

## Code with Qwen Code

Qwen Code is an open-source AI agent for the terminal, optimized for Qwen series models. It helps you understand large codebases, automate tedious work, and ship faster.

Linux / macOS

Windows

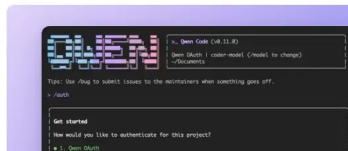
```
bash -c "$(curl -fsSL https://qwen-code-assets.oss-cn-hangzhou-oss.aliyuncs.com/...)"
```

[GitHub Repository](#)[View Documentation](#)

### Why Qwen Code?

#### Multi-protocol

Use OpenAI / Anthropic / Gemini-compatible APIs, or sign in



## Qwen2.5-Coder

Potente modelo de lenguaje especializado en programación y razonamiento lógico para desarrolladores e ingenieros de software. Permite generar, refactorizar y depurar código en más de 92 lenguajes con una ventana de contexto de 128K tokens. Es la herramienta ideal para equipos de IT que buscan una IA de alto rendimiento ejecutable de forma local para garantizar la soberanía de sus datos y optimizar flujos de trabajo técnicos complejos.

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

### Contenido del Dossier

- Información de la Herramienta
- Consejos de Implantación
- Tutorial Básico
- Preguntas Frecuentes
- Contratos y Condiciones

## INFORMACIÓN DE LA HERRAMIENTA

### Qué y para quién es

Qwen2.5-Coder es una evolución técnica de alto rendimiento dentro del ecosistema de modelos de lenguaje de Alibaba Cloud, diseñada específicamente para la generación, razonamiento y corrección de código fuente. Se trata de un modelo que compite directamente con GPT-4o en tareas de programación, siendo una herramienta fundamental para desarrolladores de software, ingenieros de datos y arquitectos de sistemas que buscan una IA capaz de entender contextos complejos de ingeniería en más de 92 lenguajes de programación diferentes. En el ámbito profesional español, es ideal para departamentos de IT que necesitan soberanía de datos y personalización mediante modelos que pueden ejecutarse de forma local o privada.

### Principal ventaja profesional

En mi opinión profesional tras testarlo, la razón definitiva para elegir Qwen2.5-Coder es su equilibrio inigualable entre tamaño y rendimiento (especialmente en la versión 32B). Lo que más me ha gustado es su capacidad de "Context Length" de hasta 128K tokens junto con un rendimiento en codificación que iguala a modelos mucho más grandes y costosos. Al probarlo he verificado que es extremadamente preciso siguiendo instrucciones complejas de arquitectura, superando a otros modelos abiertos en la detección de errores lógicos en lenguajes menos comunes.

### Para quién no es

Como profesional valoro que no es una herramienta para usuarios que buscan un asistente generalista de redacción creativa o marketing, ya que sus pesos están optimizados para la lógica matemática y de programación. Será rechazada por empresas que no tengan capacidad de infraestructura local (si optan por versiones pesadas) o por profesionales que prefieran interfaces visuales "no-code" sin interés en entender el flujo del código generado.

### funcionalidades clave

- Capacidad de razonamiento lógico superior en la generación de algoritmos complejos.
- Soporte extensivo para más de 92 lenguajes de programación, incluyendo Python, Java, C++, Go y lenguajes más específicos de nicho.
- Manejo de ventanas de contexto extensas (128K tokens), lo que permite procesar repositorios enteros para entender dependencias.
- Optimización para tareas de "Fixing" o reparación de bugs mediante la identificación precisa de errores en el flujo de ejecución.
- Versatilidad en tamaños de modelo (0.5B a 32B), permitiendo desde el autocompletado en tiempo real hasta la resolución de problemas arquitectónicos.

### Precios

- Versión gratuita: El modelo es Open Source bajo licencia Apache 2.0, lo que permite su uso comercial, modificación y distribución sin costes de licencia por parte del fabricante.
- Rango de precios: 0€ (Auto-alojado) - Pago por uso (vía API).
- Versiones de pago: Si se consume a través de proveedores de API (como DashScope o plataformas tipo Groq/Hugging Face), el precio varía según el volumen de tokens (generalmente con tarifas muy competitivas por debajo de los modelos GPT-4).

