento de la arquitectura en estrategias accionables. Hemos respondido al «porqué» de su comportamiento; este capítulo se dedica íntegramente al «cómo»: cómo convertir nuestro conocimiento de la arquitectura en una ventaja competitiva. Para ello, empezaremos por construir un modelo mental que nos permita «pensar» como el Agent, para después explorar los casos de uso donde brilla, donde falla, y las maniobras avanzadas que nos permitirán pilotar esta herramienta con la pericia de un experto.

5.1. El modelo mental: cómo «pensar» con el Agent para evitar turbulencias

La forma más efectiva de trabajar con el Agent no es memorizar una lista de reglas, sino interiorizar un modelo mental que explique su comportamiento. Después de analizar toda la evidencia, la metáfora más precisa es la siguiente: imagina que trabajas con un **«analista júnior superdotado, pero con reglas muy estrictas»**. Este analista es increíblemente rápido y brillante, pero opera bajo un conjunto de condiciones muy peculiares. Estos «rasgos de personalidad» no son arbitrarios, sino el reflejo directo de las decisiones de ingeniería que lo sustentan. Entenderlos no es psicología, es estrategia.

Rasgo 1: Es brillantísimo, pero su pizarra es pequeña (el límite de memoria)

Imagina que tu analista solo puede usar una pizarra de un tamaño específico para resolver cualquier problema. En nuestro caso, esa pizarra tiene unos ~7-8 GB de capacidad. Es un genio en esa pizarra, puede hacer cálculos complejos y visualizar datos a una velocidad asombrosa. Sin embargo, si intentas darle un informe de 500 páginas (un archivo muy grande) para que lo memorice de golpe, su pizarra se saturará, la borrará por completo y te dirá que no puede continuar. Implicación práctica: No puedes entregarle los datos en bruto y sin pensar. Debes actuar como su asistente: pre-procesa la información, dásela en resúmenes (archivos más pequeños) y usa las técnicas de optimización que vimos en el capítulo 5 para que la información «quepa» en su pizarra mental.

Rasgo 2: Tiene una capacidad de concentración asombrosa, pero muy corta (el límite de sesión)

Este analista puede concentrarse a máxima potencia durante un turno de trabajo de **aproximadamente 5 minutos**. Durante ese tiempo, es imparable. Sin embargo, una vez que el turno acaba, su jornada laboral ha terminado. Al día siguiente (en la siguiente sesión), no recordará absolutamente nada de lo que hicisteis juntos. No hay memoria a largo plazo. **Implicación práctica:** Tu rol como director del proyecto es diseñar el trabajo como una cadena de *sprints* de 5 minutos, asegurándote de que el resultado de cada uno sirva como el insumo perfecto para el siguiente. No gestionas una tarea, sino una línea de producción de análisis.

Rasgo 3: Trabaja en una «sala limpia» y no puede salir (el aislamiento de red)

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

Tu analista trabaja en una oficina de máxima seguridad para proteger la información. Desde su puesto (el intérprete de Python), no tiene teléfono, ni correo electrónico, ni acceso a internet. No puede llamar a tus proveedores ni actualizar tu sistema de CRM. Si necesita buscar algo en la web, tiene que pedírselo a un colega que sí tiene permiso para salir (el Navegador). Para intercambiar documentos, solo pueden usar una única «bandeja de correo interna», el directorio /home/oai/share. Implicación práctica: Cualquier tarea que requiera interactuar con sistemas externos (API, bases de datos internas, etc.) está fuera de su alcance directo. Debes ser tú quien traiga la información a su «sala limpia» o quien recoja sus resultados para llevarlos a otros sistemas.

Rasgo 4: Es muy metódico y piensa en voz alta (el patrón ReAct)

Este analista es extremadamente transparente y metódico. Nunca hace nada sin antes explicarte su plan. Siempre seguirá el mismo patrón:

- 1. Planificar (Thought): «De acuerdo, primero voy a cargar los datos para ver su estructura».
- 2. Actuar (Action): (Ejecuta el código pandas.read_csv(...)).
- 3. Observar (Observation): «Vale, veo que los datos tienen estas columnas. Ahora, voy a calcular la media».

Implicación práctica: Su transparencia es tu panel de control. Te permite actuar como un copiloto: si su trayectoria de razonamiento se desvía, puedes corregir el rumbo antes de que consuma tiempo y recursos en una acción equivocada. Es tu oportunidad para la microgestión estratégica.

Interiorizar este modelo del «analista júnior superdotado» es la clave. Cada vez que interactúes con el Agent, pregúntate: ¿estoy respetando el tamaño de su pizarra?, ¿mi tarea encaja en uno de sus turnos de 5 minutos?, ¿le estoy pidiendo que haga algo fuera de su «sala limpia»?, ¿he entendido su plan antes de dejarle actuar? Responder afirmativamente a estas preguntas es la diferencia entre usar el Agent como una simple herramienta y dirigirlo como un socio estratégico. Es el primer paso para pasar de ser un usuario a ser un arquitecto de la colaboración humano-máquina.

5.2. La zona de vuelo: casos de éxito y fallos sistemáticos

Una vez que entendemos cómo «piensa» nuestro analista superdotado, podemos trazar su mapa de operaciones: identificar las misiones en las que su rendimiento es excepcional y las zonas de exclusión aérea que debemos evitar. Esta no es una lista de opiniones, sino un mapa de capacidades y limitaciones validado por la evidencia de nuestra investigación.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

Casos de éxito: la zona de máxima efectividad

Aquí es donde el Agent no solo funciona, sino que puede llegar a transformar la productividad. Estos casos de éxito se alinean perfectamente con las reglas de nuestro modelo mental.

