

¿Qué es Gemini 2.5 Computer Use?

Gemini 2.5 Computer Use es un modelo especializado de inteligencia artificial desarrollado por Google DeepMind, construido sobre las capacidades de comprensión visual y razonamiento de Gemini 2.5 Pro. Este modelo representa un avance significativo en la automatización de tareas digitales, ya que permite a los agentes de IA interactuar directamente con interfaces gráficas de usuario (GUI) de manera similar a como lo haría un humano.

La característica más distintiva de este modelo es su capacidad para **"ver"** la pantalla de una computadora y **"actuar"** generando acciones específicas de interfaz de usuario como clics del mouse, entradas de teclado, desplazamiento y arrastre de elementos. A diferencia de las APIs estructuradas tradicionales o el web scraping convencional, el modelo Gemini 2.5 Computer Use observa y analiza capturas de pantalla para determinar qué acciones tomar, imitando la forma en que los humanos interactúan con las aplicaciones.

El modelo está **optimizado principalmente para navegadores web**, aunque también demuestra un rendimiento prometedor en tareas de control de interfaces móviles. Sin embargo, aún no está optimizado para el control a nivel del sistema operativo de escritorio. Esta capacidad cierra una brecha crítica en la automatización de tareas digitales que tradicionalmente requerían intervención manual, como completar formularios complejos, navegar por sitios web con autenticación o manipular elementos interactivos como menús desplegables y filtros.

Cómo Opera el Modelo

El funcionamiento del modelo Gemini 2.5 Computer Use se basa en un **mecanismo de bucle iterativo** que opera a través de la herramienta `computer_use` en la API de Gemini. Este proceso cíclico consta de cuatro pasos fundamentales que se repiten hasta completar la tarea asignada:

1. Envío de Solicitud al Modelo

El proceso comienza cuando los desarrolladores proporcionan al modelo varias entradas clave: la solicitud del usuario con el objetivo a realizar, una captura de pantalla del entorno actual, un historial de acciones recientes, y opcionalmente pueden especificar exclusiones de acciones o incluir funciones personalizadas.

2. Análisis y Generación de Respuesta

El modelo procesa estas entradas utilizando sus capacidades de comprensión visual avanzada para analizar la captura de pantalla. Luego genera una respuesta, típicamente en forma de una llamada a función (`function_call`) que representa una acción específica de UI como hacer clic, escribir texto, desplazarse o arrastrar elementos. Para acciones sensibles

como realizar una compra, el modelo puede incluir una solicitud de confirmación del usuario final.

3. Ejecución de Acciones

El código del lado del cliente recibe estas instrucciones y ejecuta las acciones correspondientes en el entorno real del navegador. Esta separación entre la propuesta del modelo y la ejecución del cliente es fundamental para mantener la seguridad y el control sobre las operaciones.

4. Captura del Nuevo Estado

Después de ejecutar cada acción, se captura una nueva captura de pantalla del GUI junto con la URL actualizada, enviándose de vuelta al modelo como una ``function_response``. Esto reinicia el bucle, permitiendo al modelo evaluar el resultado de la acción anterior y determinar el siguiente paso.

Este proceso iterativo continúa hasta que se completa la tarea, ocurre un error, o la interacción es terminada por una respuesta de seguridad o decisión del usuario. El modelo utiliza un ****sistema de coordenadas estandarizado de 1000x1000**** que se escala automáticamente al tamaño real de la pantalla, siendo recomendable usar una resolución de 1440x900 píxeles para obtener mejores resultados.

instalación y Configuración

La implementación del modelo Gemini 2.5 Computer Use requiere varios pasos de configuración. Te explicaré tanto la instalación local como las diferentes formas de acceso:

Requisitos Previos

Para implementar el modelo necesitarás configurar los siguientes elementos:

- Python 3.10 o superior con pip
- Un entorno virtual (recomendado usando Anaconda o venv)
- SDK de Google GenAI para Python (``google-genai``)
- Playwright para el control del navegador
- Una clave API de Google AI Studio

Instalación Paso a Paso

Paso 1: Configuración del Entorno

Primero, debes crear y activar un entorno virtual. Puedes usar Anaconda como se muestra en los tutoriales:

```
```bash
conda create -n computer_use python=3.12
```

```
conda activate computer_use
```
```

O alternativamente con venv:

```
```bash
python -m venv .venv
source .venv/bin/activate # En Windows: .venv\Scripts\activate
```
```

Paso 2: Instalación de Dependencias

Instala las librerías necesarias:

```
```bash
pip install google-genai playwright
playwright install chromium
```
```