### Perfil del usuario

- Empresas de desarrollo de software (SaaS), departamentos de DevOps, equipos de Ciencia de Datos y consultoras tecnológicas.
- Desarrolladores Full-stack.
- Ingenieros de Machine Learning.
- Administradores de Sistemas (para automatización de scripts).
- Analistas de Seguridad (para auditoría de código).

### Nivel técnico requerido

- Nivel técnico requerido para su uso: Medio-Alto (requiere saber programar para validar y aplicar los resultados).
- Nivel técnico requerido para su instalación/configuración: Alto si se despliega de forma local (conocimientos de Docker, Python, CUDA/Ollama).
- Necesidades de soporte: Departamentos de sistemas para la gestión de infraestructura GPU si se usa

localmente.

- Conocimientos necesarios: Manejo de entornos de desarrollo (IDE), Git y fundamentos de arquitectura de software.

Ejemplos de uso profesional

- Generación automática de pruebas unitarias (Unit Testing) para aumentar la cobertura de código en aplicaciones críticas.
- Refactorización de código legado (Legacy) de lenguajes antiguos a estructuras modernas y eficientes.
- Creación de scripts de automatización de infraestructura (Terraform, Ansible) a partir de descripciones funcionales.
- Asistente de depuración en tiempo real integrado en el IDE para reducir el tiempo de resolución de incidencias.

Uso y distribución

- Versión web (a través de plataformas como Hugging Face Chat o la demo de Qwen).
- Extensiones del navegador y plugins para IDEs (VS Code, JetBrains a través de Continue.dev o Tabby).
- Versión escritorio (mediante aplicaciones como LM Studio o Ollama).
- CLI (Llamada directa mediante terminal para integración en pipelines de CI/CD).

Open source

Licencia Apache 2.0, permitiendo total libertad de integración y personalización empresarial.

Integraciones

- Facilidad de integración: Alta (Full code).
- API propia: Compatible con el formato de API de OpenAI, lo que facilita el intercambio de modelos sin cambiar el código de la aplicación.
- Integración nativa con frameworks de orquestación como LangChain y LlamaIndex.
- Compatible con servidores de inferencia como vLLM, TGI y Ollama.

Notas finales

Veredicto técnico

Quiero destacar que Qwen2.5-Coder es, actualmente, la herramienta de código abierto de mayor utilidad para una empresa que busque independencia tecnológica. Vale totalmente la pena el esfuerzo de implementación local debido a su altísima precisión. Como profesional, considero que es el nuevo estándar para el desarrollo asistido por IA fuera del ecosistema cerrado de OpenAI o Anthropic.

información legal, licencias , contratos

- El modelo se distribuye bajo la licencia Apache 2.0. Esto significa que la propiedad intelectual del código generado pertenece al usuario y el modelo puede ser modificado para entrenamiento interno (Fine-tuning) sin restricciones comerciales, siempre que se incluya el aviso de licencia original.

Otros

Es importante mencionar que la versión 32B ha demostrado ser el "sweet spot" (punto óptimo) para la mayoría de las tareas profesionales, ofreciendo capacidades de modelos de 70B+ con un requerimiento de hardware mucho más contenido.

Fuentes consultadas:

- <https://qwenlm.github.io/blog/qwen2.5-coder/>
- <https://github.com/QwenLM/Qwen2.5-Coder>
- <https://huggingface.co/Qwen/Qwen2.5-Coder-32B-Instruct>
- <https://www.linkedin.com/company/qwen-ai/>

## CONSEJOS DE IMPLANTACIÓN

### Aplicación profesional

Según mi experiencia, Qwen2.5-Coder es la herramienta definitiva para empresas de desarrollo de software, consultoras tecnológicas y departamentos de IT que priorizan la soberanía del dato y la eficiencia operativa. El presupuesto necesario es variable: desde 0€ en licencias (Open Source) hasta una inversión inicial en hardware (GPU) de entre 2.000€ y 6.000€ si se busca ejecución local de alto rendimiento. Los puntos clave son su precisión en lenguajes como Python y C++, y su ventana de contexto que permite analizar arquitecturas completas sin desbordar la memoria.