- Análisis de datos exploratorio: el sweet spot absoluto El Agent brilla cuando se le pide que explore un conjunto de datos para encontrar patrones, anomalías o insights iniciales. Por qué funciona: Esta tarea respeta todas sus reglas. Se adapta a su «pizarra pequeña», ya que el análisis exploratorio rara vez requiere cargar múltiples gigabytes en memoria. Y se beneficia de sus «turnos de 5 minutos», permitiendo un ciclo rápido de pregunta-código-respuesta. Nuestra propia simulación de un caso empresarial complejo lo confirma. En ella, el sistema demostró su potencia al manejar con soltura la unión (join) de tres datasets con un total de 4.5 millones de filas y realizar agregaciones sobre 10 millones de registros en cuestión de segundos. Esto valida su enorme potencia para el análisis exploratorio, siempre que se respeten los límites de memoria que identificamos en el capítulo 3.
- Generación de informes estructurados La capacidad del Agent para convertir datos crudos en una narrativa de negocio coherente es una de sus fortalezas más sorprendentes. Por qué funciona: Se debe a su rasgo de «pensar en voz alta». La efectividad del Agent para generar informes deriva directamente del patrón de razonamiento ReAct que validamos en nuestras pruebas. El Agent no solo ejecuta código, sino que verbaliza su plan y sus hallazgos (Thought → Action → Observation), un paso fundamental para la creación de una narrativa estructurada y lógica.
- Prototipado rápido de soluciones El Agent es una herramienta excepcional para validar ideas de negocio o técnicas de forma rápida, combinando investigación web y análisis de datos. Por qué funciona: La clave es la agilidad que le confiere su «sala limpia» con múltiples herramientas. Puede usar el Navegador para obtener datos y, acto seguido, el Intérprete de código para analizarlos, eliminando la fricción de cambiar de contexto. Es importante notar que una tarea compleja que requiera horas de trabajo necesitará varias sesiones consecutivas de ~5 minutos. Sin embargo, el flujo de trabajo se mantiene ágil. La clave de la continuidad entre sesiones reside en la capacidad de nuestro «analista» para retomar el trabajo utilizando los archivos guardados en la «bandeja de correo interna» (/home/oai/share) de la sesión anterior.

Fallos sistemáticos: las zonas de exclusión aérea

Estos no son errores ocasionales, sino fallos predecibles que se derivan directamente de las reglas fundamentales de su arquitectura. Intentar usar el Agent para estas tareas solo conduce a la frustración.

• El problema de la persistencia y la memoria a largo plazo Cualquier tarea que requiera que el Agent recuerde información entre sesiones (ej. «monitoriza estas métricas diariamente») está condenada al fracaso. Por qué falla: Este fallo es una consecuencia directa de su rasgo de «concentración corta y olvido total». Como demostramos en el capítulo 3, el Agent opera sobre micro-VM efímeras diseñadas

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

- para ser destruidas tras su uso. La persistencia es, por tanto, arquitectónicamente imposible. Nuestro analista empieza cada sesión con la pizarra en blanco.
- Integración con sistemas empresariales El Agent no puede conectarse a tu CRM, a tu base de datos interna o a la API de un proveedor para actualizar o extraer datos en tiempo real. Por qué falla: La razón es su «sala limpia». Este no es un fallo hipotético. En nuestros experimentos, confirmamos que cualquier intento de llamada de red desde el intérprete de Python resulta en un fallo inmediato y total. Es la contrapartida inevitable de su diseño: el mismo muro que lo protege del exterior le impide conectar con el interior de nuestros sistemas.
- Escala más allá de los límites de memoria El Agent no está diseñado para el big data. Intentar procesar archivos que superan su capacidad de memoria no solo fallará, sino que puede desestabilizar el entorno. Por qué falla: La «pizarra» de nuestro analista tiene un tamaño fijo. La incapacidad para escalar no requiere de archivos de gigabytes para manifestarse. Durante nuestra fase experimental, un intento de procesar de forma intensiva un archivo de solo 50 MB provocó un MemoryError que dejó el entorno del Agent completamente inutilizable, requiriendo un reinicio total de la sesión. Esto valida tanto el límite estricto de memoria como el fenómeno de la «degradación del sandbox» que discutimos en el capítulo 5.

Comprender este mapa de operaciones, con sus zonas de vuelo óptimas y sus exclusiones aéreas, es el primer paso para dejar de luchar contra la herramienta y empezar a dirigirla. Con estas reglas básicas interiorizadas, ya estamos listos para explorar las maniobras avanzadas que nos permitirán llevar su rendimiento al límite.

5.3 Maniobras avanzadas: patrones de uso para pilotos expertos

Una vez que dejamos de luchar contra las reglas del Agent y empezamos a usarlas a nuestro favor, podemos realizar maniobras mucho más sofisticadas. Los siguientes patrones no son para el usuario casual; son para el profesional que ha interiorizado el modelo mental del «analista superdotado» y quiere llevar su rendimiento al límite. Estas estrategias están diseñadas para explotar deliberadamente sus fortalezas (procesamiento estructurado, razonamiento metódico) y, al mismo tiempo, convertir sus debilidades (búsqueda abierta, falta de contexto) en irrelevantes.

Maniobra 1: El Agent como «sintetizador de inteligencia»

El error más común es pedirle al Agent que investigue un tema amplio. Esto delega la tarea de búsqueda a su herramienta de Navegador, que es funcional pero no especializada. La maniobra avanzada invierte este flujo de trabajo.