Si estás usando el repositorio oficial de GitHub, puedes clonar el proyecto:

```
```bash
git clone [URL del repositorio de computer-use-preview]
cd computer-use-preview
pip install -r requirements.txt
```
```

Paso 3: Configuración de la Clave API

Obtén una clave API desde Google AI Studio en <https://aistudio.google.com> y configúrala como variable de entorno:

```
```bash
macOS/Linux
export GEMINI_API_KEY="tu-clave-api"

Windows PowerShell
$Env:GEMINI_API_KEY="tu-clave-api"
```
```

Paso 4: Inicialización del Cliente

Configura el cliente con el modelo específico de Computer Use:

```
```python
from google import genai
from google.genai import types

client = genai.Client()
```

```
Configuración con la herramienta Computer Use
generate_content_config = genai.types.GenerateContentConfig(
 tools=[
 types.Tool(
 computer_use=types.ComputerUse(
 environment=types.Environment.ENVIRONMENT_BROWSER,
)
),
]
)
...

```

### Paso 5: Uso del Modelo

El modelo debe especificarse como `gemini-2.5-computer-use-preview-10-2025`:

```
```python
response = client.models.generate_content(
    model="gemini-2.5-computer-use-preview-10-2025",
    contents=contents,
    config=generate_content_config,
)
...

```

Opciones de Acceso Alternativas

****Browserbase Demo****: Para probar el modelo sin instalación local, puedes acceder a la demostración en <https://gemini.browserbase.com>, aunque tiene un límite de pruebas gratuitas.

****Google AI Studio y Vertex AI****: Los desarrolladores pueden acceder directamente a través de estas plataformas con sus credenciales existentes, siendo Vertex AI más apropiado para implementaciones empresariales.

Capacidades y Características Principales

El modelo Gemini 2.5 Computer Use ofrece un conjunto completo de ****acciones de interfaz de usuario**** que incluyen:

- Hacer clic en elementos específicos
- Escribir texto en campos de entrada
- Desplazarse vertical y horizontalmente
- Arrastrar y soltar elementos
- Pasar el cursor sobre elementos (hover)

- Esperar tiempos específicos para carga de contenido dinámico

En cuanto al rendimiento, el modelo ha demostrado resultados ****líderes en múltiples benchmarks**** de la industria. Por ejemplo, en el entorno de prueba Browserbase para Online-Mind2Web, logra más del 70% de precisión con una latencia de aproximadamente 225 segundos, superando a los competidores en calidad mientras mantiene tiempos de procesamiento reducidos.

Aplicaciones Prácticas y Casos de Uso

Las capacidades del modelo se ilustran mejor a través de ejemplos concretos de implementación:

Gestión de CRM Automatizada: En una demostración, el agente accede a una página de registro de cuidado de mascotas, extrae detalles de mascotas con residencia en California y las integra automáticamente en un sistema CRM de spa, programando posteriormente citas de seguimiento con especialistas específicos.

Organización de Información Visual: El modelo puede reorganizar tableros de notas adhesivas digitales, clasificando y arrastrando notas a secciones predefinidas de manera autónoma, demostrando su capacidad de comprensión espacial y manipulación de elementos.

Automatización de Flujos de Trabajo: Los equipos de desarrollo reportan tasas de éxito superiores al 60% en la rehabilitación de ejecuciones estancadas, reduciendo los tiempos de resolución de problemas de días a minutos.

Consideraciones de Seguridad

Google ha integrado múltiples **capas de seguridad** directamente en el modelo:

- Mecanismos de rechazo incorporados para acciones potencialmente dañinas
- Solicitudes de confirmación obligatorias para operaciones sensibles como transacciones financieras
- Evaluaciones de seguridad paso a paso antes de la ejecución de acciones
- Protección contra inyección de prompts y otros vectores de ataque

Los desarrolladores deben implementar el modelo en **entornos seguros y controlados**, preferiblemente usando máquinas virtuales aisladas, contenedores o perfiles de navegador con permisos limitados.

Limitaciones Actuales

Es importante conocer las limitaciones del modelo en su estado actual de preview:

- Alcance limitado al navegador:** No está optimizado para control a nivel del sistema operativo de escritorio
- **Restricciones de cumplimiento:** No debe usarse para aceptar automáticamente términos de servicio, políticas de privacidad o resolver CAPTCHAs
- **Variabilidad en preview:** Como característica en desarrollo, puede experimentar cambios en la API y variaciones en el rendimiento

El modelo Gemini 2.5 Computer Use representa un avance significativo en la automatización inteligente de interfaces, ofreciendo a los desarrolladores una herramienta poderosa para construir agentes capaces de interactuar con aplicaciones web de manera natural y adaptativa.