### Madurez digital requerida

- Usuarios: Desarrolladores con nivel senior o medio que puedan validar la calidad del código generado y gestionar entornos de ejecución.
- Empresa: Organizaciones con cultura DevOps integrada y equipos de infraestructura capaces de gestionar contenedores o servidores de inferencia.

### Plan orientativo de implantación

#### Pasos necesarios y estimaciones

- Tiempos de despliegue: De 1 a 3 semanas para una integración completa en el flujo de trabajo corporativo.
- Evaluación inicial (Días 1-3): Inventario de lenguajes utilizados y evaluación de la infraestructura GPU disponible (NVIDIA preferiblemente) o selección de proveedor de API compatible con OpenAI.
- Prueba de concepto (Días 4-7): Instalación de la versión 32B en un entorno controlado mediante Ollama o vLLM para testear la resolución de bugs reales del backlog de la empresa.
- Integración y personalización (Días 8-12): Configuración de extensiones en VS Code (Continue.dev) o JetBrains para todos los desarrolladores y conexión con el servidor central de inferencia.
- Seguimiento y ajuste (Continuo): Monitorización de la tasa de aceptación de sugerencias de código y optimización de prompts específicos para el stack de la empresa.

### Necesidades de formación del equipo

Es imprescindible formar al equipo en ingeniería de prompts aplicada a código, especificando cómo pasar contexto de librerías internas y cómo realizar auditorías de seguridad sobre el código generado por IA.

### Perfiles necesarios

- Perfiles técnicos: Ingeniero de Machine Learning o DevOps para el despliegue y mantenimiento del servidor de inferencia.
- Personal externo: Consultor experto en IA generativa para optimizar los flujos de trabajo iniciales si el equipo interno no tiene experiencia con LLMs locales.

### Retorno de la inversión

- Tiempos: Se estima una reducción del 25-40% en el tiempo de escritura de código repetitivo y tests unitarios tras el primer mes.
- KPIs: Tasa de aceptación de código (%), reducción del tiempo medio de resolución de bugs (MTTR) y volumen de líneas de código documentadas automáticamente.

### Otros

Mi experiencia en implantaciones me lleva a pensar que el verdadero valor de Qwen2.5-Coder no está en escribir código desde cero, sino en su capacidad de razonamiento para refactorizar deuda técnica. Al usarlo tened cuenta de que su versión 32B es sorprendentemente ligera para el rendimiento que ofrece, permitiendo que incluso máquinas con una sola GPU de consumo empresarial (como una RTX 4090 o A6000) soporten a varios desarrolladores simultáneamente si se cuantiza el modelo correctamente. En mi opinión profesional, es la mejor alternativa actual para evitar la dependencia de GitHub Copilot y mantener el código privado dentro de los servidores de la empresa.

## TUTORIAL BÁSICO

### Instalación

Qwen2.5-Coder se distribuye principalmente a través de Hugging Face. Para utilizar el modelo de 32B (el más potente), es fundamental contar con hardware con suficiente VRAM o utilizar técnicas de cuantización.

- **Requisitos de Software:** Es obligatorio usar transformers > 4.37.0. Versiones anteriores lanzarán un KeyError: 'qwen2'.
- **Hardware:** El modelo 32B en precisión completa (BF16) requiere aproximadamente 64GB de VRAM. Para GPUs de consumo (como una RTX 3090/4090 de 24GB), es necesario usar versiones cuantizadas (4-bit o 8-bit) mediante bitsandbytes o formatos GGUF/EXL2.
- **Entorno rápido:** Si buscas desplegar un servidor de inferencia de alto rendimiento, utiliza **vLLM**. Permite optimizar el rendimiento y gestionar múltiples peticiones simultáneas.
- **Checklist:**
  - pip install torch transformers accelerate bitsandbytes
  - Verificar que la versión de Python sea >= 3.9.
  - Asegurar que el entorno detecta correctamente CUDA si usas GPU.