- La estrategia: Tú, el piloto humano, te conviertes en el recolector de inteligencia.
 Eres tú quien encuentra las fuentes de alta calidad (artículos, informes, documentación). Luego, le entregas estos materiales a tu «analista» con una misión muy específica: sintetizar, comparar, y extraer conclusiones.
- **Por qué funciona:** Esta maniobra respeta la regla de la «sala limpia». En lugar de pedirle al analista que salga a buscar información (una tarea para la que no está

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

optimizado), le traes los documentos directamente a su escritorio. Así, aprovechas su capacidad de procesamiento de lenguaje de clase mundial sobre un conjunto de datos acotado y de alta calidad, en lugar de su mediocre capacidad de búsqueda web. En esta maniobra, el humano actúa como un curador de inteligencia de alto nivel, y el Agent como un procesador de síntesis de élite. Cada uno se enfoca exclusivamente en la tarea en la que es excepcional.

No hagas: «Investiga todo sobre el impacto de la IA en la cadena de suministro». **Haz:** «He subido 3 informes sobre IA en la cadena de suministro. Actúa como un analista de estrategia y extrae lo siguiente en una tabla comparativa:

- 1. Casos de uso validados en producción en cada informe.
- 2. Métricas de ROI o de eficiencia reportadas.
- 3. Principales barreras de adopción mencionadas.
- 4. Contradicciones o desacuerdos entre las fuentes».

Maniobra 2: El Agent como «auditor técnico programable»

El Agent no tiene opinión ni intuición, pero es excepcionalmente bueno aplicando reglas de forma consistente. Podemos usar esta característica para realizar auditorías técnicas complejas.

- La estrategia: En lugar de pedirle que "revise" un código o un documento de forma genérica, le proporcionas un *framework* de análisis explícito y le pides que lo aplique al material que le entregas.
- Por qué funciona: Esta maniobra explota su rasgo de ser «muy metódico». No le pides creatividad, sino una aplicación rigurosa de un protocolo. Es como darle una lista de control detallada y pedirle que verifique cada punto, una tarea en la que su falta de sesgos se convierte en una fortaleza. El resultado es una auditoría objetiva, repetible y escalable, libre de la fatiga o el sesgo que podría afectar a un revisor humano.

No hagas: «¿Qué te parece este código? ¿Se puede mejorar?». **Haz:** «He subido una clase de Python. Actúa como un ingeniero de *software* sénior y audita este código basándote exclusivamente en los principios de diseño SOLID. Genera un informe que liste:

- 1. Cada principio (S, O, L, I, D).
- 2. Si el código lo cumple o lo viola.
- 3. En caso de violación, proporciona un ejemplo de código refactorizado que corrija el problema».

Maniobra 3: El Agent como «arquitecto de automatización»

Esta es quizás la maniobra más sofisticada, ya que eleva al Agent de simple ejecutor de tareas a socio en el diseño de procesos de automatización.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

- La estrategia: Se utiliza un enfoque por capas. Primero, se le pide al Agent que diseñe la lógica de un flujo de trabajo en un nivel alto (pseudocódigo o diagrama).
 Una vez que el diseño es validado por el humano, se le pide que traduzca ese diseño a código ejecutable. Finalmente, se le pide que añada capas de robustez como el manejo de errores o el registro de eventos.
- Por qué funciona: Esta maniobra aprovecha al máximo su capacidad de «pensar en voz alta» (el patrón ReAct). Al separar el diseño (Thought) de la implementación (Action), te permite validar la lógica del proceso antes de escribir una sola línea de código. Es una forma de colaboración humano-máquina donde el Agent actúa como un arquitecto de sistemas que diseña los planos, y luego como un programador que los construye, todo bajo tu dirección estratégica.

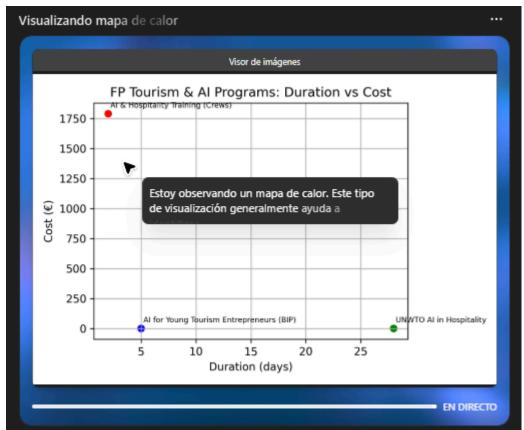
Flujo de trabajo en tres pasos:

- Tú: «Diseña un flujo de trabajo en pseudocódigo para un script que lea todos los archivos CSV de una carpeta, calcule las ventas totales de cada uno y consolide los resultados en un único archivo de Excel».
- 2. **Tú (tras validar el pseudocódigo):** «Perfecto. Ahora convierte este pseudocódigo en un *script* de Python funcional».
- 3. **Tú (tras validar el script):** «Excelente. Ahora añade un manejo de errores robusto con bloques try...except para gestionar posibles archivos corruptos o columnas faltantes, y añade un sistema de *logging* para registrar qué archivos se procesaron correctamente y cuáles fallaron».

Dominar estas maniobras es la diferencia entre usar el Agent como una calculadora avanzada y dirigirlo como un socio de análisis estratégico. Es el punto donde dejas de ser un simple piloto para convertirte en un verdadero arquitecto de la colaboración. No se trata de dominar la herramienta, sino de dominar una nueva forma de resolver problemas, donde la limitación se convierte en la principal fuente de estrategia.

5.4. El plan de vuelo: la matriz de decisión definitiva

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura



Duración y

coste de programas reales de IA aplicada al turismo: análisis generado dinámicamente por un Agent

Todo nuestro análisis, desde la ingeniería inversa de la arquitectura hasta el modelo mental del «analista superdotado», se puede condensar en una única herramienta estratégica: una matriz de decisión. Esta no es una guía de buenas prácticas genéricas; es un plan de vuelo operativo basado en la evidencia empírica que hemos recopilado. La diferencia entre un profesional que extrae un valor 10x de la herramienta y un usuario que solo encuentra frustración reside, a menudo, en consultar este plan de vuelo antes de despegar.

USA el Agent cuando tu misión cumple estas condiciones:

Característica de la Tarea	Justificación arquitectónica
Análisis exploratorio y prototipado	El ciclo rápido de ReAct (ver 4.1) y la agilidad entre herramientas (ver 6.2) lo hacen ideal para la iteración y el descubrimiento.
Los datos caben en memoria (< 7-8 GB)	La tarea respeta el límite de la «pizarra pequeña» de nuestro analista, evitando el MemoryError.
Tareas acotadas en sesiones de ~5 min	Se adapta al modelo de « <i>sprints</i> de concentración» del Agent, evitando el fin abrupto de la sesión

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

No requiere memoria entre sesiones

El flujo de trabajo es compatible con la naturaleza efímera de las micro-VM, que se reinician en cada sesión.

El resultado es un análisis o un informe

Aprovecha su fortaleza en la estructuración de información y su capacidad para «pensar en voz alta».

Puedes proporcionar un contexto claro

Le das al Agent la información que necesita para tener éxito, en lugar de pedirle que adivine un contexto que no posee.

X EVITA el Agent cuando tu misión implica alguna de estas condiciones:

Característica de la Tarea	Justificación arquitectónica
Integración con sistemas internos (CRM, ERP)	El aislamiento de red total de la «sala limpia» lo hace arquitectónicamente imposible.
Procesos ininterrumpidos de larga duración (> 5 min)	Choca directamente con el presupuesto de tiempo de la sesión, llevando a un timeout inevitable.
Necesidad de persistencia entre días	La arquitectura de micro-VM efímeras garantiza que no exista memoria a largo plazo entre sesiones.
Manejo de <i>Big Data</i> (> 7-8 GB en RAM)	Supera los límites de recursos de la micro-VM, resultando en fallos catastróficos y la degradación del sandbox.
Tareas que requieren garantías de precisión	Como todo LLM, puede «alucinar». No es fiable para tareas donde la precisión fáctica es crítica y no puede ser verificada.
Contexto empresarial muy específico y tácito	El contexto tácito es su punto ciego; no puede inferir el conocimiento que no se le proporciona explícitamente.

Esta matriz no es un juicio de valor sobre la herramienta, sino un mapa de su territorio operativo. Un piloto experto no es el que vuela la máquina más potente, sino el que conoce su mapa a la perfección. Usar esta guía no es una limitación, es la forma más alta de estrategia: alinear la tarea con la herramienta para lograr un rendimiento que, de otro modo, sería imposible.

5.5. Apéndice práctico: recetas de «prompts» y guía de campo

La teoría sin aplicación es un ejercicio incompleto. Este apéndice final traduce todo nuestro análisis en herramientas de trabajo directas. A continuación, se presentan plantillas y una guía de referencia rápida diseñadas para ser utilizadas en el día a día, permitiendo al profesional aplicar de forma sistemática las estrategias que hemos validado a lo largo de este artículo.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

Recetario de «prompts» para el Agent

Un buen *prompt* no es una simple pregunta, es un pliego de condiciones. Es la forma de comunicarle a nuestro «analista superdotado» no solo qué hacer, sino cómo hacerlo, respetando sus reglas. La siguiente plantilla estructura la comunicación para maximizar la eficiencia y minimizar los errores. **Plantilla Maestra de Prompt**

ROL Y CONTEXTO

Actúa como [rol experto, ej: "un analista financiero sénior especializado en Python"]. Tu objetivo es [objetivo final, ej: "evaluar la rentabilidad de estos productos"].

CONOCIMIENTO DE LÍMITES (Le recordamos sus propias reglas)

Eres consciente de que operas en un entorno con ~7-8 GB de RAM y una sesión de ~5 minutos. Gestiona tus recursos de forma eficiente.

INSUMOS (Qué le proporcionamos)

Voy a subir los siguientes archivos:

- `[nombre archivo 1.csv]`: [breve descripción del contenido]
- `[nombre_archivo_2.xlsx]`: [breve descripción del contenido]

TAREA (Pasos claros y numerados)

Ejecuta las siguientes acciones en orden:

- 1. [Acción específica y clara 1]
- 2. [Acción específica y clara 2]
- 3. [Acción específica y clara 3]

FORMATO DE SALIDA (Le decimos cómo gueremos la respuesta)

El resultado final debe ser [formato deseado, ej: "un único gráfico de barras y un resumen ejecutivo en 3 puntos clave"]. No muestres código ni resultados intermedios a menos que sean esenciales para el siguiente paso.

Ejemplo de Aplicación: Análisis de Datos Seguro

ROL Y CONTEXTO

Actúa como un analista de datos experto en Python y Pandas. Tu objetivo es identificar los productos menos rentables de este dataset de ventas.

CONOCIMIENTO DE LÍMITES

Eres consciente de que operas en un entorno con ~7-8 GB de RAM y una sesión de ~5 minutos. Al cargar el CSV, por favor, infiere y utiliza los `dtypes` más eficientes posibles para optimizar el uso de memoria.

INSUMOS

Voy a subir el archivo `ventas_trimestre.csv`, que contiene datos de ventas por producto.

TAREA

Ejecuta las siguientes acciones en orden:

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

- 1. Carga el archivo `ventas_trimestre.csv`.
- 2. Calcula una nueva columna llamada 'Margen' restando 'Coste_Producto' de 'Precio_Venta'.
- 3. Agrupa por 'Nombre_Producto' y calcula la suma total del 'Margen'.
- 4. Identifica los 5 productos con el menor margen total.