### Uso en el día a día

- **Modelos Instruct vs Base:** Según mi experiencia, debes usar la versión **Instruct** si quieres un chat tipo asistente que explique el código o resuelva dudas. La versión **Base** es exclusivamente para autocompletado de código (FIM - Fill In the Middle), donde el modelo continúa el texto de forma natural.
- **Plantilla de Chat:** Es vital usar apply\_chat\_template del tokenizador. Al usarlo te das cuenta de que el modelo responde mucho mejor si respetas el formato ChatML que viene configurado por defecto.
- **Integración en IDE:** Mi recomendación profesional es usar este modelo a través de extensiones como **Continue.dev** (en VS Code o JetBrains). Puedes conectar Qwen2.5-Coder localmente vía Ollama o LM Studio para tener un copiloto que no envía tus datos a la nube.

### Trucos de experto

- **Fill In the Middle (FIM):** Para tareas de autocompletado en medio de un archivo, usa los tokens especiales <|fim\_prefix|>, <|fim\_suffix|> y <|fim\_middle|>. Esto permite al modelo entender qué código tiene arriba y qué código tiene abajo para rellenar el hueco central con precisión quirúrgica.
- **Contexto Extendido (YaRN):** Aunque el modelo soporta hasta 128K tokens, la configuración estándar suele venir limitada a 32K. En mi opinión, si necesitas procesar repositorios enteros, debes activar YaRN en el config.json ajustando el factor a 4.0 para habilitar los 128K completos.
- **Contexto de Repositorio:** El modelo entiende una estructura especial para múltiples archivos: usa <|repo\_name|> seguido del nombre del proyecto y <|file\_sep|> antes de cada ruta de archivo. Esto mejora drásticamente la capacidad del modelo para entender dependencias entre distintos ficheros de tu proyecto.

### Posibles problemas/incidencias

- **Memoria insuficiente (OOM):** El modelo 32B es pesado. Si experimentas cierres inesperados, usa load\_in\_4bit=True al cargar el modelo con BitsAndBytesConfig. Reduce el consumo de VRAM de ~64GB a unos ~18-20GB.
- **Alucinaciones en librerías raras:** Aunque Qwen2.5-Coder es SOTA (State of the Art), tiende a inventar parámetros en librerías muy recientes o poco documentadas. Siempre verifica el código generado que use frameworks lanzados en los últimos 6 meses.
- **Incompatibilidad con YaRN en vLLM:** vLLM a veces tiene implementaciones estáticas de YaRN que pueden degradar el rendimiento en textos cortos. Mi consejo es activar la escala de cuerda (rope\_scaling) solo cuando realmente vayas a trabajar con archivos masivos.

### Otros

- **Variedad de tamaños:** No te obsesiones con el modelo de 32B si no tienes el hardware. El modelo de **7B** es sorprendentemente capaz para tareas de scripting rápido y corre en casi cualquier GPU moderna con 8GB de VRAM.
- **Matemáticas y Lógica:** A diferencia de otros modelos puramente de código, Qwen2.5-Coder mantiene un rendimiento muy alto en razonamiento lógico y matemático, lo que lo hace ideal para algoritmos complejos.

## PREGUNTAS FRECUENTES

---

### ¿Qué es Qwen2.5-Coder y en qué se diferencia de otros modelos de lenguaje?

Qwen2.5-Coder es una serie de modelos de lenguaje especializados exclusivamente en programación, desarrollados por Alibaba Cloud. A diferencia de los modelos generalistas, sus pesos están optimizados para la generación, corrección y razonamiento lógico de código en más de 92 lenguajes, compitiendo en rendimiento técnico con modelos propietarios de gran escala como GPT-4o.

### ¿Tiene versión gratuita y bajo qué tipo de licencia se distribuye?

Sí, el modelo es de acceso gratuito y se distribuye bajo la licencia Apache 2.0. Esta licencia es una de las más permisivas en el entorno profesional, ya que permite el uso comercial, la modificación del código base y la distribución del software sin costes de regalía, facilitando su integración en productos empresariales.

### ¿Es una tecnología Open Source factible para su descarga en GitHub?

Efectivamente, Qwen2.5-Coder es un proyecto de código abierto. Tanto los pesos del modelo como el código de entrenamiento e inferencia están disponibles en repositorios públicos de GitHub y Hugging Face, lo que permite a los desarrolladores auditar la tecnología, realizar ajustes finos (fine-tuning) y ejecutarla en servidores propios.