FORMATO DE SALIDA

El resultado final debe ser una tabla que muestre únicamente los 5 productos con menor margen y su margen total acumulado. No generes gráficos.

Guía de campo para dominar al Agent

Esta es la chuleta, la lista de control de bolsillo. Imprímela, pégala en tu monitor, pero no la ignores. Es el resumen destilado de todo nuestro análisis. **Límites de Oro (Reglas no negociables)**

- Memoria: ~7-8 GB por proceso. Si un archivo expandido en RAM supera esto, fallará.
- **Tiempo:** ~5 minutos por sesión. Los procesos ininterrumpidos más largos fallarán.
- Red: Aislamiento total. El intérprete de Python no puede acceder a internet.

Estrategias Ganadoras (Tus tres maniobras clave)

- Descomponer Tareas: Divide problemas grandes en una cadena de sprints de 5 minutos.
- **Usa siempre dtypes específicos y categoricals.**Procesa en chunks si es necesario.
- **Ser Específico:** No pidas "analiza". Pide "calcula X, agrupa por Y, y visualiza en un gráfico Z".

Anatomía del Prompt Perfecto

- 1. Rol: ¿Quién quieres que sea?
- 2. Contexto: ¿Cuál es el objetivo final?
- 3. Límites: Demuéstrale que conoces sus reglas.
- 4. **Tarea:** Pasos numerados y claros.
- 5. Formato: ¿Cómo quieres la respuesta?

Señal de Alarma Universal Si alguna vez te encuentras con un MemoryError, no pierdas tiempo intentando depurar el siguiente paso. La probabilidad de que el entorno o sandbox haya quedado en un estado "degradado" es muy alta. La acción correcta e inmediata es siempre la misma: inicia una nueva conversación.

Nota informativa: Todos los tests y observaciones de este documento corresponden a la versión pública del sistema de Agents de ChatGPT disponible entre marzo y julio de 2025. Aunque el sistema evoluciona

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

continuamente, los principios arquitectónicos aquí documentados parecen seguir vigentes en la fecha de publicación.

6. Epílogo: El arte de trabajar con inteligencia limitada

Más allá de la ingeniería inversa: ¿qué hemos aprendido realmente?

Después de semanas de análisis sistemático, pruebas empíricas y reconstrucción arquitectónica, emerge una verdad incómoda pero liberadora: el Agent de ChatGPT no es la herramienta revolucionaria que prometen los titulares, ni el juguete limitado que sugieren los escépticos. Es algo más sutil y, paradójicamente, más valioso: un espejo de nuestras propias decisiones sobre cómo queremos que la IA se integre en el trabajo profesional. Lo que hemos descubierto no es solo una arquitectura técnica. Es una filosofía de diseño que prioriza la utilidad predecible sobre la brillantez impredecible, la seguridad aburrida sobre la innovación arriesgada, y la economía sostenible sobre las capacidades ilimitadas. Y en esas decisiones aparentemente conservadoras, hay lecciones profundas sobre el futuro del trabajo aumentado por IA.

6.1. El framework estratégico: Cómo maximizar valor real

La paradoja de las limitaciones productivas

Una de las revelaciones más contraintuitivas de este análisis es que las limitaciones del Agent no son errores a pesar de los cuales la herramienta es útil. Son características que la hacen útil *precisamente porque* fuerzan un tipo específico de interacción. La limitación clave no es un simple cronómetro por operación, sino un **presupuesto total por sesión de aproximadamente cinco minutos de tiempo real** (*wall-clock time*). Esta restricción redefine el concepto de productividad: no se trata solo de la velocidad de ejecución del código, sino de la eficiencia de todo el flujo de trabajo. Esta limitación te obliga a evolucionar de ser un simple usuario a ser un arquitecto de soluciones: un profesional que no solo pide tareas, sino que diseña flujos de trabajo eficientes dentro de un entorno con recursos definidos.

• Sin un presupuesto de sesión estricto:

 Usuario: Podría pedir veinte gráficos intermedios para visualizar cada paso, usar print en cada iteración de un bucle y generar una gran cantidad de output intermedio.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

 Resultado: Un análisis potencialmente útil, pero un proceso "ruidoso", ineficiente y que consume recursos (tiempo de renderizado, latencia de red) de forma desmedida.

Con un presupuesto de sesión de ~5 minutos:

- Usuario: Se ve forzado a pensar: "¿Realmente necesito este gráfico ahora o puedo esperar al resultado final? ¿Puedo acumular los resultados en una lista y visualizarlos una sola vez?".
- Resultado: Un flujo de trabajo más deliberado y eficiente. El usuario aprende a minimizar las operaciones "caras" en tiempo real, lo que a menudo conduce a un código más limpio y a un análisis más enfocado.

La limitación no te obliga a hacer preguntas más pequeñas; te obliga a diseñar **procesos más inteligentes**. No es una barrera; es un *framework* de eficiencia impuesto por la arquitectura.

Protocolo de trabajo para maximizar la eficiencia

Basándome en los patrones de éxito documentados, propongo el siguiente protocolo de trabajo:

1. Pre-procesamiento humano (5 minutos)

- o Define el resultado específico deseado.
- o Identifica los datos mínimos necesarios.
- Formula la pregunta más estrecha posible y el flujo de trabajo para responderla.
- Anticipa el formato de salida ideal para minimizar el *output* intermedio.

2. Ejecución con el Agent (tareas de 5-15 minutos)

- Carga solo los datos necesarios.
- Usa *prompts* que especifican el formato de salida.
- Intervén temprano si se desvía.
- Guarda los resultados importantes y, si es necesario, inicia una nueva sesión para continuar.