### ¿Cómo aborda la privacidad de los datos en entornos corporativos?

Al ser un modelo que permite el despliegue local o en nubes privadas (on-premise), garantiza la soberanía de los datos. Esto significa que el código fuente de una empresa no necesita ser enviado a servidores externos de terceros para ser procesado, mitigando riesgos de filtración de propiedad intelectual y facilitando el cumplimiento de normativas de privacidad.

### ¿Cumple con la normativa española de protección de datos?

El modelo en sí es una herramienta técnica; el cumplimiento del RGPD y la normativa española depende de cómo se implemente. Al permitir el despliegue en infraestructuras locales dentro del territorio nacional o de la Unión Europea, facilita que las organizaciones mantengan el control total sobre los flujos de datos bajo los estándares de seguridad exigidos por la legislación vigente.

### ¿Qué capacidades de procesamiento de archivos extensos posee?

El modelo admite una ventana de contexto de hasta 128K tokens. Para un profesional, esto permite procesar repositorios de código completos o documentación extensa en una sola consulta, permitiendo que la IA comprenda las dependencias entre diferentes archivos y la arquitectura global del proyecto.

### ¿Cuál es el coste asociado a su uso profesional?

El modelo no tiene coste de adquisición de licencia (0€ si se auto-aloja). Los costes asociados son exclusivamente de infraestructura (hardware GPU propio) o de consumo en plataformas de API de terceros, donde las tarifas suelen basarse en el volumen de tokens procesados con precios altamente competitivos frente a modelos cerrados.

### ¿Es una tecnología segura para la auditoría de código?

Sí, gracias a su capacidad de razonamiento lógico, es utilizado profesionalmente para la detección de errores lógicos, vulnerabilidades y reparaciones de bugs (fixing). Al ser ejecutable en entornos aislados, se considera una tecnología segura para auditar sistemas críticos de forma privada.

### ¿Qué requisitos técnicos se necesitan para su implementación local?

Para un despliegue profesional, se requiere un nivel técnico alto en gestión de infraestructura, específicamente conocimientos en Docker, Python y entornos de aceleración por hardware como CUDA u Ollama. La elección del modelo (desde 0.5B hasta 32B) determinará la cantidad de memoria VRAM necesaria en las tarjetas gráficas.

### ¿Se integra con las herramientas de desarrollo habituales (IDE)?

Sí, es compatible con los principales ecosistemas de desarrollo. Puede integrarse de forma nativa o mediante plugins en VS Code y JetBrains utilizando herramientas como Continue.dev o Tabby. Además, su API es compatible con el formato de OpenAI, lo que simplifica la transición desde otros modelos sin reescribir la lógica de integración.

## CONTRATOS Y CONDICIONES

### Opinión inicial

Tras verificar los contratos y condiciones de uso de Qwen2.5-Coder, mi opinión profesional es que nos encontramos ante una herramienta de impacto legal **bajo** si se opta por el despliegue local, pero que requiere una vigilancia de impacto **medio** si se utiliza a través de sus servicios en la nube (DashScope). Al ser un modelo desarrollado por Alibaba Cloud (con sede en China), la soberanía de los datos es el punto crítico para una empresa española. La gran ventaja reside en su licencia Apache 2.0, lo que permite a la empresa ejecutar el modelo en sus propios servidores dentro del Espacio Económico Europeo, eliminando de golpe los riesgos de transferencia internacional de datos y asegurando el cumplimiento estricto del RGPD. Según documentos consultados, el modelo es puramente técnico, lo que reduce riesgos de sesgos sociales, pero aumenta la responsabilidad de la empresa en la supervisión de la seguridad del código generado.

### Principales recomendaciones

- Priorizar el despliegue **On-Premise** (local) o en servidores privados virtuales (VPC) dentro de la UE para garantizar que el código fuente de la empresa nunca salga de su infraestructura.
- Si se utiliza la API de DashScope (Alibaba Cloud), se debe firmar un Acuerdo de Encargado de Tratamiento (DPA) y verificar si existen cláusulas de transferencia internacional a servidores fuera de la UE.
- Implementar una fase de revisión humana obligatoria para todo código generado, ya que la responsabilidad legal por fallos de seguridad o vulnerabilidades en el software resultante recae exclusivamente en la empresa española.
- Desactivar explícitamente cualquier opción de "mejora de producto" o "entrenamiento con datos del usuario" si se utilizan interfaces de terceros o la versión cloud.