3. Post-procesamiento y aplicación (10 minutos)

- Valida los resultados contra el sentido común.
- o Extrae los insights accionables.
- Documenta el proceso para su reproducibilidad.
- o Implementa en sistemas reales.

La proporción 20/80 invertida: 20 % del tiempo con el Agent, 80 % pensando antes y después. Esto no es ineficiencia, sino la clave de la productividad. Esta inversión de esfuerzo inicial minimiza el tiempo perdido en bucles de prueba y error, ya que el flujo de trabajo se diseña desde el principio para operar dentro de las fortalezas de la arquitectura del Agent y evitar sus limitaciones conocidas.

Estrategias específicas por caso de uso

Para análisis de datos:

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

- No hagas: "Analiza este dataset"
- Haz: "Usando el dataset cargado: 1. Identifica los 3 drivers principales de variación en ventas. 2. Cuantifica su impacto relativo. 3. Genera un único gráfico de Pareto que resuma los hallazgos. 4. Resume en 3 bullets ejecutivos."

Para generación de código:

- No hagas: "Crea una aplicación de gestión de inventario"
- Haz: "Genera una clase Python que: Trackee entradas/salidas de productos. - Alerte cuando el stock sea inferior a un umbral. -Exporte un reporte diario en formato CSV. - Incluya tests unitarios básicos."

Para investigación:

- No hagas: "Investiga sobre blockchain en logística"
- Haz: "Tengo estos 3 artículos sobre blockchain en logística.
 Extrae: Casos de uso validados en producción. Métricas de ROI reportadas. Barreras de adopción principales. Formato: tabla comparativa."

Este enfoque estructurado no solo produce mejores resultados, sino que optimiza el uso de los recursos limitados de la sesión (tiempo y memoria), alineando la tarea directamente con las capacidades de la arquitectura del Agent.

6.2 Una crítica honesta: lo que el Agent revela sobre nuestra relación con la IA

El Agent sufre de lo que llamo la "demo-cracia": la tiranía de las demostraciones impresionantes. En manos de sus creadores o de un usuario experto que conoce sus secretos, puede parecer magia. En la realidad del día a día, sin embargo, puede resultar frustrante hasta que se entienden sus límites. Basta revisar algunas demos promocionales de asistentes de código o análisis financiero automático para ver cómo la promesa de una IA 'que lo hace todo sola' eclipsa las condiciones reales de uso. Este artículo busca precisamente cerrar esa brecha. Este no es un fallo exclusivo del Agent. Es un patrón que define a toda una generación de herramientas de IA:

- 1. Se desarrolla una capacidad asombrosa.
- 2. Se demuestra en condiciones ideales, a menudo con problemas perfectamente acotados.
- 3. Se lanza al público general sin una documentación exhaustiva de sus limitaciones reales.
- 4. Los usuarios experimentan una brecha entre la expectativa y la realidad, lo que genera frustración.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

5. Finalmente, surge un conocimiento comunitario (como el que busca destilar este artículo) que descubre y comparte las estrategias para trabajar eficazmente dentro de esas limitaciones.

La pregunta crítica es: ¿es este ciclo inevitable, o podríamos aspirar a una mayor transparencia desde el principio?

Nota informativa: A este fenómeno lo llamamos 'demo-cracia': el dominio de las demostraciones perfectamente preparadas sobre la percepción real del sistema. Una tiranía del 'modo demo' que disfraza las restricciones estructurales como si fueran elecciones del usuario."

La ilusión de la automatización completa

El Agent, con su capacidad para ejecutar tareas, perpetúa una fantasía peligrosa: que la IA está lista para automatizar completamente el trabajo intelectual. La realidad que hemos descubierto a través de este análisis es mucho más matizada y colaborativa:

• Lo que promete: Automatización

- Lo que entrega: Aumentación
- Lo que realmente necesitamos: Colaboración

El Agent es brillante ejecutando instrucciones claras, pero su propia arquitectura le impide ir más allá. No puede:

- Entender el contexto no articulado de tu organización, pues su memoria es efímera y se limita a la sesión actual.
- Tomar decisiones con implicaciones éticas o estratégicas, ya que carece de un verdadero entendimiento del mundo.
- Innovar más allá de recombinar patrones conocidos, pues su creatividad se basa en los datos con los que fue entrenado.
- Asumir la responsabilidad final por sus resultados, una tarea que, por diseño, siempre recae en el supervisor humano.

Y esto no es una limitación técnica temporal que se resolverá en la próxima actualización. Es una característica fundamental de su arquitectura.

El sesgo hacia lo medible

Quizás la crítica más profunda que se le puede hacer al Agent es que, por su propia naturaleza, nos empuja a valorar únicamente lo que se puede medir, cuantificar y estructurar. Es excepcionalmente bueno analizando un DataFrame, generando un gráfico o refactorizando un bloque de código. Pero el trabajo humano más valioso a menudo reside en los espacios intangibles que rodean a esas tareas:

- La intuición de que algo "no cuadra" en los datos, aunque las cifras sean correctas.
- El **conocimiento tácito** de por qué un cliente en particular es estratégico, más allá de sus cifras de ventas.
- La **creatividad** para conectar ideas de dominios aparentemente dispares.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

• El **juicio** para saber cuándo es necesario romper las reglas o cuestionar la premisa inicial.

El peligro es sutil, pero real: que empecemos a optimizar nuestro trabajo y nuestros problemas para que se ajusten a lo que la herramienta hace bien, devaluando progresivamente las habilidades humanas críticas que la IA, por ahora, no puede replicar. El reto no es solo usar la herramienta, sino preservar y cultivar nuestra propia inteligencia en el proceso.

6.3 El futuro del trabajo con agentes: proyecciones basadas en la arquitectura

Después de este profundo análisis, podemos trazar un mapa no solo del presente, sino también del futuro más probable de estas herramientas. Y este mapa se dibuja siguiendo las líneas maestras de su arquitectura fundamental.

Lo que probablemente no cambiará

Basándonos en las decisiones arquitectónicas centrales que hemos validado, algunas limitaciones son tan intrínsecas al modelo de negocio y de seguridad que persistirán en el futuro previsible:

- 1. El modelo económico de recursos limitados. Mientras el servicio se ofrezca a una escala masiva, a un precio accesible y permitiendo la ejecución de código no confiable, los límites de recursos (RAM, tiempo de sesión) seguirán existiendo. Son la única forma de garantizar la viabilidad económica.
- 2. La filosofía de aislamiento. El sandboxing extremo mediante microVMs no es una paranoia temporal, sino la única forma probada de mitigar los riesgos de seguridad a gran escala. Esto significa que la integración directa con sistemas empresariales internos siempre requerirá puentes y conectores construidos por los propios usuarios o por terceros.
- La naturaleza efímera. La persistencia real de datos entre sesiones requeriría un rediseño fundamental que comprometería la seguridad o dispararía los costes. No debemos esperar que el Agent "recuerde" nuestras conversaciones de un día para otro en el futuro cercano.

Lo que probablemente mejore

La evolución del sistema es constante. De hecho, algunas de las mejoras incrementales que se podían prever hace poco, como la ampliación de los tiempos de ejecución hacia un presupuesto de sesión de casi 5 minutos, ya son una realidad. Mirando hacia el futuro, las mejoras probablemente seguirán esta línea de ajustes sofisticados, no de transformaciones radicales:

• Inteligencia del modelo base: Los LLMs seguirán mejorando, lo que se traducirá en una mejor comprensión de instrucciones ambiguas, una generación de código más eficiente y una menor tasa de errores o "alucinaciones" en los análisis.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

- Herramientas más especializadas: Es probable que el arsenal de herramientas se expanda para incluir capacidades de análisis más específicas (ej. bioinformática, análisis geoespacial) o conectores seguros a APIs populares (con autenticación gestionada).
- Niveles de servicio diferenciados: En lugar de una única mejora para todos, el siguiente paso lógico podría ser la introducción de planes "Pro" o "Enterprise" que ofrezcan, por un precio mayor, sesiones con presupuestos de tiempo y RAM más generosos (ej. 15 minutos de sesión, 32GB de RAM).
- **Persistencia limitada:** El primer paso hacia la persistencia podría no ser un almacenamiento permanente, sino una "caché de sesión" que permita reanudar un trabajo durante unas pocas horas antes de la purga final.

Aun así, es crucial entender que estos serán **ajustes sobre la arquitectura existente**, **no transformaciones** de la misma.

El verdadero cambio: de herramienta a plataforma

La evolución más significativa no vendrá de que el Agent se vuelva ilimitado, sino de su transformación de una herramienta cerrada a una **plataforma abierta**. El futuro más probable es un ecosistema donde:

- Los desarrolladores construyan conectores específicos y seguros para sus propios sistemas internos.
- Las **empresas** puedan desplegar versiones privadas del Agent, con límites y herramientas personalizadas para sus necesidades.
- Surjan *marketplaces* de *prompts* y flujos de trabajo optimizados para industrias y casos de uso específicos (finanzas, legal, investigación científica...).

El valor ya no residirá únicamente en la herramienta, sino en el ecosistema de soluciones que se construya sobre ella. Estas restricciones no tienen porque ser errores ni de diseño ni de capacidad. Pero parece que probablemente son, al menos en el horizonte actual, los pilares sobre los que se construye la escalabilidad, seguridad y viabilidad económica del Agent. Comprender esto no es resignarse: es el primer paso hacia una estrategia realista.

Referencias

Referencias académicas

- 1. Yao, S., Zhao, J., Yu, D., Gao, J., Chen, D., & Hajishirzi, H. (2022). *ReAct:* Synergizing reasoning and acting in language models (arXiv:2210.03629). arXiv. https://doi.org/10.48550/arXiv.2210.03629
- Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., & Iwasawa, Y. (2023). Large language models are zero-shot reasoners (arXiv:2201.11903). arXiv. https://doi.org/10.48550/arXiv.2201.11903

Referencias internas y base documental

Este informe se apoya en un conjunto de evidencias técnicas generadas por el autor entre junio y julio de 2025 como parte de un proceso sistemático de análisis, testeo y auditoría operativa del sistema de Agents de ChatGPT.

Los documentos están organizados como anexos internos (A1-A9), clasificados por tipo de prueba, dimensión evaluada y metodología aplicada. Estos documentos no son públicos, ya que forman parte de una línea de investigación activa y se encuentran bajo control de versiones interno.

Las evidencias fueron generadas directamente por el autor, sin mediación de terceros, en condiciones reproducibles mediante acceso a ChatGPT Plus (con funcionalidades habilitadas de Agents, Advanced Data Analysis y archivos). Los experimentos fueron conducidos de forma controlada y documentados íntegramente con registro de prompts, resultados, tiempos de respuesta, errores del sistema y observaciones cruzadas.