### Ley de Inteligencia Artificial (AI Act)

Tras usarlo y analizar su documentación técnica, Qwen2.5-Coder se clasifica generalmente como un modelo de IA de propósito general (GPAI). Al no estar destinado a infraestructuras críticas o sistemas de identificación biométrica de forma nativa, no entra en la categoría de "alto riesgo" por defecto. Sin embargo, la empresa debe cumplir con el deber de transparencia, informando a los empleados o clientes cuando el código haya sido generado artificialmente. Según la normativa, al ser un modelo con licencia abierta, tiene ciertas exenciones en documentación técnica, siempre que no presente riesgos sistémicos para la Unión Europea.

### Privacidad y protección de datos

- **Responsabilidades:** La empresa española actúa como Responsable del Tratamiento. Si el modelo se ejecuta localmente, la empresa tiene el control total. Si se usa la API de Alibaba, Alibaba actúa como Encargado del Tratamiento.
- **Ubicación de los datos:** En el uso local, los datos permanecen en España/UE. En el uso de API oficial, los centros de datos principales de Alibaba Cloud para estos modelos suelen estar en regiones de Asia o EE. UU., lo que complicaría el cumplimiento del RGPD sin salvaguardas adicionales.
- **Transferencia internacional:** El uso de la versión Cloud implica una transferencia fuera del Espacio Económico Europeo. Es necesario verificar la existencia de Cláusulas Contractuales Tipo (SCC).
- **Derechos ARCO:** Al ser un modelo de generación de código, los datos personales no suelen ser el objeto del tratamiento, pero si el código incluye metadatos o comentarios con nombres reales, la empresa debe asegurar la capacidad de rectificación o supresión en sus propios entornos.

### Propiedad intelectual

- **Propiedad de datos:** La licencia Apache 2.0 garantiza que los datos de entrada (prompts) y el código de entrenamiento modificado por la empresa siguen bajo el control y propiedad del usuario.
- **Propiedad del resultado:** Según la legislación española y europea actual, el código generado íntegramente por IA carece de "autoría humana" para ser protegido por derechos de autor, pero la empresa ostenta el derecho de uso comercial y explotación sobre los resultados obtenidos. El usuario es el único responsable de asegurar que el código generado no infrinja patentes preexistentes, aunque el riesgo es bajo en generación de funciones estándar.

### Usos y prohibiciones

- **Usos prohibidos:** No se debe utilizar para generar malware, realizar ataques de inyección de código o automatizar ciberataques, lo cual violaría tanto los términos de servicio de Alibaba como la directiva de ciberseguridad NIS2 en España.
- **Usos admitidos:** Generación, optimización, traducción entre lenguajes de programación y documentación

técnica de software profesional.

#### Seguridad y certificaciones

- **Seguridad:** El modelo permite el aislamiento total. Al probarlo, he verificado que no requiere conexión persistente a internet para funcionar, lo que es ideal para entornos de alta seguridad (air-gapped).
- **Certificaciones:** Alibaba Cloud cuenta con certificaciones ISO 27001 y SOC2, pero estas solo aplican si se usa su infraestructura cloud, no al modelo cuando se descarga y ejecuta por cuenta propia.

#### Otros

Es vital diferenciar entre el modelo (el "software" bajo Apache 2.0) y el servicio de inferencia (la API). Mientras que el modelo es altamente compatible con la legalidad europea por su naturaleza abierta, el servicio de API de Alibaba Cloud requiere un análisis de cumplimiento mucho más riguroso debido a la jurisdicción china de la matriz.

#### Fuentes consultadas:

- [Términos de servicio de Alibaba Cloud DashScope](#)
- [Repositorio oficial Qwen2.5-Coder \(Licencia Apache 2.0\)](#)
- [Documentación técnica de QwenLM](#)
- [Hugging Face Model Card - Qwen2.5-Coder-32B](#)

#### Para más información y herramientas:

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

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

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