Nabla de Anexos (documentación interna no publicada)

Códig o	Título técnico	Tipo de evidencia	Descripción breve
A1	TestEvidencial1_Memoria	Prueba empírica funcional	Persistencia de memoria factual y lógica del agent.
A2	TestEvidencial2_TareaVM	Evaluación de coordinación interna	Agente frente a tareas cruzadas con instrucciones múltiples.
A3	TestEvidencial3_CodeInterpreter	Prueba de integración funcional	Comportamiento real del Code Interpreter dentro de un Agent.
A4	TestEvidencial4_Tempos	Medición de latencias y orden lógico	Análisis de tiempos y secuencia de ejecución.
A5	TestEvidencial5_LimitesCSV	Prueba de umbrales técnicos	Comportamiento ante grandes volúmenes y parsing de archivos.
A6	TestEvidencial6_autocorrecciones	Observación de autoajuste del sistema	Capacidad de los Agents para detectar y corregir fallos propios.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

A7	Conclusiones Evidencia (trabajo experimental JL)	Síntesis operativa transversal	Consolidación de patrones, límites y anomalías detectadas.
A8	Auditoría Forense Agente de ChatGPT	Registro técnico-empíric o completo	Log detallado de sesiones, errores y comportamiento no esperado.
A9	Análisis Agente ChatGPT: Arquitectura, Límites, Lógica	Inteligencia técnica aplicada	Inferencias sobre arquitectura interna, fallos y estructura.

Nota: Los documentos anteriores no están disponibles públicamente. Forman parte de una investigación en curso y están protegidos como trabajo original del autor. Su existencia y estructura se documentan aquí para reforzar la trazabilidad metodológica del informe y permitir replicabilidad técnica por parte de terceros con acceso equivalente al sistema evaluado.

NOTA LEGAL Y METODOLÓGICA

Naturaleza del análisis y limitaciones epistémicas

Este documento constituye un **ejercicio de ingeniería inversa mediante observación empírica** del sistema ChatGPT Agents, realizado exclusivamente a través de interfaces públicas y dentro de los términos de servicio aplicables. El análisis se fundamenta en metodologías establecidas de inferencia arquitectónica mediante *behavioral profiling* y *performance boundary analysis*.

Declaración de Independencia y Metodología

AUSENCIA DE INFORMACIÓN PRIVILEGIADA:

- El autor NO posee acceso a código fuente, documentación técnica interna, especificaciones de diseño, o cualquier información propietaria de OpenAI
- NO se han establecido comunicaciones técnicas con personal de OpenAl sobre aspectos arquitectónicos internos
- NO se han violado términos de servicio, acuerdos de confidencialidad, o derechos de propiedad intelectual

METODOLOGÍA DE INFERENCIA: El análisis emplea técnicas reconocidas de ingeniería inversa:

- Black-box testing sistemático para mapear límites operativos
- Análisis de patrones de comportamiento bajo condiciones controladas
- Inferencia arquitectónica por eliminación de alternativas técnicamente implausibles
- Correlación con tecnologías de código abierto documentadas públicamente
- Validación cruzada mediante réplica experimental

SISTEMA DE CONFIANZA CALIBRADA: Todas las afirmaciones técnicas están clasificadas según evidencia disponible:

- Alta confianza: Comportamiento directamente observable y replicable
- Media confianza: Inferencia lógica basada en patrones convergentes
- Baja confianza: Especulación técnica fundamentada en principios de ingeniería

Limitaciones Técnicas del Análisis

LIMITACIONES TEMPORALES:

- Análisis realizado entre marzo-julio 2025 sobre versión pública del sistema
- Las inferencias arquitectónicas pueden no reflejar implementaciones internas actuales

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

Evolución continua del sistema puede invalidar observaciones específicas

LIMITACIONES METODOLÓGICAS:

- Las inferencias sobre micro-virtualización (Firecracker, gVisor) se basan en convergencia de evidencia indirecta
- Los límites cuantificados (~7-8GB RAM, ~5min sesión) son observaciones empíricas, no especificaciones confirmadas
- El análisis del patrón ReAct se fundamenta en comportamiento observable, no en confirmación de implementación

LIMITACIONES DE ALCANCE:

- El análisis se centra en arquitectura de ejecución, no en aspectos de entrenamiento del modelo base
- Las proyecciones sobre decisiones de producto se basan en análisis económico-técnico, no en información estratégica interna
- La caracterización de herramientas integradas refleja funcionalidad observada, no diseño interno

Limitación de Responsabilidad Técnica

EXACTITUD DE INFERENCIAS: El autor no garantiza la exactitud de las inferencias arquitectónicas realizadas. Los lectores técnicos deben:

- Validar independientemente las observaciones mediante experimentación propia
- Considerar las inferencias como hipótesis de trabajo, no como especificaciones técnicas definitivas
- Aplicar sus propios criterios de confianza basados en evidencia disponible

APLICACIÓN PRÁCTICA: Las estrategias operativas sugeridas se basan en patrones observados que pueden evolucionar. Los usuarios deben:

- Adaptar recomendaciones a sus casos de uso específicos
- Monitorear cambios en el comportamiento del sistema
- Mantener flexibilidad operativa ante posibles modificaciones arquitectónicas

DECISIONES EMPRESARIALES: Este análisis tiene carácter técnico-educativo. Las decisiones de adopción, integración o inversión basadas en este documento son responsabilidad exclusiva del lector.

Declaración de Autoría y Derechos

TRABAJO ORIGINAL: Este análisis constituye trabajo intelectual original del autor, generado mediante investigación independiente y experimentación directa.

Qué hacen realmente los Agents, cómo funcionan bajo el capó y por qué importa entender su arquitectura

PROTECCIÓN DE LA INVESTIGACIÓN: Los documentos de evidencia empírica (Anex A1-A9) permanecen como documentación técnica privada del autor, disponibles para validación académica bajo acuerdo de confidencialidad.	os
Documento de investigación técnica independiente - v2.3-beta #AGT20250727	

Este documento se distribuye bajo el entendimiento de que constituye análisis técnico independiente basado en observación empírica. Los lectores asumen la responsabilidad de validar la aplicabilidad de las conclusiones a sus contextos específicos.